

# Number Systems, Computer Arithmetic and Character Systems

Joannah Nanjekye

July 09, 2024

# Decimal System

Consider a whole number:  $62_{10}$

▶  $62_{10} = (60 \times 10^1) + (2 \times 10^0)$

▶  $62_{10} = 111110_2$

Consider a fraction number:  $0.456_{10}$

▶  $0.456_{10} = (4 \times 10^{-1}) + (5 \times 10^{-2}) + (6 \times 10^{-3})$

▶  $0.81_{10} = 0.110011_2$

The right hand side is the **Least Significant Bit** while the left hand side is the **Most Significant Bit**

# Binary System

- ▶  $111110_2 = (1 \times 2^5) + (1 \times 2^4) + (1 \times 2^3) + (1 \times 2^2) + (1 \times 2^1) + (0 \times 2^0) = 62_{10}$
- ▶  $1001.101 = 2^3 + 2^0 + 2^{-1} + 2^{-3} = 9.625_{10}$

# Hexadecimal System

- ▶ Binary digits are grouped into sets of four bits, called a nibble
- ▶  $00010001_2 = 11_{16}$
- ▶  $41_{16} = 00100001_2$
- ▶  $41_{16} = (4 \times 16^1) + (1 \times 16^0) = 65_{10}$

0000 = 0	0100 = 4	1000 = 8	1100 = C
0001 = 1	0101 = 5	1001 = 9	1101 = D
0010 = 2	0110 = 6	1010 = A	1110 = E
0011 = 3	0111 = 7	1011 = B	1111 = F

# Hexadecimal System

- ▶ Why the hexadecimal notation is preferred:
  - ▶ It is more compact than binary notation
  - ▶ In most computers, binary data occupy some multiple of 4 bits, and hence some multiple of a single hexadecimal digit
  - ▶ It is extremely easy to convert between binary and hexadecimal notation

# Integer Representation

- ▶ Consider an 8-bit word
- ▶ It can represent positive numbers from 0 to 255
- ▶ And negative numbers from -127 to +128
- ▶ This is because the the Most Significant Bit will represent the sign

# Sign Magnitude Representation

- ▶ In an n-bit word, the right most n - 1 bits hold the magnitude of the integer:
  - ▶ +18 = 00010010
  - ▶ -18 = 10010010
- ▶ Limitations of the sign magnitude representation
  - ▶ Arithmetic has to consider both the sign and magnitude
  - ▶ There are two representations of zero (+0 and -0)

# twos Complement

- ▶ +3 = 00000011
- ▶ +2 = 00000010
- ▶ +1 = 00000001
- ▶ +0 = 00000000
- ▶ -1 = 11111111
- ▶ -2 = 11111110
- ▶ -3 = 11111101



# Integer Arithmetic: Negation

- ▶ In two's complement negation is achieved in two steps:
  - ▶ compute the Boolean complement of each bit of the integer
  - ▶ Treat the result as an unsigned binary integer, add 1
- ▶  $+18 = 00010010$
- ▶ bit wise complement =  $11101101$
- ▶  $11101101 + 1 = 11101110 = -18$
- ▶ Negation is referred to as the **twos complement operation**
- ▶ There are two special cases 0 and 128 for an 8-bit representation:
  - ▶  $0 = 00000000$  twos complement + 1 =  $[1]00000000 = 0$   
(discard carried 1)
  - ▶  $+128 = 10000000$  twos complement + 1 =  $10000000 = -128$

# Integer Arithmetic: Addition

- ▶ Change negative numbers to twos complement
- ▶ Treat positive numbers as unsigned
- ▶ The result is in twos complement form
- ▶ Whether it is positive or negative

$\begin{array}{r} 1001 = -7 \\ +\underline{0101} = 5 \\ 1110 = -2 \\ (a) (-7) + (+5) \end{array}$	$\begin{array}{r} 1100 = -4 \\ +\underline{0100} = 4 \\ \underline{1}0000 = 0 \\ (b) (-4) + (+4) \end{array}$
$\begin{array}{r} 0011 = 3 \\ +\underline{0100} = 4 \\ 0111 = 7 \\ (c) (+3) + (+4) \end{array}$	$\begin{array}{r} 1100 = -4 \\ +\underline{1111} = -1 \\ \underline{1}1011 = -5 \\ (d) (-4) + (-1) \end{array}$
$\begin{array}{r} 0101 = 5 \\ +\underline{0100} = 4 \\ 1001 = \text{Overflow} \\ (e) (+5) + (+4) \end{array}$	$\begin{array}{r} 1001 = -7 \\ +\underline{1010} = -6 \\ \underline{1}0011 = \text{Overflow} \\ (f) (-7) + (-6) \end{array}$

# Integer Arithmetic: Addition and Overflow

- ▶ The result may be larger than the word size
- ▶ This is called an **overflow**
- ▶ When adding two numbers, if they are both positive or both negative, an overflow occurs if and only if the result has the opposite sign

$\begin{array}{r} 1001 = -7 \\ +\underline{0101} = 5 \\ 1110 = -2 \\ (a) (-7) + (+5) \end{array}$	$\begin{array}{r} 1100 = -4 \\ +\underline{0100} = 4 \\ \underline{1}0000 = 0 \\ (b) (-4) + (+4) \end{array}$
$\begin{array}{r} 0011 = 3 \\ +\underline{0100} = 4 \\ 0111 = 7 \\ (c) (+3) + (+4) \end{array}$	$\begin{array}{r} 1100 = -4 \\ +\underline{1111} = -1 \\ \underline{1}1011 = -5 \\ (d) (-4) + (-1) \end{array}$
$\begin{array}{r} 0101 = 5 \\ +\underline{0100} = 4 \\ 1001 = \text{Overflow} \\ (e) (+5) + (+4) \end{array}$	$\begin{array}{r} 1001 = -7 \\ +\underline{1010} = -6 \\ \underline{1}0011 = \text{Overflow} \\ (f) (-7) + (-6) \end{array}$

# Integer Arithmetic: Subtraction

Subtraction is achieved by performing addition. The number being subtracted (subtrahend) should be converted to twos complement before being added to the minuend

$\begin{array}{r} 0010 = 2 \\ +1001 = -7 \\ \hline 1011 = -5 \end{array}$ <p>(a) M = 2 = 0010 S = 7 = 0111 -S = 1001</p>	$\begin{array}{r} 0101 = 5 \\ +1110 = -2 \\ \hline 0011 = 3 \end{array}$ <p>(b) M = 5 = 0101 S = 2 = 0010 -S = 1110</p>
$\begin{array}{r} 1011 = -5 \\ +1110 = -2 \\ \hline 1001 = -7 \end{array}$ <p>(c) M = -5 = 1011 S = 2 = 0010 -S = 1110</p>	$\begin{array}{r} 0101 = 5 \\ +0010 = 2 \\ \hline 0111 = 7 \end{array}$ <p>(d) M = 5 = 0101 S = -2 = 1110 -S = 0010</p>
$\begin{array}{r} 0111 = 7 \\ +0111 = 7 \\ \hline 1110 = \text{Overflow} \end{array}$ <p>(e) M = 7 = 0111 S = -7 = 1001 -S = 0111</p>	$\begin{array}{r} 1010 = -6 \\ +1100 = -4 \\ \hline 0110 = \text{Overflow} \end{array}$ <p>(f) M = -6 = 1010 S = 4 = 0100 -S = 1100</p>

## Unsigned: Multiplication

- ▶ Generate partial products which are summed to get the final product
- ▶ Each successive partial product is shifted one position to the left relative to the preceding partial product
- ▶ The multiplication of two n-bit binary integers results in a product of up to 2n bits in length (e.g.,  $11 * 11 = 1001$ )

1011		<b>Multiplicand (11)</b>
×1101		<b>Multiplier (13)</b>
<hr/>		
1011	}	<b>Partial products</b>
0000		
1011		
1011		
<hr/>		
10001111		<b>Product (143)</b>

# Unsigned: Multiplication

There are two approaches:

- ▶ perform a running addition on the partial products (saves storage)
- ▶ For each 1 on the multiplier, an add and a shift operation are required; but for each 0, only a shift is required

C	A	Q	M		
0	0000	1101	1011	<b>Initial values</b>	
0	1011	1101	1011	<b>Add</b>	} <b>First cycle</b>
0	0101	1110	1011	<b>Shift</b>	
0	0010	1111	1011	<b>Shift</b>	} <b>Second cycle</b>
0	1101	1111	1011	<b>Add</b>	} <b>Third cycle</b>
0	0110	1111	1011	<b>Shift</b>	
1	0001	1111	1011	<b>Add</b>	} <b>Fourth cycle</b>
0	1000	1111	1011	<b>Shift</b>	

# Signed: Multiplication

One alternative is to:

- ▶ convert both multiplier and multiplicand to positive numbers
- ▶ perform the multiplication
- ▶ compute the two's complement of the result

We can use Booth's algorithm instead.

$\begin{array}{r} 0111 \\ \times 0011 \\ \hline 11111001 \\ 00000000 \\ 000111 \\ \hline 00010101 \end{array}$	$\begin{array}{r} (0) \\ 1-0 \\ 1-1 \\ 0-1 \\ \hline (21) \end{array}$
$\begin{array}{r} 0111 \\ \times 1101 \\ \hline 11111001 \\ 0000111 \\ 111001 \\ \hline 11101011 \end{array}$	$\begin{array}{r} (0) \\ 1-0 \\ 0-1 \\ 1-0 \\ \hline (-21) \end{array}$

(a)  $(7) \times (3) = (21)$

(b)  $(7) \times (-3) = (-21)$

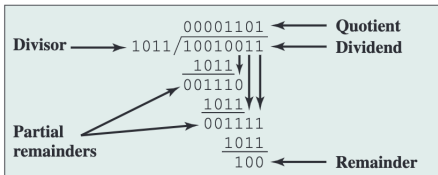
$\begin{array}{r} 1001 \\ \times 0011 \\ \hline 00000111 \\ 00000000 \\ 111001 \\ \hline 11101011 \end{array}$	$\begin{array}{r} (0) \\ 1-0 \\ 1-1 \\ 0-1 \\ \hline (-21) \end{array}$
$\begin{array}{r} 1001 \\ \times 1101 \\ \hline 00000111 \\ 1111001 \\ 000111 \\ \hline 00010101 \end{array}$	$\begin{array}{r} (0) \\ 1-0 \\ 0-1 \\ 1-0 \\ \hline (21) \end{array}$

(c)  $(-7) \times (3) = (-21)$

(d)  $(-7) \times (-3) = (21)$

# Unsigned: Division

- ▶ Examine the bits of the dividend from left to right
- ▶ Stop when the bits examined correspond to number greater than or equal to the divisor
- ▶ Before the condition, place 0s in the quotient from left to right
- ▶ After the condition, place a 1 in the quotient and the divisor is subtracted from the partial dividend





# Signed: Division

1. Compute the two's complement of the divisor
2. Express the dividend as a  $2n$ -bit positive number, e.g. 4-bit 0111 becomes 00000111
3. Shift the dividend left 1 bit position
4. Subtract the divisor from the dividend
5. If the result is positive, MSB = 0, then dividend is 1
6. If the result is negative, MSB = 1, then dividend is 0
7. Repeat steps 2 through 4 as many times as there are bit positions in the dividend

# Signed: Division

<b>A</b> 0000	<b>Q</b> 0111	<b>Initial value</b>
0000 <u>1101</u> 1101 0000	1110  1110	<b>Shift</b> <b>Use twos complement of 0011 for subtraction</b> <b>Subtract</b> <b>Restore, set <math>Q_0 = 0</math></b>
0001 <u>1101</u> 1110 0001	1100  1100	<b>Shift</b>  <b>Subtract</b> <b>Restore, set <math>Q_0 = 0</math></b>
0011 <u>1101</u> 0000	1000  1001	<b>Shift</b>  <b>Subtract, set <math>Q_0 = 1</math></b>
0001 <u>1101</u> 1110 0001	0010  0010	<b>Shift</b>  <b>Subtract</b> <b>Restore, set <math>Q_0 = 0</math></b>

# Limitations of the twos Complement

- ▶ Very large and very small numbers can not be represented
- ▶ The fractional part of the quotient in a division of two large numbers could be lost

# The Scientific Notation

- ▶ Move the decimal point to a convenient location
- ▶ Use the exponent of 10 to keep track of that decimal point
- ▶ This approach can be used on binary numbers as follows:

$$\pm S \times B^{\pm E}$$

Three fields are used to store the binary word of a number:

- ▶ Sign: plus or minus
- ▶ Significand S
- ▶ Exponent E
- ▶ Base B

# Floating-point Representation

- ▶ The leftmost bit stores the sign of the number
- ▶ The exponent value is stored in the next 8 bits
- ▶ The representation used is known as a biased representation
- ▶ Bias =  $(2^{k-1} - 1) = 2^7 - 1 = 127$ , k is the number of bits
- ▶ The final portion of the word (23 bits in this case) is the significand



$$\begin{aligned} 1.1010001 \times 2^{10100} &= 0 \ 10010011 \ 101000100000000000000000 = 1.6328125 \times 2^{20} \\ -1.1010001 \times 2^{10100} &= 1 \ 10010011 \ 101000100000000000000000 = -1.6328125 \times 2^{20} \\ 1.1010001 \times 2^{-10100} &= 0 \ 01101011 \ 101000100000000000000000 = 1.6328125 \times 2^{-20} \\ -1.1010001 \times 2^{-10100} &= 1 \ 01101011 \ 101000100000000000000000 = -1.6328125 \times 2^{-20} \end{aligned}$$

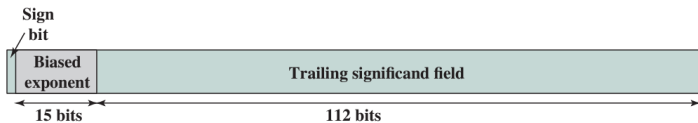
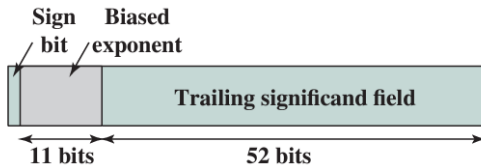
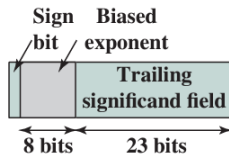
# IEEE Binary Floating-point Representation

Defines the following format:

- ▶ **Arithmetic format:** all the mandatory operations defined by the standard are supported by the format
- ▶ **Basic format:** covers five floating-point representations, three binary and two decimal, whose encodings are specified by the standard, and which can be used for arithmetic
- ▶ **Interchange format:** fixed-length binary encoding that allows data interchange between different platforms and that can be used for storage

Each of the formats have bit lengths of 32, 64, and 128 bits and exponents of 8, 11, and 15 bits

# IEEE Binary Floating-point Representation



# Floating-point Arithmetic

Floating-Point Numbers	Arithmetic Operations
$X = X_S \times B^{X_E}$ $Y = Y_S \times B^{Y_E}$	$\left. \begin{aligned} X + Y &= (X_S \times B^{X_E - Y_E} + Y_S) \times B^{Y_E} \\ X - Y &= (X_S \times B^{X_E - Y_E} - Y_S) \times B^{Y_E} \end{aligned} \right\} X_E \leq Y_E$ $X \times Y = (X_S \times Y_S) \times B^{X_E + Y_E}$ $\frac{X}{Y} = \left( \frac{X_S}{Y_S} \right) \times B^{X_E - Y_E}$

Examples:

$$X = 0.3 \times 10^2 = 30$$

$$Y = 0.2 \times 10^3 = 200$$

$$X + Y = (0.3 \times 10^{2-3} + 0.2) \times 10^3 = 0.23 \times 10^3 = 230$$

$$X - Y = (0.3 \times 10^{2-3} - 0.2) \times 10^3 = (-0.17) \times 10^3 = -170$$

$$X \times Y = (0.3 \times 0.2) \times 10^{2+3} = 0.06 \times 10^5 = 6000$$

$$X \div Y = (0.3 \div 0.2) \times 10^{2-3} = 1.5 \times 10^{-1} = 0.15$$



# Addition and Subtraction

- ▶ Check for zeros
- ▶ Align the significands
- ▶ Add or subtract the significands
- ▶ Normalize the result

# Guard Bits

- ▶ The ALU loads the exponent and significand before a floating point operation
- ▶ The length of the register is almost always greater than the length of the significand plus an implied bit
- ▶ The register contains additional bits, called guard bits
- ▶ To pad out the right end of the significand with 0s

$$\begin{aligned}x &= 1.000\dots00 \times 2^1 \\ \underline{-y} &= \underline{0.111\dots11} \times 2^1 \\ z &= 0.000\dots01 \times 2^1 \\ &= 1.000\dots00 \times 2^{-22}\end{aligned}$$

(a) Binary example, without guard bits

$$\begin{aligned}x &= .100000 \times 16^1 \\ \underline{-y} &= \underline{.0FFFFFF} \times 16^1 \\ z &= .000001 \times 16^1 \\ &= .100000 \times 16^{-4}\end{aligned}$$

(c) Hexadecimal example, without guard bits

$$\begin{aligned}x &= 1.000\dots00\ 0000 \times 2^1 \\ \underline{-y} &= \underline{0.111\dots11\ 1000} \times 2^1 \\ z &= 0.000\dots00\ 1000 \times 2^1 \\ &= 1.000\dots00\ 0000 \times 2^{-23}\end{aligned}$$

(b) Binary example, with guard bits

$$\begin{aligned}x &= .100000\ 00 \times 16^1 \\ \underline{-y} &= \underline{.0FFFFFF\ F0} \times 16^1 \\ z &= .000000\ 10 \times 16^1 \\ &= .100000\ 00 \times 16^{-5}\end{aligned}$$

(d) Hexadecimal example, with guard bits

# Rounding

- ▶ This is where extra bits in a floating-point format number are removed
- ▶ To generate a number that is close to the original number
- ▶ Alternative approaches to rounding:
  - ▶ **Round to nearest:** rounded to the nearest representable number
  - ▶ **Round toward  $+\infty$ :** rounded up toward plus infinity
  - ▶ **Round toward  $-\infty$ :** rounded up toward minus infinity
  - ▶ **Round toward 0:** rounded toward zero

# Character Systems

- ▶ We deal with both symbolic alphabetic and numeric data
- ▶ Systems to encode characters as binary numbers are required for a computer
- ▶ Common systems include:
  - ▶ ASCII (American Standard Code for Information Interchange)
  - ▶ EBCDIC (Extended Binary Coded Decimal Interchange Code)
  - ▶ UNICODE (A unique number is provided for each character)

# ASCII

"One" = O = 0x4F, n = 0x6E, e = 0x65

Dec	Hx	Oct	Char	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr
0	0	000	<b>NUL</b> (null)	32	20	040	&#32;	Space	64	40	100	&#64;	@	96	60	140	&#96;	`
1	1	001	<b>SOH</b> (start of heading)	33	21	041	&#33;	!	65	41	101	&#65;	A	97	61	141	&#97;	a
2	2	002	<b>STX</b> (start of text)	34	22	042	&#34;	"	66	42	102	&#66;	B	98	62	142	&#98;	b
3	3	003	<b>ETX</b> (end of text)	35	23	043	&#35;	#	67	43	103	&#67;	C	99	63	143	&#99;	c
4	4	004	<b>EOT</b> (end of transmission)	36	24	044	&#36;	&	68	44	104	&#68;	D	100	64	144	&#100;	d
5	5	005	<b>ENQ</b> (enquiry)	37	25	045	&#37;	%	69	45	105	&#69;	E	101	65	145	&#101;	e
6	6	006	<b>ACK</b> (acknowledge)	38	26	046	&#38;	&	70	46	106	&#70;	F	102	66	146	&#102;	f
7	7	007	<b>BEL</b> (bell)	39	27	047	&#39;	'	71	47	107	&#71;	G	103	67	147	&#103;	g
8	8	010	<b>BS</b> (backspace)	40	28	050	&#40;	(	72	48	110	&#72;	H	104	68	150	&#104;	h
9	9	011	<b>TAB</b> (horizontal tab)	41	29	051	&#41;	)	73	49	111	&#73;	I	105	69	151	&#105;	i
10	A	012	<b>LF</b> (NL line feed, new line)	42	2A	052	&#42;	*	74	4A	112	&#74;	J	106	6A	152	&#106;	j
11	B	013	<b>VT</b> (vertical tab)	43	2B	053	&#43;	+	75	4B	113	&#75;	K	107	6B	153	&#107;	k
12	C	014	<b>FF</b> (NP form feed, new page)	44	2C	054	&#44;	,	76	4C	114	&#76;	L	108	6C	154	&#108;	l
13	D	015	<b>CR</b> (carriage return)	45	2D	055	&#45;	-	77	4D	115	&#77;	M	109	6D	155	&#109;	m
14	E	016	<b>SO</b> (shift out)	46	2E	056	&#46;	.	78	4E	116	&#78;	N	110	6E	156	&#110;	n
15	F	017	<b>SI</b> (shift in)	47	2F	057	&#47;	/	79	4F	117	&#79;	O	111	6F	157	&#111;	o
16	10	020	<b>DLE</b> (data link escape)	48	30	060	&#48;	0	80	50	120	&#80;	P	112	70	160	&#112;	p
17	11	021	<b>DC1</b> (device control 1)	49	31	061	&#49;	1	81	51	121	&#81;	Q	113	71	161	&#113;	q
18	12	022	<b>DC2</b> (device control 2)	50	32	062	&#50;	2	82	52	122	&#82;	R	114	72	162	&#114;	r
19	13	023	<b>DC3</b> (device control 3)	51	33	063	&#51;	3	83	53	123	&#83;	S	115	73	163	&#115;	s
20	14	024	<b>DC4</b> (device control 4)	52	34	064	&#52;	4	84	54	124	&#84;	T	116	74	164	&#116;	t
21	15	025	<b>NAK</b> (negative acknowledge)	53	35	065	&#53;	5	85	55	125	&#85;	U	117	75	165	&#117;	u
22	16	026	<b>SYN</b> (synchronous idle)	54	36	066	&#54;	6	86	56	126	&#86;	V	118	76	166	&#118;	v
23	17	027	<b>ETB</b> (end of trans. block)	55	37	067	&#55;	7	87	57	127	&#87;	W	119	77	167	&#119;	w
24	18	030	<b>CAN</b> (cancel)	56	38	070	&#56;	8	88	58	130	&#88;	X	120	78	170	&#120;	x
25	19	031	<b>EM</b> (end of medium)	57	39	071	&#57;	9	89	59	131	&#89;	Y	121	79	171	&#121;	y
26	1A	032	<b>SUB</b> (substitute)	58	3A	072	&#58;	:	90	5A	132	&#90;	Z	122	7A	172	&#122;	z
27	1B	033	<b>ESC</b> (escape)	59	3B	073	&#59;	:	91	5B	133	&#91;	[	123	7B	173	&#123;	{
28	1C	034	<b>FS</b> (file separator)	60	3C	074	&#60;	<	92	5C	134	&#92;	\	124	7C	174	&#124;	
29	1D	035	<b>GS</b> (group separator)	61	3D	075	&#61;	=	93	5D	135	&#93;	]	125	7D	175	&#125;	}
30	1E	036	<b>RS</b> (record separator)	62	3E	076	&#62;	>	94	5E	136	&#94;	^	126	7E	176	&#126;	~
31	1F	037	<b>US</b> (unit separator)	63	3F	077	&#63;	?	95	5F	137	&#95;	_	127	7F	177	&#127;	DEL