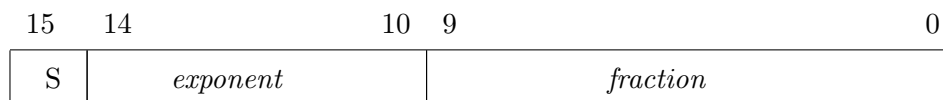


**Question 3 (25 points):**

NVIDIA has defined a “half-precision” floating pointing format for use in its Graphics Processing Units (GPUs). A floating-point number is represented in this format in 16 bits as follows: the most significant bit is the sign bit, next there are 5 bits used for the exponent, and 10 bits for the fraction. This format is illustrated below:



The exponent is expressed in excess-16 format (also known as a bias representation). Given the binary representation above, the decimal value of the number represented can be computed by the following expression:

$$N = \begin{cases} (-1)^S \times 0.0 & \text{if } exponent = 0 \text{ and } fraction = 0 \\ (-1)^S \times 0.fraction \times 2^{-14} & \text{if } exponent = 0 \text{ and } fraction \neq 0 \\ (-1)^S \times 1.fraction \times 2^{exponent-15} & \text{if } 0 < exponent < 31 \\ (-1)^S \times \infty & \text{if } exponent = 31 \text{ and } fraction = 0 \\ NaN & \text{if } exponent = 31 \text{ and } fraction \neq 0 \end{cases}$$

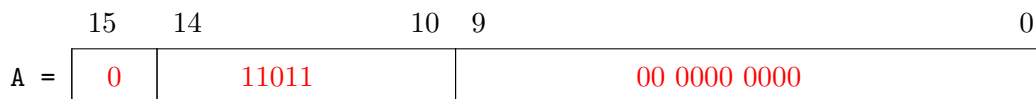
- a. (5 points) Give the bit pattern for the representation of the number  $-4.125_{10}$  in this notation.



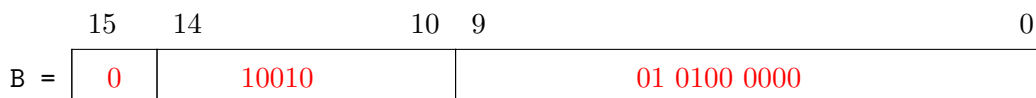
$$4.125_{10} = 4.0 + 0.125 = 4 + \frac{1}{8} = (-1)^0 \times 100.001 = 1.00001 \times 2^2 \\ \Rightarrow exponent - 15 = 2 \Rightarrow exponent = 17$$

Let A = 0x6C00 and B = 0x4940 be two floating pointing numbers in this format.

- b. (5 points) What is the value of A and what is the value of B? Express each of these values both in normalized base-two notation and in decimal notation.



$$A = (-1)^0 \times 1.0 \times 2^{27-15} = 1.0 \times 2^{12} \\ A = 1\ 0000\ 0000\ 0000_2 = 2^{10} \times 2^2 = 1024 \times 4 = 4096_{10}$$



$$B = (-1)^0 \times 1.0101 \times 2^{18-15} = 1.0101 \times 2^3$$

$$B = 1010.1_2 = 8 + 2 + 0.5 = 10.5_{10}$$

- c. **5 points**) What is the true value of  $A + B$  expressed in decimal notation? In other words, what is the value of  $A + B$  if an infinite precision could be used to compute the addition and to store the result?

$$A + B = 4096 + 10.5 = 4106.5_{10}$$

- d. **(5 points)** Assume a floating-point unit that uses the NVIDIA format presented above. This unit has no guard, no round, and no sticky bits. What is the value of  $A + B$ , expressed both in normalized base-two notation and in decimal notation, computed by this machine?

To align A with B, we need to move the binary point of B nine positions to the left. Therefore:

$$B = 0.0000\ 0000\ 1010\ 001 \times 2^{12}$$

mantissa
A = + 1.0000 0000 00
B = + 0.0000 0000 10
-----
A+B = 1.0000 0000 10

Therefore  $A + B = 1.0000\ 0000\ 10 \times 2^{12} = 1\ 0000\ 0000\ 1000 = 2^{12} + 2^3 = 1024 \times 4 + 8 = 4096 + 8 = 4104_{10}$

- e. **(5 points)** Now assume a floating-point unit that uses the NVIDIA format presented above. This unit has one guard, one round, and one sticky bit. What is the value of  $A + B$ , expressed in normalized base-two notation, computed by this machine?

mantissa	Guard	Round	Sticky
A = + 1.0000 0000 00	0	0	0
B = + 0.0000 0000 10	1	0	1
-----			
A+B = 1.0000 0000 10	1	0	1

Now we have to round up because of the sticky bit. Therefore the result is:

$$A + B = 1.0000\ 0000\ 11 \times 2^{12} = 1\ 0000\ 0000\ 1100_2 = 4096 + 8 + 4 = 4108_{10}$$