Yousef Jarrar

CSE 401 – Dr. Gomez

Homework 3


## 3.8

*Assume 185 and 122 are signed 8-bit decimal integers stored in*

*sign-magnitude format. Calculate 185 - 122. Is there overflow, underflow, or*

*neither?*


The sign-magnitude representation for these two given 8-bit decimal integers are:

$(185)_{10}=(1011\ 1001)_2$

= -57

Similarly,

$(122)_{10} = (0111\ 1010)_2$

= 122

In sign-magnitude form, the result of: $- A – (+B) = - (A + B)$

Therefore,

185-122 = -57 – (122)

= -57 + 122

The calculation for 185 – 122 is:


0111001 (57)

+

1111010 (122)

10110011


Neglecting the 8<sup>th</sup> bit, the result is: $(011011)_2 = (51)_{10}$

The sign of the operation will be negative.

**<u>Thus, the calculation 185 – 122 = (-57) – (122) results in -51</u>**.

Since, there was a carry generated during the addition, hence, an **overflow** occurs in the operation.

## 3.24

*Write down the binary representation of the decimal number 63.25 assuming the IEEE 754 double precision format.*

Double precision uses two 32-bit words for representing a floating-point value. The numbers are 53 bit long in double precision (1 + 52).

Following steps must be taken for converting 63.25 from base 10 to IEEE 754 double precision:

• Convert 63 to base 2 which is $(111111)_2$

• Convert (.25) to base 2 which is $(.01)_2$

• Add both:

$$(63) + (.25) = (111111) + (0.01)$$

$$= (1111111.01)_2$$

Writing it to binary scientific notation:

$$63.25 \times 10^0 = (111111.01) * (2^0)$$

Normalize and move the binary point 5 times to the left

$$(1.1111101) * (2^5)$$

This number is written in IEEE 754 Double Precision:

$$(-1)^s * (1+\text{Fraction}) * 2^{(\text{exponent} - 1023)}$$

$$= (-1)^0 * ( 1 + (.1111\ 1010\ 0000) ) * 2^{(1028 - 1023)}$$

5 is converted to the correct bias; Since the bias is 1023, we add 5 to it. Which gives us 1028. Binary form of 1028 is: $10000000100_2$

The Binary Representation assuming IEEE 754 double precision is: 0 1000000100 1111 1010 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000