

Homework 2

(Due date: February 16th @ 11:59 pm)

Presentation and clarity are very important! Show your procedure!

PROBLEM 1 (15 PTS)

- Multiply the following signed fixed-point numbers (6 pts):

1000.001 × 01.100101	01.10001 × 10.1101	01.101 × 1.001011
-------------------------	-----------------------	----------------------

- Get the division result (with $x = 4$ fractional bits) for the following signed fixed-point numbers:

1.010110 ÷ 010.1011	101.1001 ÷ 1.0101	11.011 ÷ 1.10111
------------------------	----------------------	---------------------

PROBLEM 2 (11 PTS)

- We want to represent numbers between -251.5 and 256.7 . What is the fixed-point format that requires the fewest number of bits for a resolution better or equal than 0.0015 ? (4 pts).
- We want to represent numbers between -127.69 and 120.69 . What is the fixed-point format that requires the fewest number of bits for a resolution better or equal than 0.0025 ? (4 pts).
- Represent these numbers in Fixed Point Arithmetic (signed numbers). Select the minimum number of bits in each case.

-128.6875	-147.3125	79.125
-----------	-----------	--------

PROBLEM 3 (10 PTS)

- Complete the table for the following fixed point formats (signed numbers): (4 pts)

Fractional bits	Integer Bits	FX Format	Range	Dynamic Range (dB)	Resolution
7	5				
12	4				
17	7				

- Complete the table for these floating point formats (which resemble the IEEE-754 standard). Only consider ordinary numbers.

Exponent bits (E)	Significant bits (p)	Min	Max	Range of e	Range of significand
7	8				
8	15				
11	36				

PROBLEM 4 (20 PTS)

- Calculate the decimal values of the following floating point numbers represented as hexadecimals. Show your procedure.

Single (32 bits)		Double (64 bits)	
✓ 7FCE4710	✓ 803ACBAC	✓ FEAAFC0FEE000000	✓ 000ABBAF25C00000
✓ BDE32856	✓ 7BEAD360	✓ 7A09D3784D039800	✓ FFFECE4710ABCDEF

PROBLEM 5 (44 PTS)

- Perform the following 32-bit floating point operations. For fixed-point division, use 8 fractional bits. Truncate the result when required. Show your work: how you got the significand and the biased exponent bits of the result. Provide the 32-bit result.

✓ 801A8000 + 33CEC000	✓ 40D90000 - 42EAC000	✓ 0E2CE000 × 8B092000	✓ C9744000 ÷ 40C90000
✓ ECE4710A + FF800000	✓ CF4A8000 - 30A90000	✓ AD0BEBED × 7F800000	✓ 000C0000 ÷ C94A0000