

## Team- AA

### Team Members-

Ansh Chaudhary (22095133)

Aditya C Singh (22095143)

In this problem statement, our task was to design and optimize a Processing Element (PE) that will serve as a fundamental building block for the AI accelerator.

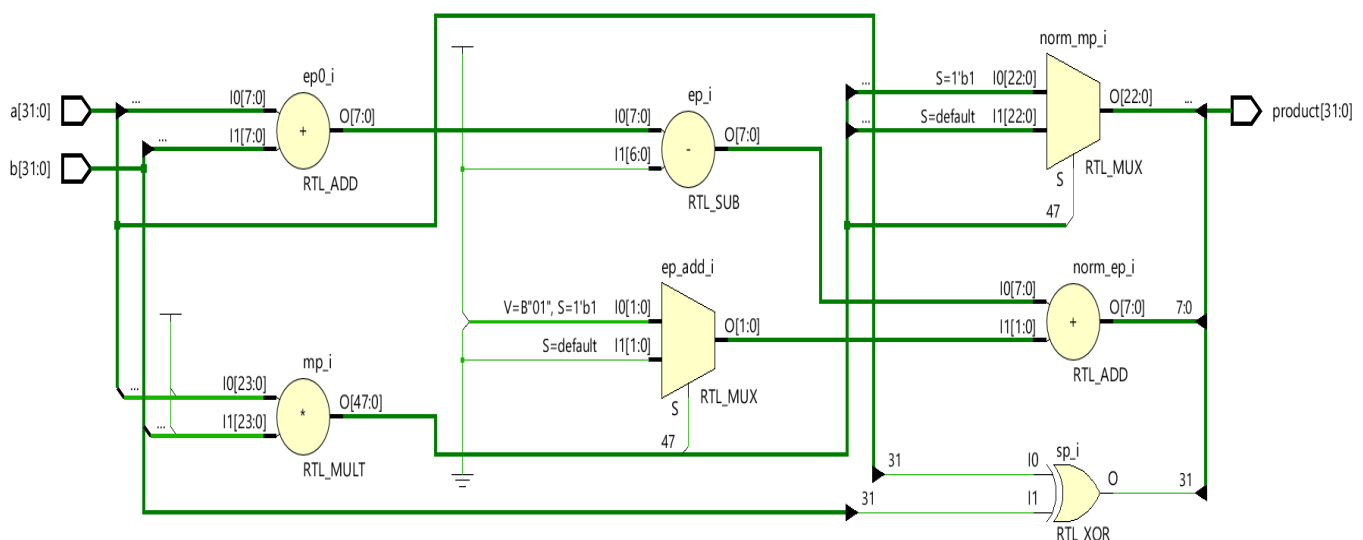
The PE consists of the following components:

- 1) Floating Point Multiplier: This component is dedicated to performing multiplication operations on 32-bit floating-point numbers.
- 2) Floating Point Adder: Responsible for executing addition and subtraction operations on 32-bit floating-points numbers.
- 3) Register: The PE includes a register for storing intermediate results and operands during arithmetic operations.

### Floating Point Multiplier:

Let us consider we have two 32-bit floating-point numbers  $a$  and  $b$ . Let us define a module `fp_mult`, consisting of input as 32-bit  $a$  and 32-bit  $b$ , while output is 32 bit variable product. Let  $s_a$  and  $s_b$  represent the sign bit of  $a$  and  $b$ . Let  $e_a$  and  $e_b$  represent the 8-bit exponent of  $a$  and  $b$ . Let  $m_a$ ,  $m_b$  represent the mantissas of  $a$  and  $b$  concatenated with 1. Let  $s_p$  represent the sign bit of final product which is obtained by having XOR operation between  $s_a$  and  $s_b$ . The 8-bit exponents  $e_a$  and  $e_b$  are added and a bias value of 127 in 8-bit binary is removed from their sum. Also, the multiplication operation is performed between  $m_a$  and  $m_b$  and assigned to a variable  $m_p$ ,

which is 48 bit long. Another 8-bit variable `ep_add` is created which is assigned with 8-bit binary representation of decimal 1 or 0, if the MSB of `mp` is 1 or 0 respectively. Another 23-bit variable `norm_mp` is created, which is assigned with 46<sup>th</sup> to 24<sup>th</sup> bit of `mp`, if MSB of `mp` is 1, else it is assigned with 45<sup>th</sup> to 23<sup>rd</sup> bit of `mp`. The final product is represented using a 32-bit variable `product`, which is formed by concatenation of `sp`, `norm_ep`, `norm_mp`.



### Floating Point Adder:

A module `FloatingAddition` is created with inputs 32-bit variable `A` and 32-bit variable `B` and output variable `result` which is also 32-bit long. 24-bit variables `A_Mantissa`, `B_Mantissa`, `Temp_Mantissa` are created. A 22-bit variable `Mantissa` is created. 7-bit variables `A_Exponent`, `B_Exponent`, `Temp_Exponent`, `diff_Exponent`, `Exponent`, `exp_adjust` are created. Single bit variables `Sign`, `carry`, `comp`, `A_sign`, `B_sign` and `Temp_sign` are created. Next, we compare the exponents of `A` and `B` and assign single bit 1 to `comp` if exponent of `A` is greater than or equal to the exponent of `B`, or else we assign single bit 0 to it.

Basically, the variables A\_Mantissa, A\_Exponent, A\_sign represent the mantissa, exponent and the sign bit of the greater number out of A and B, while the variables B\_Mantissa, B\_Exponent, B\_sign represent the mantissa, exponent and the sign bit of the smaller number out of A and B. The comparison between the two is decided by the variable comp. The variable diff\_Exponent is assigned with the difference between the exponents A\_Exponent and B\_Exponent. The B\_Mantissa is then right shifted by diff\_Exponent bits. If the A\_sign and B\_sign are same, then addition is performed or else subtraction between them is performed and assigned to carry, Temp\_Mantissa concatenated. The variable exp\_adjust is assigned with A\_Exponent. If after addition operation, carry produced is 1, then Temp\_Mantissa is right shifted by 1 bit and exp\_adjust incremented by 1. Else if Temp\_Mantissa is 0, it means that the output is 0 and so Temp\_Mantissa is assigned with 23-bit 0, exp\_adjust is assigned with 8-bit 0 while A\_sign is assigned with 1-bit 0. Else, if 23<sup>rd</sup> bit of exponent is 0, it represents the underflow condition. Here the Temp\_Mantissa is left shifted by 1 and exp\_adjust is decremented by 1. The variable Mantissa is assigned with the 22<sup>nd</sup> to 0<sup>th</sup> bit of Temp\_Mantissa. The variable exp\_adjust is assigned to the variable Exponent. The variable result is obtained after the concatenation of the variables A\_sign, Exponent and Mantissa. This variable result is the final output of the Floating-Point Adder.