

CSci370 Computer Architecture: Homework 3

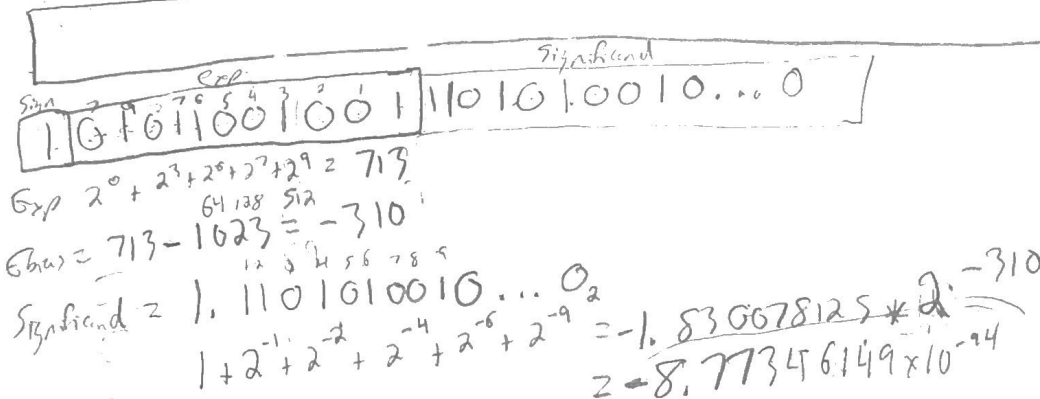
Due date: On or before Monday, April 06, 2020 in class

Absolutely no copying others' works

Name: Dech Trom

- The purpose of homeworks is for students to practice for the exams without others' help, so the penalty of mistakes will be minor.
- Without practicing for the exams properly, students would not be able to do well on the exams.

1. (Floating-point representation: 20%) What decimal number does the hexadecimal number (AC9D 4800 0000 0000)₁₆ represent if it is a floating point number? Use the IEEE 754 standard.



Ans = -3.619201739

A = 1.601110610 x 2^6

B = 1.6100001100 x 2^1

C = -1.1616111110 x 2^2

Convert A, B, C to binary keep 10 bits of number

2. (Floating-point arithmetic: 80%) IEEE 754-2008 contains a half precision that is only 16 bits wide. The leftmost bit is still the sign bit, the exponent is 5 bits wide and has a bias of 15, and the mantissa is 10 bits long. A hidden 1 is assumed. Assuming the three numbers, A=1.237219, B=2.524723, and C=6.742837, are stored in the 16-bit IEEE 754-2008 format, calculate A×B-C by hand. Assume 1 guard bit, 1 round bit, and 1 sticky bit, and round to the nearest even (only on the Step 4. Rounding the Significand). Show all the steps in normalized floating point numbers with 11 bits of precision (or 10 bits of fraction); i.e., $(-1)^S \times (1.F)_2 \times 2^E$, where S is the sign bit, |F|=10, and E is without including the bias. Note that

- have to show all five steps step-by-step,
- the calculation should be similar to those in [Slide 9.11](#) and [Slide 9.16](#); i.e., the numbers should use a base-2 scientific notation and include 3 extra bits each (e.g., IEEE format and bias need NOT be used),
- to answer this question, simulate the hardware by including 3 extra bits for each operand and result and using the 3 extra bits as possible as you could,
- if the number has to be shifted to the right on the Step 1. Making Exponents Equal, round (just using the regular rounding method) the number to 10 bits of fraction and 3 extra bits before calculation,
- rounding to the nearest even will be applied only once after finishing all calculation because the 3 extra bits need to be kept for further calculation, and
- 2's complement representation has to be used if needed.

Work attached

-1.1100111110×2^1

$= -3.62109375_{10}$

Converting to Floating Binary

22

A 1.237219

11

$$0.237219 \times 2 = 0.474438 = 0$$

$$0.474438 \times 2 = 0.948876 = 0.$$

$$0.948876 \times 2 = 1.897752 = 1.$$

$$0.897752 \times 2 = 1.795504 = 1$$

$$0.795504 \times 2 = 1.591008 = 1$$

$$0.991008 \times 2 = 1.982016 = 1$$

$$0.182016 \times 2 = 0.364032 = 0.$$

$$0.364032 \times 2 = 0.728064 = 0.$$

$$2 \quad 1,456128 = 1 \dots$$

$$= 0.912256 = 0$$

$$A \approx 1.001110010 \cdot 2^{10}$$

h2 2.524723

$$10.10000110010 = 1.0100001100 \times 2^1$$

$$0.524723 \times 2 = 1.049446 = 1$$

$$0,49446 \times 2 = 0,98892 = 0$$

$$.098882 \times 2 = 0.197764 \approx 0$$

$$.197784 \times 2 = 0.395568 = 0$$

$$.395568 \times 2 = 0.791136 = 0$$

$$= 1.582272 \approx 1$$

$$= 1.164544 \approx 1$$

$$0.329088 = 0$$

$$0.658176 \approx 0$$

1.31635221

G.632764E 0

C2 6.74 28 37

110. 1011110001

$$1.10101110001 \times 2^2$$

GR 521

Step 1 calc exponent

$$0 + 1 = 1$$

Step 2 sign and multiplication

$$\begin{array}{r} 1.0011110010 \\ 1.0100001101 \\ \hline \end{array}$$

$$\begin{array}{r} 100001000000000 \\ 110011110010000 \\ 110011110010000 \\ 100000000000000 \\ 100000000000000 \\ 100000000000000 \\ 100000000000000 \\ 110011110010000 \\ 000000000000000 \\ 100111100100000000 \\ \hline 1.10001111001001010 \end{array}$$

Step 3 normalize

$$1.10001111001001010 \times 2^1$$

Step 4 Round sign and

$$\begin{array}{r} 1.10001111001001010 \\ \hline \end{array}$$

Step 5 Check overflow

No overflow

$$\text{Subtract } C = -1.1010111110$$

Step 1 Exponents equal

$$\begin{array}{r} 0.1100011110100 \times 2^2 \\ 1.1010111100000 \times 2^2 \\ \hline \end{array}$$

Step 2 Perform subtraction

$$\begin{array}{r} 0.1100011110100 \\ 0.0101000000000 \\ \hline 1.0001110000100 \\ 0.1110011110010 \text{ (complement)} \\ \hline 0.1110011110010 \end{array}$$

Step 3 normalize

$$-1.1100011110010 \times 2^1$$

Step 4 round the sign and

$$-1.1100011110010 \times 2^1$$

Step 5 Check overflow / underflow

= None

Answer

$$-1.1100011110010 \times 2^1 = -3.62809375_{10}$$