**Name: KEY**                                                    **Id#**

# COE 301/ICS 233, Term 151

# Computer Architecture & Assembly Language

## Quiz# 4

Date: Tuesday, Nov. 3, 2015

**Q1.**

**(i)** What is the decimal value of the following single-precision floating-point number?

**0100 0011 0110 1001 1000 0100 0000 0000**.

$= + (1.1101001100001000...0)_2 * 2^{(134-127)} = + (1.1101001100001000…0)_2 * 2^7$
$= + (11101001.100001000….0)_2$
$= + 233.515625$

**(ii)** Show the single-precision floating-point binary representation for: **555.9375**.

$555.9375 = (1000101011.1111)_2 = (1.0001010111111)_2 * 2^9$
Exp. $= 9 + 127 = 136$
<u>Single precision binary representation:</u>

**0100 0100 0000 1010 1111 1100 0000 0000**

**(iii)** Perform the following floating-point operation rounding the result to the **<u>nearest even</u>**. Perform the operation using **guard**, **round** and **sticky** bits.

We add three bits for each operand representing G, R, S bits as follows.

```
     1.000 0000 1000 0000 0000 0000 000 x 2³⁷
 –   1.000 0000 0000 0000 0100 0000 000 x 2²⁹
─────────────────────────────────────────────
     1.000 0000 1000 0000 0000 0000 000 x 2³⁷
 –   0.000 0000 1000 0000 0000 0000 010 x 2³⁷  (align)
─────────────────────────────────────────────
    01.000 0000 1000 0000 0000 0000 000 x 2³⁷
 + 11.111 1111 0111 1111 1111 1111 110 x 2³⁷  (2's complement)
─────────────────────────────────────────────
    00.111 1111 1111 1111 1111 1111 110 x 2³⁷
= +  0.111 1111 1111 1111 1111 1111 110 x 2³⁷
= +  1.111 1111 1111 1111 1111 1111 100 x 2³⁶  (normalize)
= + 10.000 0000 0000 0000 0000 0000     x 2³⁶  (round)
= +  1.000 0000 0000 0000 0000 0000     x 2³⁷  (renormalize)
```

**Q2.** Consider a simplified 8-bit floating point representation following the general guidelines of the IEEE format in representing normalized, denormalized, Nan, infinity and 0. Suppose that the number of bits used for the exponent is 3 and for the fraction is 4 bits.

      **(i)**      Determine the smallest and largest positive values of normalized numbers.

      **(ii)**     Determine the smallest and largest positive values of denormalized numbers.

      **(iii)**   Determine the representation used for +0 and +∞.

      **(iv)**   What is the largest and smallest error in this representation?

(i)

| Number | S | Exp | Fraction | E | Value |
|---|---|---|---|---|---|
| Smallest normalize number | 0 | 001 | 0000 | 1-3=-2 | 1*1/4=1/4 |
| Largest normalize number | 0 | 110 | 1111 | 6-3=3 | 31/2=15.5 |

(ii)

| Number | S | Exp | Fraction | E | Value |
|---|---|---|---|---|---|
| Smallest denormalize number | 0 | 000 | 0001 | -2 | 1/16*1/4=1/64 |
| Largest denormalize number | 0 | 000 | 1111 | -2 | 15/64≈1/4 |

(iii)

| Number | S | Exp | Fraction |
|---|---|---|---|
| +0 | 0 | 000 | 0000 |
| +∞ | 0 | 111 | 0000 |

(iii)

If we consider the largest values with E=110=3, we can see that these values range from 8 to 15.5. The values in this range are: 8, 8.5, 9, 9.5 …., 15.5. Thus, the largest error is 0.5/2=0.25.

If we consider the smallest normalized values with E=001=-2, we can see that these values range from 1/4 to 31/16. The values in this range are: 1/4=16/64, 17/64, 18/64, 19/64 …., 31/16. Thus, the smallest error is 1/64*2=1/128.

Note that the same magnitude of error occurs in the representation of denormalized numbers as they are in the range of 1/64, 2/64,..,15/64.