

# Bits, Data Types, and Operations



# How do we represent data in a computer?

---

- At the lowest level, a computer is an electronic machine.
  - works by controlling the flow of electrons
  
- Easy to recognize two conditions:
  1. presence of a voltage – we'll call this state “1”
  2. absence of a voltage – we'll call this state “0”
  
- Could base state on *value* of voltage, but control and detection circuits more complex.
  - compare turning on a light switch to measuring or regulating voltage



# Computer is a binary digital system

## Digital system:

- finite number of symbols

## Binary (base two) system:

- has two states: 0 and 1



- Basic unit of information is the *binary digit*, or *bit*.
- Values with more than two states require multiple bits.
  - A collection of **two** bits has **four** possible states:  
**00, 01, 10, 11**
  - A collection of **three** bits has **eight** possible states:  
**000, 001, 010, 011, 100, 101, 110, 111**
  - A collection of  $n$  bits has  $2^n$  possible states.



# What kinds of data do we need to represent?

---

- **Numbers** – signed, unsigned, integers, floating point, complex, rational, irrational, ...
- **Text** – characters, strings, ...
- **Images** – pixels, colors, shapes, ...
- **Sound**
- **Logical** – true, false
- **Instructions**
- ...
- **Data type:**
  - *representation and operations* within the computer
- We'll start with numbers...



# Unsigned Integers

- Non-positional notation

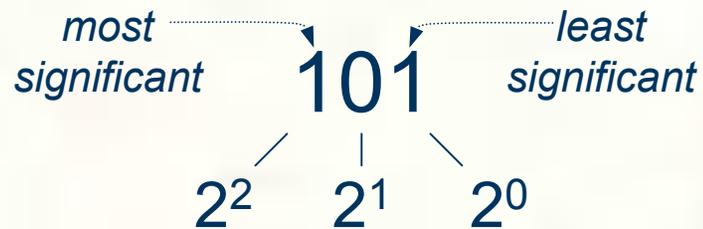
- could represent a number (“5”) with a string of ones (“11111”)
- problems?

- Weighted positional notation

- like decimal numbers: “329”
- “3” is worth 300, because of its position, while “9” is only worth 9



$$3 \times 100 + 2 \times 10 + 9 \times 1 = 329$$



$$1 \times 4 + 0 \times 2 + 1 \times 1 = 5$$



# Unsigned Integers (cont.)

- An  $n$ -bit unsigned integer represents  $2^n$  values:  
from 0 to  $2^n - 1$ .

$2^2$	$2^1$	$2^0$	
0	0	0	0
0	0	1	1
0	1	0	2
0	1	1	3
1	0	0	4
1	0	1	5
1	1	0	6
1	1	1	7



# Unsigned Binary Arithmetic

- Base-2 addition – just like base-10!
  - add from right to left, propagating carry

$$\begin{array}{r} 10010 \\ + \underline{1001} \\ \hline 11011 \end{array}$$
$$\begin{array}{r} \text{carry} \\ \downarrow \\ 10010 \\ + \underline{1011} \\ \hline 11101 \end{array}$$
$$\begin{array}{r} \downarrow \downarrow \downarrow \downarrow \\ 1111 \\ + \underline{1} \\ \hline 10000 \end{array}$$
$$\begin{array}{r} 10111 \\ + \underline{111} \\ \hline \end{array}$$

Subtraction, multiplication, division,...



# Signed Integers

- With  $n$  bits, we have  $2^n$  distinct values.
  - assign about half to positive integers (1 through  $2^{n-1}$ ) and about half to negative ( $-2^{n-1}$  through  $-1$ )
  - that leaves two values: one for 0, and one extra
- Positive integers
  - just like unsigned – zero in *most significant* (MS) bit  
**00101 = 5**
- Negative integers
  - sign-magnitude – set MS bit to show negative, other bits are the same as unsigned  
**10101 = -5**
  - one's complement – flip every bit to represent negative  
**11010 = -5**
  - in either case, MS bit indicates sign: 0=positive, 1=negative



# Two's Complement

- Problems with sign-magnitude and 1's complement
  - two representations of zero (+0 and -0)
  - arithmetic circuits are complex
    - How to add two sign-magnitude numbers?
      - e.g., try  $2 + (-3)$
    - How to add to one's complement numbers?
      - e.g., try  $4 + (-3)$
  - *Two's complement* representation developed to make circuits easy for arithmetic.
    - for each positive number (X), assign value to its negative (-X), such that  $X + (-X) = 0$  with "normal" addition, ignoring carry out

	<b>00101</b>	(5)		<b>01001</b>	(9)
<b>+</b>	<b><u>11011</u></b>	(-5)	<b>+</b>	<b><u>10111</u></b>	(-9)
	<b>00000</b>	(0)		<b>00000</b>	(0)



# Two's Complement Representation

- If number is positive or zero,
  - normal binary representation, zeroes in upper bit(s)
- If number is negative,
  - start with positive number
  - flip every bit (i.e., take the one's complement)
  - then add one

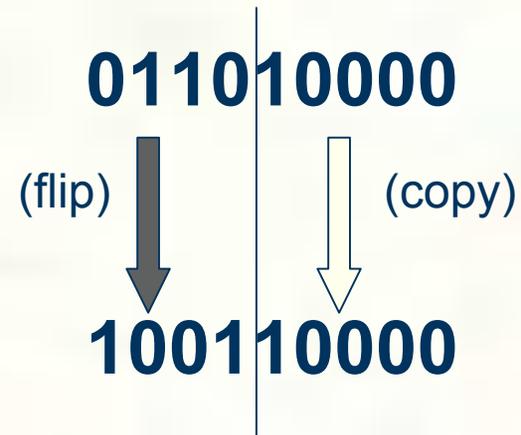
	<b>00101</b>	(5)		<b>01001</b>	(9)
	<b>11010</b>	(1's comp)		<b>10110</b>	(1's comp)
<b>+</b>	<b>1</b>		<b>+</b>	<b>1</b>	
	<b>11011</b>	(-5)		<b>10111</b>	(-9)



# Two's Complement Shortcut

- To take the two's complement of a number:
  - copy bits from right to left until (and including) the first "1"
  - flip remaining bits to the left

 **011010000**  
**100101111** (1's comp)  
**+** **1**  
**100110000**





# Two's Complement Signed Integers

- MS bit is sign bit – it has weight  $-2^{n-1}$ .
- Range of an n-bit number:  $-2^{n-1}$  through  $2^{n-1} - 1$ .
  - The most negative number ( $-2^{n-1}$ ) has no positive counterpart.

-	$2^2$	$2^1$	$2^0$		-	$2^2$	$2^1$	$2^0$	
$2^3$	0	0	0	0	$2^3$	0	0	0	-8
0	0	0	1	1	1	0	0	1	-7
0	0	1	0	2	1	0	1	0	-6
0	0	1	1	3	1	0	1	1	-5
0	1	0	0	4	1	1	0	0	-4
0	1	0	1	5	1	1	0	1	-3
0	1	1	0	6	1	1	1	0	-2
0	1	1	1	7	1	1	1	1	-1



# Converting Binary (2's C) to Decimal

1. If leading bit is one, take two's complement to get a positive number.
2. Add powers of 2 that have "1" in the corresponding bit positions.
3. If original number was negative, add a minus sign.

$$\begin{aligned} X &= 01101000_{\text{two}} \\ &= 2^6 + 2^5 + 2^3 = 64 + 32 + 8 \\ &= 104_{\text{ten}} \end{aligned}$$

$n$	$2^n$
0	1
1	2
2	4
3	8
4	16
5	32
6	64
7	128
8	256
9	512
10	1024

*Assuming 8-bit 2's complement numbers.*



# More Examples

$$\begin{aligned} X &= 00100111_{\text{two}} \\ &= 2^5 + 2^2 + 2^1 + 2^0 = 32 + 4 + 2 + 1 \\ &= 39_{\text{ten}} \end{aligned}$$

$$\begin{aligned} X &= 11100110_{\text{two}} \\ -X &= 00011010 \\ &= 2^4 + 2^3 + 2^1 = 16 + 8 + 2 \\ &= 26_{\text{ten}} \\ X &= -26_{\text{ten}} \end{aligned}$$

$n$	$2^n$
0	1
1	2
2	4
3	8
4	16
5	32
6	64
7	128
8	256
9	512
10	1024

*Assuming 8-bit 2's complement numbers.*



# Converting Decimal to Binary (2's C)

## ■ First Method: *Division*

1. Find magnitude of decimal number. (Always positive.)
2. Divide by two – remainder is least significant bit.
3. Keep dividing by two until answer is zero, writing remainders from right to left.
4. Append a zero as the MS bit; if original number was negative, take two's complement.

$$X = 104_{\text{ten}}$$

$$104/2 = 52 \text{ r}0 \quad \textit{bit 0}$$

$$52/2 = 26 \text{ r}0 \quad \textit{bit 1}$$

$$26/2 = 13 \text{ r}0 \quad \textit{bit 2}$$

$$13/2 = 6 \text{ r}1 \quad \textit{bit 3}$$

$$6/2 = 3 \text{ r}0 \quad \textit{bit 4}$$

$$3/2 = 1 \text{ r}1 \quad \textit{bit 5}$$

$$X = 01101000_{\text{two}}$$

$$1/2 = 0 \text{ r}1 \quad \textit{bit 6}$$



# Converting Decimal to Binary (2's C)

## ■ Second Method: *Subtract Powers of Two*

1. Find magnitude of decimal number.
2. Subtract largest power of two less than or equal to number.
3. Put a one in the corresponding bit position.
4. Keep subtracting until result is zero.
5. Append a zero as MS bit;  
if original was negative, take two's complement.

$n$	$2^n$
0	1
1	2
2	4
3	8
4	16
5	32
6	64
7	128
8	256
9	512
10	1024

$X = 104_{\text{ten}}$	$104 - 64 = 40$	<i>bit 6</i>
	$40 - 32 = 8$	<i>bit 5</i>
	$8 - 8 = 0$	<i>bit 3</i>

$$X = 01101000_{\text{two}}$$



# Operations: Arithmetic and Logical

---

- Recall:
  - a data type includes *representation* and *operations*.
- We now have a good representation for signed integers, so let's look at some arithmetic operations:
  - Addition
  - Subtraction
  - Sign Extension
- We'll also look at overflow conditions for addition.
- Multiplication, division, etc., can be built from these basic operations.
- Logical operations are also useful:
  - AND
  - OR
  - NOT



# Addition

- As we've discussed, 2's comp. addition is just binary addition.
  - assume all integers have the same number of bits
  - ignore carry out
  - for now, assume that sum fits in n-bit 2's comp. representation

$$\begin{array}{r} \mathbf{01101000} \quad (104) \\ + \mathbf{11110000} \quad (-16) \\ \hline \mathbf{01011000} \quad (98) \end{array} \quad \begin{array}{r} \mathbf{11110110} \quad (-10) \\ + \mathbf{11110111} \quad (-9) \\ \hline \mathbf{11101101} \quad (-19) \end{array}$$

*Assuming 8-bit 2's complement numbers.*



# Subtraction

- Negate subtrahend (2nd no.) and add.
  - assume all integers have the same number of bits
  - ignore carry out
  - for now, assume that difference fits in n-bit 2's comp. representation

	<b>01101000</b> (104)		<b>11110110</b> (-10)
-	<u><b>00010000</b></u> (16)	-	<u><b>11110111</b></u> (-9)
	<b>01101000</b> (104)		<b>11110110</b> (-10)
+	<u><b>11110000</b></u> (-16)	+	<u><b>00001001</b></u> (9)
	<b>01011000</b> (88)		<b>11111111</b> (-1)

*Assuming 8-bit 2's complement numbers.*



# Sign Extension

- To add two numbers, we must represent them with the same number of bits.
- If we just pad with zeroes on the left:

**4-bit**

**0100** (4)

**1100** (-4)

**8-bit**

**00000100** (still 4)

**00001100** (12, not -4)

- Instead, replicate the MS bit -- the sign bit:

**4-bit**

**0100** (4)

**1100** (-4)

**8-bit**

**00000100** (still 4)

**11111100** (still -4)



# Overflow

- If operands are too big, then sum cannot be represented as an  $n$ -bit 2's comp number.

$$\begin{array}{r} \mathbf{01000} \quad (8) \\ + \mathbf{01001} \quad (9) \\ \hline \mathbf{10001} \quad (-15) \end{array} \qquad \begin{array}{r} \mathbf{11000} \quad (-8) \\ + \mathbf{10111} \quad (-9) \\ \hline \mathbf{01111} \quad (+15) \end{array}$$

- We have overflow if:
  - signs of both operands are the same, and
  - sign of sum is different.
- Another test -- easy for hardware:
  - carry into MS bit does not equal carry out



# Logical Operations

- Operations on logical TRUE or FALSE

- two states -- takes one bit to represent: TRUE=1, FALSE=0

A	B	A AND B
0	0	0
0	1	0
1	0	0
1	1	1

A	B	A OR B
0	0	0
0	1	1
1	0	1
1	1	1

A	NOT A
0	1
1	0

- View  $n$ -bit number as a collection of  $n$  logical values

- operation applied to each bit independently



# Examples of Logical Operations

## ■ AND

- useful for clearing bits
  - AND with zero = 0
  - AND with one = no change

AND     11000101  
00001111  
00000101

## ■ OR

- useful for setting bits
  - OR with zero = no change
  - OR with one = 1

OR     11000101  
00001111  
11001111

## ■ NOT

- unary operation -- one argument
- flips every bit

NOT     11000101  
00111010



# Hexadecimal Notation

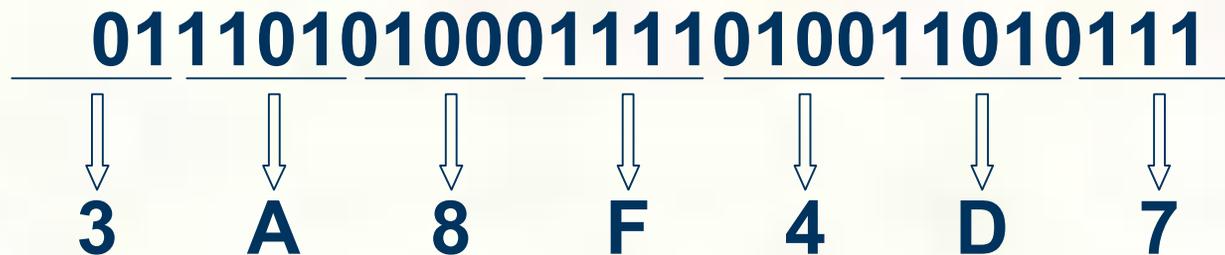
- It is often convenient to write binary (base-2) numbers as hexadecimal (base-16) numbers instead.
  - fewer digits -- four bits per hex digit
  - less error prone -- easy to corrupt long string of 1's and 0's

Binary	Hex	Decima	Binary	Hex	Decima
0000	0	0	1000	8	8
0001	1	1	1001	9	9
0010	2	2	1010	A	10
0011	3	3	1011	B	11
0100	4	4	1100	C	12
0101	5	5	1101	D	13
0110	6	6	1110	E	14
0111	7	7	1111	F	15



# Converting from Binary to Hexadecimal

- Every four bits is a hex digit.
  - start grouping from right-hand side



*This is not a new machine representation,  
just a convenient way to write the number.*



# Fractions: Fixed-Point

- How can we represent fractions?
  - Use a “binary point” to separate positive from negative powers of two -- just like “decimal point.”
  - 2’s comp addition and subtraction still work.

- if binary points are aligned

$$\begin{array}{r}
 \phantom{+} 00101000.101 \quad (40.625) \\
 + \underline{11111110.110} \quad (-1.25) \\
 \hline
 00100111.011 \quad (39.375)
 \end{array}$$

$2^{-1} = 0.5$   
 $2^{-2} = 0.25$   
 $2^{-3} = 0.125$

*No new operations -- same as integer arithmetic.*



# Very Large and Very Small: Floating-Point

- Large values:  $6.023 \times 10^{23}$  -- requires 79 bits
- Small values:  $6.626 \times 10^{-34}$  -- requires >110 bits
- Use equivalent of “scientific notation”:  $F \times 2^E$
- Need to represent  $F$  (*fraction*),  $E$  (*exponent*), and sign.
- IEEE 754 Floating-Point Standard (32-bits):



$$N = (-1)^S \times 1.\text{fraction} \times 2^{\text{exponent}-127}, \quad 1 \leq \text{exponent} \leq 254$$

$$N = (-1)^S \times 0.\text{fraction} \times 2^{-126}, \quad \text{exponent} = 0$$





# Floating-Point Addition

- Will regular 2's complement arithmetic work for Floating Point numbers?
- (*Hint: In decimal, how do we compute  $3.07 \times 10^{12} + 9.11 \times 10^8$ ?*)
  - Step 1: match the exponents. We'll prefer to do this by making the small exponent match the large one since that means shifting the mantissa to the right, and with finite precision representations any lost significant digits will be at the low order end of the number
  - Step 2: Add the mantissas.
  - Step 3: Normalize the result if necessary.

$$\begin{aligned} & 3.75 \times 10^{12} + 9.125 \times 10^8 \\ &= 3.75 \times 10^{12} + .0009125 \times 10^{12} \\ &= 3.7509125 \times 10^{12} \end{aligned}$$



# Floating-Point Addition

- Will regular 2's complement arithmetic work for Floating Point numbers?
- Same algorithm in binary
  - Let's do it in IEEE 754

$$\begin{aligned}
 & 3.75 \times 10^{12} + 9.125 \times 10^8 \\
 &= 3.75 \times 10^{12} + .0009125 \times 10^{12} \\
 &= 3.7509125 \times 10^{12}
 \end{aligned}$$

$$3.75 \times 10^{12} \approx 11.11 \times 2^{40} = 1.111 \times 2^{41}$$

$$= 0 \mid 10101000 \mid 11100000000000000000000000000000 = \text{x54700000}$$

$$9.125 \times 10^8 \approx 1001.001 \times 2^{27} = 1.001001 \times 2^{30}$$

$$= 0 \mid 10000010 \mid 00100100000000000000000000000000 = \text{x41200000}$$

$$1.111 \times 2^{41} + 1.001001 \times 2^{30} = 1.111 \times 2^{41} + 0.00000000001001001 \times 2^{41}$$

$$= 1.11100000001001001 \times 2^{41}$$

$$= 0 \mid 10101000 \mid 111000000010010010000000 = \text{x54701240}$$

(note: binary and decimal not quite equal due to inexact conversion of decimal exponents to binary)



# Floating-Point Multiplication

- In decimal

- Step 1: add exponents
- Step 2: multiply mantissas
- Step 3: normalize result

- Same algorithm in binary

- Let's do it in IEEE 754

$$\begin{aligned} & 3.75 \times 10^{12} \times 9.125 \times 10^8 \\ &= 3.75 \times 9.125 \times 10^{20} \\ &= 34.21875 \times 10^{20} \\ &= 3.421875 \times 10^{21} \end{aligned}$$

$$\begin{aligned} & 1.111 \times 2^{41} \times 1.001001 \times 2^{30} \\ &= 1.111 \times 1.001001 \times 2^{71} \\ &= 10.001000111 \times 2^{71} \\ &= 1.0001000111 \times 2^{72} \\ &= 0 \mid 11000111 \mid 000100011100000000000000 = \text{x6388E000} \end{aligned}$$

(note: binary and decimal not quite equal due to inexact conversion of decimal exponents to binary)



# Text: ASCII Characters

- ASCII: Maps 128 characters to 7-bit code.
  - both printable and non-printable (ESC, DEL, ...) characters

00	nul	10	dle	20	sp	30	0	40	@	50	P	60	`	70	p
01	soh	11	dc1	21	!	31	1	41	A	51	Q	61	a	71	q
02	stx	12	dc2	22	"	32	2	42	B	52	R	62	b	72	r
03	etx	13	dc3	23	#	33	3	43	C	53	S	63	c	73	s
04	eot	14	dc4	24	\$	34	4	44	D	54	T	64	d	74	t
05	enq	15	nak	25	%	35	5	45	E	55	U	65	e	75	u
06	ack	16	syn	26	&	36	6	46	F	56	V	66	f	76	v
07	bel	17	etb	27	'	37	7	47	G	57	W	67	g	77	w
08	bs	18	can	28	(	38	8	48	H	58	X	68	h	78	x
09	ht	19	em	29	)	39	9	49	I	59	Y	69	i	79	y
0a	nl	1a	sub	2a	*	3a	:	4a	J	5a	Z	6a	j	7a	z
0b	vt	1b	esc	2b	+	3b	;	4b	K	5b	[	6b	k	7b	{
0c	np	1c	fs	2c	,	3c	<	4c	L	5c	\	6c	l	7c	
0d	cr	1d	gs	2d	-	3d	=	4d	M	5d	]	6d	m	7d	}
0e	so	1e	rs	2e	.	3e	>	4e	N	5e	^	6e	n	7e	~
0f	si	1f	us	2f	/	3f	?	4f	O	5f	_	6f	o	7f	del



# Interesting Properties of ASCII Code

---

- What is relationship between a decimal digit ('0', '1', ...) and its ASCII code?
- What is the difference between an upper-case letter ('A', 'B', ...) and its lower-case equivalent ('a', 'b', ...)?
- Given two ASCII characters, how do we tell which comes first in alphabetical order?
- Are 128 characters enough?  
(<http://www.unicode.org/>)

*No new operations -- integer arithmetic and logic.*



# Other Data Types

---

## ■ Text strings

- sequence of characters, terminated with NULL (0)
- typically, no hardware support

## ■ Image

- array of pixels
  - monochrome: one bit (1/0 = black/white)
  - color: red, green, blue (RGB) components (e.g., 8 bits each)
  - other properties: transparency
- hardware support:
  - typically none, in general-purpose processors
  - MMX -- multiple 8-bit operations on 32-bit word

## ■ Sound

- sequence of fixed-point numbers



# LC-3 Data Types

---

- Some data types are supported directly by the instruction set architecture.
- For LC-3, there is only one hardware-supported data type:
  - 16-bit 2's complement signed integer
  - Operations: ADD, AND, NOT
- Other data types are supported by interpreting 16-bit values as logical, text, fixed-point, etc., in the software that we write.

# Chapter 2

## Bits, Data Types & Operations

- *Integer Representation*
- *Floating-point Representation*
- *Logic Operations*



# Data types

---

- Our first requirement is to find a way to represent information (data) in a form that is mutually comprehensible by human and machine.
  - Ultimately, we will have to develop schemes for representing all conceivable types of information - language, images, actions, etc.
  - We will start by examining different ways of representing *integers*, and look for a form that suits the computer.
  - Specifically, the devices that make up a computer are switches that can be on or off, i.e. at high or low voltage. Thus they naturally provide us with two symbols to work with: we can call them *on & off*, or (more usefully) *0* and *1*.



# Decimal Numbers

---

- “decimal” means that we have ten digits to use in our representation (the symbols 0 through 9)
- What is 3,546?
  - it is *three* thousands plus *five* hundreds plus *four* tens plus *six* ones.
  - i.e.  $3,546 = 3 \cdot 10^3 + 5 \cdot 10^2 + 4 \cdot 10^1 + 6 \cdot 10^0$
- How about negative numbers?
  - we use two more symbols to distinguish positive and negative:  
+ and -



# Unsigned Binary Integers

$$Y = \text{“abc”} = a.2^2 + b.2^1 + c.2^0$$

(where the digits a, b, c can each take on the values of 0 or 1 only)

**N = number of bits**  
**Range is:**  
 $0 \leq i < 2^N - 1$

	<b>3-bits</b>	<b>5-bits</b>	<b>8-bits</b>
<b>0</b>	<b>000</b>	<b>00000</b>	<b>00000000</b>
<b>1</b>	<b>001</b>	<b>00001</b>	<b>00000001</b>
<b>2</b>	<b>010</b>	<b>00010</b>	<b>00000010</b>
<b>3</b>	<b>011</b>	<b>00011</b>	<b>00000011</b>
<b>4</b>	<b>100</b>	<b>00100</b>	<b>00000100</b>
			<b>0</b>

## Problem:

- How do we represent negative numbers?



# Signed Magnitude

- Leading bit is the sign bit

$$Y = \text{"abc"} = (-1)^a (b \cdot 2^1 + c \cdot 2^0)$$

Range is:  
 $-2^{N-1} + 1 < i < 2^{N-1} - 1$

## Problems:

- How do we do addition/subtraction?
- We have two numbers for zero (+/-0)

-4	10100
-3	10011
-2	10010
-1	10001
-0	10000
+0	00000
+1	00001
+2	00010
+3	00011
+4	00100



# One's Complement

- Invert all bits

If msb (most significant bit) is 1 then the number is negative (same as signed magnitude)

Range is:

$$-2^{N-1} + 1 < i < 2^{N-1} - 1$$

## Problems:

- Same as for signed magnitude

-4	11011
-3	11100
-2	11101
-1	11110
-0	11111
+0	00000
+1	00001
+2	00010
+3	00011
+4	00100



# Two's Complement

## ■ Transformation

- To transform  $a$  into  $-a$ , invert all bits in  $a$  and add 1 to the result

Range is:  
 $-2^{N-1} < i < 2^{N-1} - 1$

-16	10000
...	...
-3	11101
-2	11110
-1	11111
0	00000
+1	00001
+2	00010
+3	00011
...	...
+15	01111



# Manipulating Binary numbers - 1

## ■ Binary to Decimal conversion & vice-versa

- A 4 bit binary number  $A = a_3a_2a_1a_0$  corresponds to:

$$a_3 \cdot 2^3 + a_2 \cdot 2^2 + a_1 \cdot 2^1 + a_0 \cdot 2^0 = a_3 \cdot 8 + a_2 \cdot 4 + a_1 \cdot 2 + a_0 \cdot 1$$

(where  $a_i = 0$  or  $1$  only)

- A decimal number can be broken down by iterative division by 2, assigning bits to the columns that result in an odd number:

$$\text{e.g. } (13)_{10} \Rightarrow (((((13 - 1)/2 - 0)/2 - 1)/2 - 1) = 0 \Rightarrow (01101)_2$$

- In the 2's complement representation, leading zeros do not affect the value of a positive binary number, and leading ones do not affect the value of a negative number. So:

$$01101 = 00001101 = 13 \quad \text{and} \quad 11011 = 11111011 = -5$$



# Manipulating Binary numbers - 2

- Binary addition simply consists of applying, to each column in the sum, the rules:

$$\underline{0 + 0 = 0}$$

$$\underline{1 + 0 = 0 + 1 = 1}$$

$$\underline{1 + 1 = 10}$$

- With 2's complement representation, this works for both positive and negative integers *so long as both numbers being added are represented with the same number of bits.*

e.g. to add the number 13 => 00001101 (8 bits) to -5 => 1011 (4 bits):  
we have to sign-extend (SEXT) the representation of -5 to 8 bits:

00001101

11111011

00001000 => 8 (as expected!)



# Manipulating Binary numbers - 3

## ■ Overflow

- If we add the two (2's complement) 4 bit numbers representing 7 and 5 we get :

$$0111 \quad \Rightarrow +7$$

$$\underline{0101} \quad \Rightarrow \underline{+5}$$

$$1100 \quad \Rightarrow -4 \text{ (in 4 bit 2's comp.)}$$

- We get -4, not +12 as we would expect !!
- We have *overflowed* the range of 4 bit 2's comp. (-8 to +7), so the result is invalid.
- Note that if we add 16 to this result we get back  $16 - 4 = 12$ 
  - this is like “stepping up” to 5 bit 2's complement representation
- In general, if the sum of two positive numbers produces a negative result, or vice versa, an overflow has occurred, and the result is invalid in that representation.



# Limitations of fixed-point

---

- Fixed point numbers are not limited to representing only integers
  - But there are other considerations:
- Range:
  - The magnitude of the numbers we can represent is determined by how many bits we use:
    - e.g. with 32 bits the largest number we can represent is about +/- 2 billion, far too small for many purposes.
- Precision:
  - The exactness with which we can specify a number:
    - e.g. a 32 bit number gives us 31 bits of precision, or roughly 9 figure precision in decimal representation.
- We need another data type!



# Real numbers

- Our decimal system handles non-integer *real* numbers by adding yet another symbol - the decimal point (.) to make a *fixed point* notation:
  - e.g.  $3,456.78 = 3 \cdot 10^3 + 5 \cdot 10^2 + 4 \cdot 10^1 + 6 \cdot 10^0 + 7 \cdot 10^{-1} + 8 \cdot 10^{-2}$
- The *floating point*, or scientific, notation allows us to represent very large and very small numbers (integer or real), with as much or as little precision as needed:
  - Unit of electric charge  $e = 1.602\ 176\ 462 \times 10^{-19}$  Coul.
  - Volume of universe  $= 1 \times 10^{85}$  cm<sup>3</sup>
    - the two components of these numbers are called the mantissa and the exponent



# Floating point numbers in binary

- We mimic the decimal floating point notation to create a “hybrid” binary floating point number:
  - We first use a “binary point” to separate whole numbers from fractional numbers to make a fixed point notation:
    - e.g.  $00011001.110 = 1.2^4 + 1.10^3 + 1.10^1 + 1.2^{-1} + 1.2^{-2} \Rightarrow 25.75$   
( $2^{-1} = 0.5$  and  $2^{-2} = 0.25$ , etc.)
  - We then “float” the binary point:
    - $00011001.110 \Rightarrow 1.1001110 \times 2^4$   
mantissa = 1.1001110, exponent = 4
  - Now we have to express this without the extra symbols ( x, 2, . )
    - by convention, we divide the available bits into three fields:  
sign, mantissa, exponent
  - These are still fixed-precision, only approximate real numbers



# IEEE-754 fp numbers - 1



$$N = (-1)^s \times 1.\mathbf{fraction} \times 2^{(\mathbf{biased\ exp.} - 127)}$$

- Sign: 1 bit

- Mantissa: 23 bits

- We “normalize” the mantissa by dropping the leading 1 and recording only its fractional part (why?)

- Exponent: 8 bits

- In order to handle both +ve and -ve exponents, we add 127 to the actual exponent to create a “biased exponent”:

- $2^{-127} \Rightarrow$  biased exponent = 0000 0000 (= 0)

- $2^0 \Rightarrow$  biased exponent = 0111 1111 (= 127)

- $2^{+127} \Rightarrow$  biased exponent = 1111 1110 (= 254)



# IEEE-754 fp numbers - 2

## ■ Example:

- $25.75 \Rightarrow 00011001.110 \Rightarrow 1.1001110 \times 2^4$
- sign bit = 0 (+ve)
- normalized mantissa (fraction) = 100 1110 0000 0000 0000 0000
- biased exponent =  $4 + 127 = 131 \Rightarrow 1000\ 0011$
- so  $25.75 \Rightarrow 0\ 1000\ 0011\ 100\ 1110\ 0000\ 0000\ 0000\ 0000 \Rightarrow \text{x41CE0000}$

## ■ Values represented by convention:

- Infinity (+ and -): exponent = 255 (1111 1111) and fraction = 0
- NaN (not a number): exponent = 255 and fraction  $\neq 0$
- Zero (0): exponent = 0 and fraction = 0
  - note: exponent = 0  $\Rightarrow$  fraction is *de-normalized*, i.e no hidden 1



# IEEE-754 fp numbers - 3

## ■ Double precision (64 bit) floating point



$$N = (-1)^s \times 1.\text{fraction} \times 2^{(\text{biased exp.} - 1023)}$$

## ◆ Range & Precision

### ◆ 32 bit:

- mantissa of 23 bits + 1 => approx. 7 digits decimal
- $2^{+/-127}$  => approx.  $10^{+/-38}$

### ◆ 64 bit:

- mantissa of 52 bits + 1 => approx. 15 digits decimal
- $2^{+/-1023}$  => approx.  $10^{+/-306}$



# Other Data Types

---

- Other numeric data types
  - e.g. BCD
- Bit vectors & masks
  - sometimes we want to deal with the individual bits themselves
- Text representations
  - ASCII: uses 8 bits to represent main Western alphabetic characters & symbols, plus several “control codes”,
  - Unicode: 16 bit superset of ASCII providing representation of many different alphabets and specialized symbol sets.
  - EBCDIC: IBM’s mainframe representation.



# Hexadecimal Representation

## ■ Base 16 (hexadecimal)

- More a convenience for us humans than a true data type
- 0 to 9 represented as such
- 10, 11, 12, 13, 14, 15 represented by A, B, C, D, E, F
- $16 = 2^4$ : i.e. every hexadecimal digit can be represented by a 4-bit binary (unsigned) and vice-versa.

## ■ Example

$$\begin{aligned}(16AB)_{16} &= x16AB \\ &= 1 \cdot 16^3 + 6 \cdot 16^2 + 10 \cdot 16^1 + 11 \cdot 16^0 \\ &= (5803)_{10} = \#5803 \\ &= b0001\ 0110\ 1010\ 1011\end{aligned}$$



# Another use for bits: Logic

---

## ■ Beyond numbers

- *logical variables* can be *true* or *false*, *on* or *off*, etc., and so are readily represented by the binary system.
- A logical variable  $A$  can take the values *false* = 0 or *true* = 1 only.
- The manipulation of logical variables is known as Boolean Algebra, and has its own set of operations - which are not to be confused with the arithmetical operations of the previous section.
- Some basic operations: NOT, AND, OR, XOR



# Basic Logic Operations

<u>NOT</u>		<u>AND</u>			<u>OR</u>		
<u>A</u>	<u>A'</u>	<u>A</u>	<u>B</u>	<u>A.B</u>	<u>A</u>	<u>B</u>	<u>A+B</u>
0	1	0	0	0	0	0	0
1	0	0	1	0	0	1	1
		1	0	0	1	0	1
		1	1	1	1	1	1

## ■ Equivalent Notations

- not A =  $A' = \overline{A}$
- A and B =  $A.B = A \wedge B = A$  intersection B
- A or B =  $A+B = A \vee B = A$  union B



# More Logic Operations

<u>XOR</u>			<u>XNOR</u>		
<u>A</u>	<u>B</u>	<u><math>A \oplus B</math></u>	<u>A</u>	<u>B</u>	<u><math>(A \oplus B)'</math></u>
0	0	0	0	0	1
0	1	1	0	1	0
1	0	1	1	0	0
1	1	0	1	1	1

- Exclusive OR (XOR): either A or B is 1, not both
- $A \oplus B = A.B' + A'.B$