

# Floating Point Addition / Subtraction

# A and B both are floating point number.

$\Rightarrow A + B$  (Make sure the number is in binary)

i) Normalize both A and B.

ii) Align the bin point so that the lower exponent match with the higher exponent.

iii) Now add / sub accordingly.

iv) Normalize the result.

iv) Round if necessary.

Ex:  $0.999 \times 10^1 + 1.610 \times 10^{-1}$  ; size of exponent field is 3 bits

$$= 99.99 + 0.1610$$

$$= 1100011.111110101 + 0.0010100100$$

$$= 1.10001111110101 \times 2^6 + 1.0100100 \times 2^{-3}$$

$$= 1.10001111110101 \times 2^6 + 0.0000000010100100 \times 2^6$$

$$= 1.10010 \times 2^6 \quad (\text{Ans})$$

$$\text{Bias} = 2^{3-1} - 1 = 3$$

$$\text{Biased Exp.} = 3 + 6 = 9$$

$$\text{Range} = 0 \text{ to } 2^3 - 1$$

$$= 0 \text{ to } 7$$

$$= 1 \text{ to } 6 \text{ [reserved} \\ \nearrow \text{0 and 7]} \\ \text{upper} \\ \text{range}$$

$$\# 110100.111011 \times 2^8 + 10110.11111 \times 2^7$$

$$= 1.10100111011 \times 2^{13} + 1.011011111 \times 2^{11}$$

$$= 1.10100111011 \times 2^{13} + 0.0101101111 \times 2^{13}$$

$$= 10.00000011010 \times 2^{13}$$

$$= 1.000000011010 \times 2^{14} \quad (\text{Ans})$$

$$9 > 6 \Rightarrow \text{So,}$$

overflow

Given number is too small to represent using the mentioned system.

Overflow / Underflow detection:

Step 1: Find the biased exponent of the answer.

Step 2: " " range of the biased exponent of the given system. (1 to upper Range)

Step 3: Detection:

if (Biased exponent  $< 1$ ):

underflow

else if (Biased exponent  $>$  upper Range)

overflow

else :  $[1 \leq \text{Biased Exp} \leq \text{upper Range}]$

No over/under flow

## Floating Point Multiplication

# A and B both are floating point numbers.

$\Rightarrow A \times B$  (Make sure the number is in binary)

i) Normalize both A and B. Ex:  $1.110 \times 2^5 \times 1.11 \times 2^{-5}$

ii) Add the exponents.  $= 1.110 \times 1.11 \times 2^{5+(-5)}$

iii) Now multiply accordingly.  $= 11.0001 \times 2^0$

iv) Normalize the result.  $= 1.10001 \times 2^1$  (Ans)

iv) Round if necessary.

v) Determine the sign from the operation.

## F.P instructions in RiscV

# Suppose, two single prec. floating point numbers A, B are stored in memory. The memory locations are directly stored in registers  $X_{10}, X_{11}$ .

Write necessary code to store the result of  $A+B$  in the memory address that is stored in  $X_{13}$ .

Sol<sup>n</sup>:

fload  $f_1, 0(X_{10}) ; f_1 = A$

fload  $f_2, 0(X_{11}) ; f_2 = B$

fadd.s  $f_3, f_1, f_2 ; f_3 = A+B$

fsw  $f_3, 0(X_{13})$