# IEEE-754 FLOATING-POINT ADDER

AN IMPLEMENTATION OF THE FP ADDER IN C AND VERILOG

MOMEN ANANI AND MOHAMMAD MATAR

*12112318 - 12112195*

# INTRODUCTION

Floating point operations are critical in various computational tasks but are often not natively supported by HDLs like Verilog. This assignment aims to implement a floating-point adder first in C and then Verilog. The objective is to understand and experiment with the algorithm and functionality of floating-point arithmetic operations in both software and hardware contexts.

## Methodology

For the C code, it was straight forward implementation; just follow the algorithm and how the hardware unit works to achieve the desired functionality of the IEEE-754 adder.

We first implemented the MyFloat struct that was the basic block of the prototypes in the program

```
struct MyFloat {
unsigned int fraction;
unsigned int exponent;
unsigned int sign;
};
```

The main function prototypes for the code are initiated as follows

```
struct MyFloat toMyFloat(float f);
struct MyFloat MyFloatAdder(struct MyFloat input1, struct MyFloat input2);
float toFloat(struct MyFloat mf);
float random_float_range(float min, float max);
```
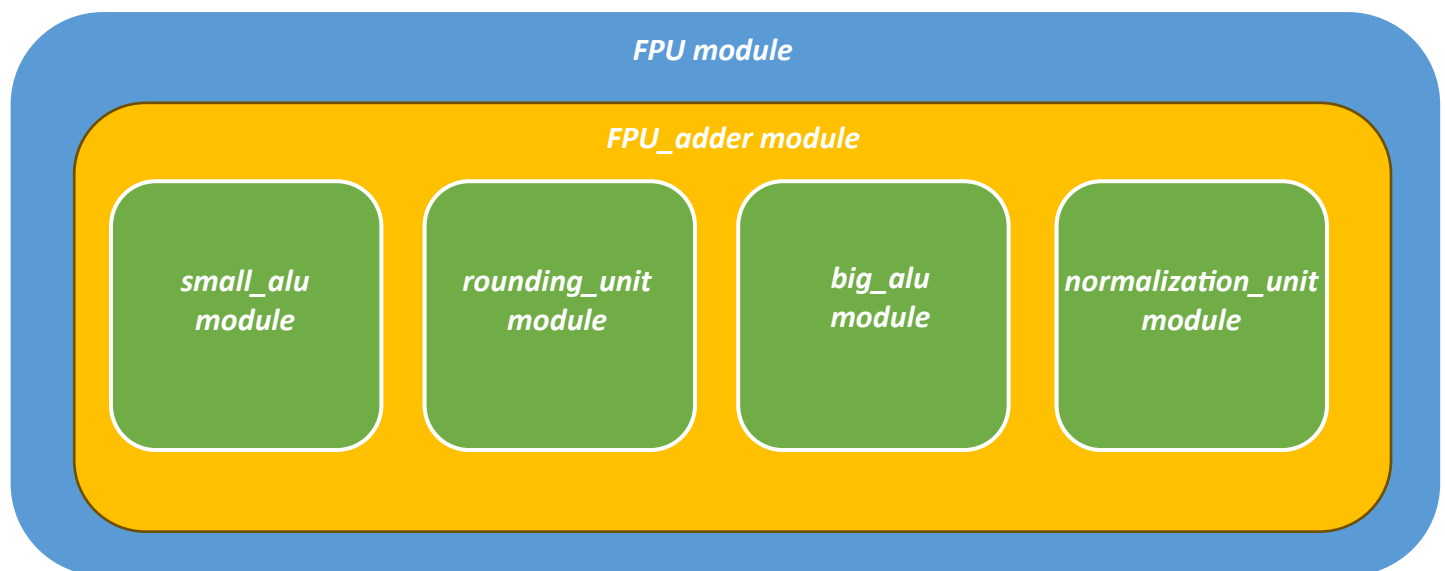
the toMyFloat function converts the standard float type to the user defined struct type.

The MyFloatAdder follows multiple steps mainly:

1. Checking for zero summation
2. Determine which exponent is larger and which is smaller to align them for summation
3. Perform addition/subtraction based on the sign bit
4. Normalizing and checking for overflow and underflow and handle NaN/inf exceptions

For the Verilog part, we took the same logic and functionality and made it resemble the hardware unit a bit more.

The main structure of the code consisted of the following hierarchical design:

As the name suggests for each module, the corresponding functionality is applied to it.

You can notice that <u>there is no control unit</u>. Well, yes and no. The control unit is implemented implicitly in each unit defined in the FPU_adder module. After analysis of the diagram, implementing an explicit control unit would make the code more complex since it is already simple. *The control unit controls signals that are simple enough to be calculated and figured within other modules easily.*

## Results and Discussion

For the C code, a random floating number was given as an input each time to test:

Here is the code result after running it:



*Figure 1. C code output for FP adder MyFloat.*

```
Test 1:
a: -3.930479e+029, b: -7.103183e+029
MyFloat Result - Sign: 1, Exponent: 226, Fraction: 6214322
MyFloat as float: -1.103366e+030  <---
Expected Result: -1.103366e+030  <---
Test 2:
a: -6.375012e+029, b: -8.724327e+029
MyFloat Result - Sign: 1, Exponent: 227, Fraction: 1603296
MyFloat as float: -1.509934e+030
Expected Result: -1.509934e+030
Test 3:
a: 8.419752e+029, b: 7.900326e+029
MyFloat Result - Sign: 0, Exponent: 227, Fraction: 2411114
MyFloat as float: 1.632008e+030
Expected Result: 1.632008e+030
Test 4:
a: 3.120518e+029, b: 7.285684e+029
MyFloat Result - Sign: 0, Exponent: 226, Fraction: 5383886
MyFloat as float: 1.040620e+030
Expected Result: 1.040620e+030
Test 5:
a: 2.882474e+029, b: 9.107638e+029
MyFloat Result - Sign: 0, Exponent: 226, Fraction: 7480173
MyFloat as float: 1.199011e+030
Expected Result: 1.199011e+030
Test 6:
a: 9.501327e+029, b: 6.766869e+029
MyFloat Result - Sign: 0, Exponent: 227, Fraction: 2376781
MyFloat as float: 1.626820e+030
Expected Result: 1.626820e+030
Test 7:
a: 2.668233e+029, b: -7.243568e+029
MyFloat Result - Sign: 1, Exponent: 225, Fraction: 3722204
MyFloat as float: -4.575335e+029
Expected Result: -4.575335e+029
Test 8:
a: -3.552660e+029, b: -8.497879e+029
MyFloat Result - Sign: 1, Exponent: 226, Fraction: 7560147
MyFloat as float: -1.205054e+030
Expected Result: -1.205054e+030
Test 9:
a: -6.036866e+029, b: -4.310740e+029
MyFloat Result - Sign: 1, Exponent: 226, Fraction: 5306334
MyFloat as float: -1.034761e+030
Expected Result: -1.034761e+030
Test 10:
a: -2.767114e+029, b: 2.755821e+028
MyFloat Result - Sign: 1, Exponent: 224, Fraction: 4801450
MyFloat as float: -2.491532e+029
Expected Result: -2.491532e+029
```

For the Verilog testbench that was run on Active-HDL software, this was the waveform:
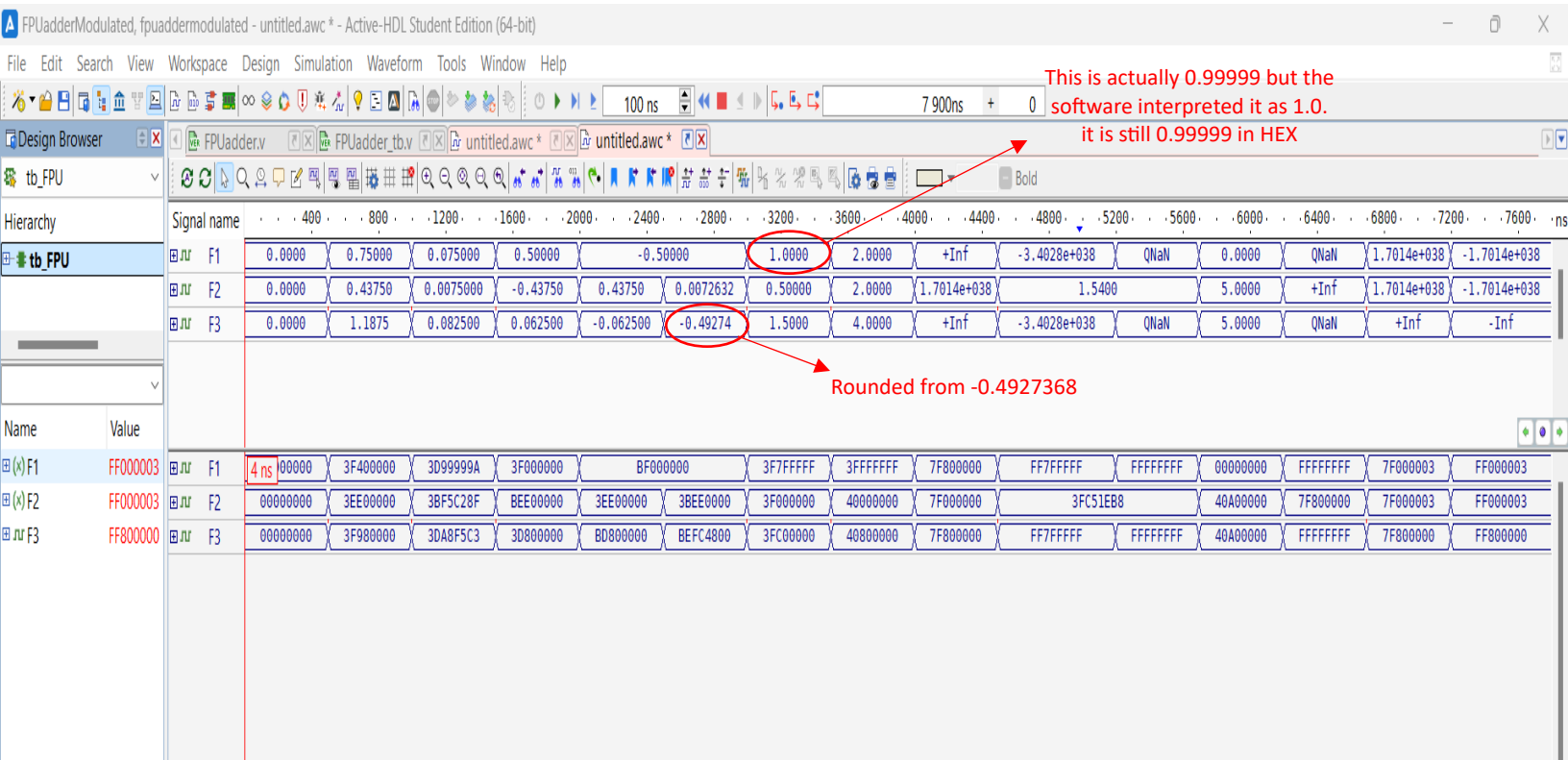


*Figure 2. waveform output for floating and HEX radix*

The test cases are chosen to test certain cases, all of which are explained in the testbench Verilog file.

## Conclusion

Implementing the floating-point adder logic in C was straightforward due to the simplicity of the language and the ability to directly translate the block diagram into code, on the HDL side, the modular approach in Verilog proved to be efficient, allowing each part of the operation to be handled within its own unit. This made the code more organized and manageable it can also help in adding other arithmetic units in the future. *The implicit control* within each unit streamlined the design, reducing complexity while maintaining functionality. Overall, this modular approach enhances the scalability and flexibility of the floating-point arithmetic unit design.