

Make a reasonable effort to show your work. Clearly indicate your answers.

Problem 1 (12 points)

Express each of the following signed decimal numbers using **8-bit signed magnitude representation** and **8-bit 2's complement representation**.

	Signed-magnitude	2's complement
a. +41	00101001	00101001
b. +84	01010100	01010100
c. +103	01100111	01100111
d. -37	10100101	11011011
e. -88	11011000	10101000
f. -106	11101010	10010110

a)  $\begin{array}{r} 41 \\ -32 \\ \hline 9 \\ -8 \\ \hline 1 \\ -1 \\ \hline 0 \end{array}$

$00101001$

b)  $\begin{array}{r} 84 \\ -64 \\ \hline 20 \\ -16 \\ \hline 4 \\ -4 \\ \hline 0 \end{array}$

$01010100$

c)  $\begin{array}{r} 103 \\ -64 \\ \hline 39 \\ -32 \\ \hline 7 \\ -4 \\ \hline 3 \\ -2 \\ \hline 1 \\ -1 \\ \hline 0 \end{array}$

$01100111$

d) -37

$+37 = 00100101$   
 $\text{Comp} = 11011010$   
 $\begin{array}{r} 1 \\ 1 \\ 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ 0 \end{array}$

$2's \text{ Comp} = 11011011 = -37$   
 $S.M. = 10100101 = -37$

$\begin{array}{r} 37 \\ -32 \\ \hline 5 \\ -4 \\ \hline 1 \\ -1 \\ \hline 0 \end{array}$

e) -88

$$+88 = 01011000$$

$$\begin{array}{r} -64 \\ \hline 24 \\ -16 \\ \hline 8 \\ -8 \\ \hline 0 \end{array}$$

$$S_0 \quad S_M =$$

$$Z_5 \text{ comp} =$$

$$11011000$$

$$01011000$$

$$10100111$$

$$\begin{array}{r} + \\ \hline 10101000 \end{array}$$

f) -106

$$\begin{array}{r} +106 \\ -64 \\ \hline 42 \\ -32 \\ \hline 10 \\ -8 \\ \hline 2 \end{array}$$

$$= 01101010$$

$$S_0 \quad S_M =$$

$$Z_5 \text{ comp} =$$

$$11101010$$

$$01101010$$

$$10010101$$

$$\begin{array}{r} + \\ \hline 10010110 \end{array}$$

## Problem 2 (18 points)

Perform the following arithmetic operations in **binary** using **8-bit 2's complement representation**.

- Perform subtraction by finding the 2's complement of the subtrahend and then performing addition.
- Give an appropriate indicator of overflow, if it occurs. *Make sure that you identify the specific sign in your arithmetic that indicates overflow.*
- I did not choose the numbers in this problem by accident; where appropriate, you need not repeat your work to represent the values you use in solving this problem.

- $(+41) + (+84)$
- $(-106) + (+41)$
- $(-37) + (+83)$
- $(+103) - (+41)$
- $(+41) + (+103)$
- $(-88) - (+84)$

$$\begin{array}{r} a) 00101001 = 41 \\ + 01010100 = 84 \\ \hline 01111101 \end{array}$$

$$\begin{array}{r} b) -106 = 10010110 \\ 41 = +00101001 \\ \hline 10111111 \end{array}$$

2's comp  $\rightarrow$

$$\begin{array}{r} c) -37 + 83 \\ +37 \\ = 00100101 \\ -37 = \overset{\textcircled{1}}{1}1011011 \\ + 01010011 \\ \hline \textcircled{1}00101110 \end{array}$$

over  
f/w

$$\begin{array}{r} d) 41 = 00101001 \\ \text{comp} = 11010110 \\ + \\ \hline 11010111 = -41 \\ \textcircled{1}01100111 = 103 \\ + 11010111 = -41 \\ \hline \textcircled{1}00111110 \end{array}$$

over f/w

e)  $41 + 103$

$$\begin{array}{r}
 00101001 = 41 \\
 + 01100111 = 103 \\
 \hline
 \end{array}$$

10010000

f)  $+88$

$$\begin{array}{r}
 01011000 \\
 \text{comp} = 10100111 \\
 + \\
 \hline
 10101000
 \end{array}$$

$+84$

$$\begin{array}{r}
 01010100 \\
 \text{comp} = 10101011 \\
 + \\
 \hline
 10101100
 \end{array}$$

$$\begin{array}{r}
 ① 10101000 \\
 + 10101100 \\
 \hline
 \end{array}$$

① 101010100

Overflow

Problem 3 (8 points)

Represent the number  $(+26.6875)_{10}$  as a 16-bit floating point number. This particular floating point format, uses signed-magnitude to store numbers in the mantissa (a normalized fraction) and the exponent (an integer). The mantissa has 10 bits (including the mantissa's sign bit) and the exponent has 6 bits (including the exponent's sign bit).

$$26 = 11010.1011$$

$$= 26.6875$$

$$\begin{array}{r} 26 \\ -16 \\ \hline 10 \\ -8 \\ \hline 2 \\ -0 \\ \hline 0 \end{array}$$

$$.6875 \times 2 = 1$$

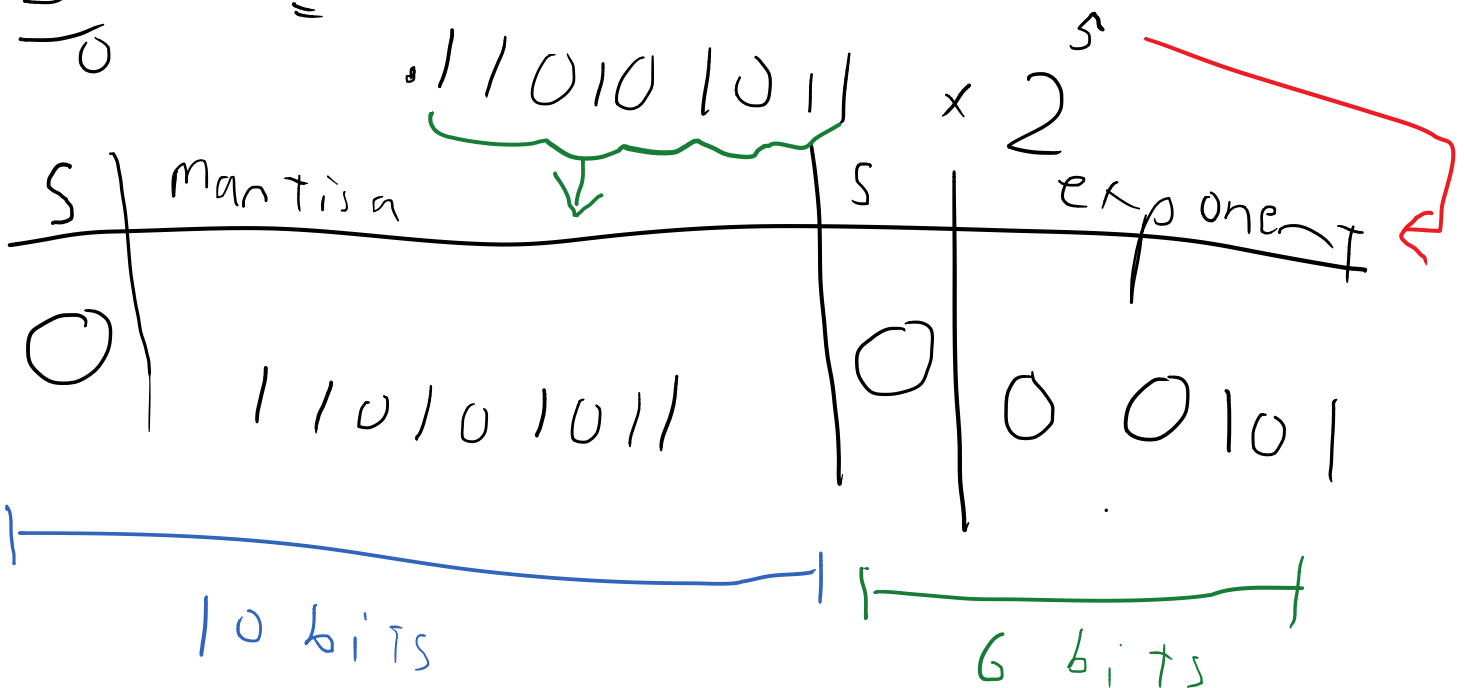
$$.375 \times 2 = 0$$

$$.75 \times 2 = 1$$

$$.5 \times 2 = 1$$

$$11010.1011$$

$$110101011$$



$$= 0110101011000101 \text{ FP}$$

#### Problem 4 (2 points)

One day, I overheard Michael Mistake-Maker expressing wondrous amazement about binary floating point numbers:

*"You really can get something for nothing! For a particular number of bits, you can define a floating-point representation that gives a greater range and a greater precision than the same number of bits give in a fixed-point representation. Why would you ever use a fixed-point representation?"*

Based on your own knowledge or upon research, gently deflate Michael's enthusiasm. I am not looking for essay-length answers; a few sentences will suffice. Answer the question in your own words, and briefly cite any sources that you use.

---

Because as floating point utilizes a set number of allocated bits for the mantissa and then exponent. What you get actually relates to rounding (Advantages). For numbers like say the 26.6875, it is fine, but when you have very large numbers and are utilizing floating point. You are pretty much using a form of rounding where you are cutting off some of the bits in the mantissa – or if it is big enough – the exponent. This is because you are using only a set number of bits for both mantissa and exponent (Advantages).

A fixed point representation, although inefficient for some smaller cases. Will provide exact certain accuracy for representing extremely large numbers exactly as they are without leaving it to rounding error (Benefits)

In the days of the past, when space exploration defined a nation, extremely large and accurate numbers were required. Fixed point is also easier to compute (Benefits). And as such in the days of space exploration, may be preferred in some cases to provide an 'absolute' value that was also faster to calculate at the cost of say memory – which is often cheaper to have more of back then - as opposed to a floating point one (Benefits).

#### Works Cited

"Advantages and Disadvantages of Floating Point and Fixed Point Representations." *Stack Overflow*. Stack

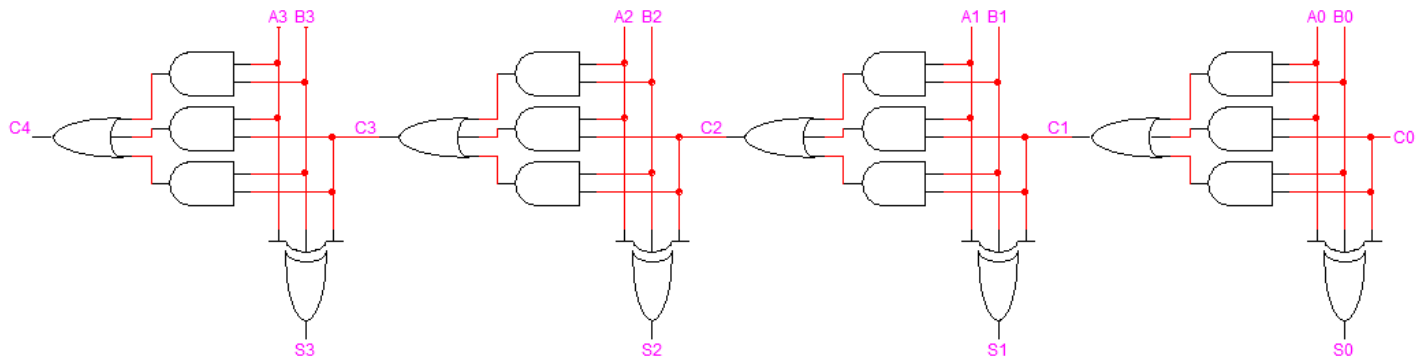
Exchange Inc., 10 Mar. 2012. Web. 22 Mar. 2015.

"Benefits of Using Fixed-Point Hardware." *MathWorks*. The MathWorks Inc. Web. 22 Mar. 2015.

**Problem 5 (20 points)**

Use the “design by contraction” strategy to derive the equations that implement a 4-bit increment-by-3 circuit from the equations that implement a standard adder.

If it helps to visualize the circuit, imagine that you are starting with a 4-bit adder:



Your strategy should consist of choosing specific values for B (B3 B2 B1 B0) and C0 that will cause the output S (S3 S2 S1 S0) to always equal A (A3 A2 A1 A0) plus 3, and then using Boolean algebra to reduce the equations that contain constants.

$$S_i = A_i \oplus B_i \oplus C_i$$

$$C_{i+1} = A_i B_i + A_i C_i + B_i C_i$$

$$S = A + B + C$$

where  $B = 3$   
 $B = 0011$

where  $C = 0$

where  $B(1) = 1$

$$S_1 = A_1 \oplus B_1 \oplus C_1$$

$$= A_1 \oplus 1 \oplus A_0$$

$$S_1 = \bar{A}_1 \oplus A_0$$

$$C_2 = A_1 B_1 + A_1 C_1 + B_1 C_1$$

$$= A_1(1) + A_1 A_0 + (1) A_0$$

$$= A_1 + A_1 A_0 + A_0$$

$$= A_1(1 + A_0) + A_0$$

$$C_2 = A_1 + A_0$$

$$S_0 = A_0 \oplus B_0 \oplus C_0 = A_0 \oplus 1 \oplus 0$$

$$= A_0 \oplus 1 = \bar{A}_0$$

so  $S_0 = \bar{A}_0$

$$C_1 = A_0 B_0 + A_0 C_0 + B_0 C_0$$

$$= A_0(1) + A_0(0) + (1)(0)$$

$$= A_0 + 0 + 0 = A_0$$

so  $C_1 = A_0$

$$S_2 = A_2 \oplus B_2 \oplus C_2$$

where  $B(2) = 0$

$$= A_2 \oplus 0 \oplus (A_1 + A_0) = A_2 \oplus (A_1 + A_0) = S_2$$

$$C_3 = A_2 B_2 + A_2 C_2 + B_2 C_2$$

$$= A_2(0) + A_2(A_1 + A_0) + (0)C_2$$

$$= A_2(A_1 + A_0) = C_3$$

$$S_3 = A_3 \oplus B_3 \oplus C_3$$

$$= A_3 \oplus 0 \oplus A_2(A_1 + A_0)$$

$$S_3 = A_3 \oplus A_2(A_1 + A_0)$$

where  $B(3) = 0$

$$C_4 = \text{don't need}$$

