

Department of Computer Science  
Technical University of Cluj-Napoca



**Floating Point Arithmetic Logic Unit**  
*Laboratory activity 2024-2025*

Name: Suciu Andrei  
Group: 30431  
Email: [suciu.se.an@student.utcluj.ro](mailto:suciu.se.an@student.utcluj.ro)

Structure of Computer Systems



# Contents

<b>1</b>	<b>Project Overview</b>	<b>1</b>
1.1	Specifications . . . . .	1
1.2	Design . . . . .	1
<b>2</b>	<b>Bibliographic Research</b>	<b>2</b>
2.1	Arithmetic Logic Unit . . . . .	2
2.2	Floating Point . . . . .	2
2.3	IEEE 754 . . . . .	2
<b>3</b>	<b>Development Plan</b>	<b>3</b>
3.1	Development Overview . . . . .	3
3.2	Project Lab 1 . . . . .	3
3.3	Project Lab 2 . . . . .	3
3.4	Project Lab 3 . . . . .	3
3.5	Project Lab 4 . . . . .	3
3.6	Project Lab 5 . . . . .	3
<b>4</b>	<b>Floating Point ALU Structure</b>	<b>4</b>
4.1	Algorithm . . . . .	4
4.1.1	Floating Point Adder . . . . .	4
4.1.2	Floating Point Multiplier . . . . .	4
4.2	Design . . . . .	5
4.2.1	ALU . . . . .	5
4.2.2	Comparator . . . . .	6
4.2.3	Multiplier . . . . .	6
4.2.4	Control Unit . . . . .	6
4.2.5	Top Level Design . . . . .	8
4.3	Implementation . . . . .	9
<b>5</b>	<b>Testing</b>	<b>9</b>
5.1	Addition . . . . .	9
5.1.1	Regular Behaviour . . . . .	9
5.1.2	Special Cases . . . . .	10
5.2	Multiplication . . . . .	10
5.2.1	Regular Behaviour . . . . .	10
5.2.2	Special Cases . . . . .	11
<b>6</b>	<b>Code</b>	<b>12</b>
6.1	ALU . . . . .	12
6.2	Comparator . . . . .	12
6.3	Register with multiplexer . . . . .	12
6.4	Wallace Tree Multiplier . . . . .	12
6.5	Control Unit . . . . .	13
6.6	ALU Top level . . . . .	14
6.7	ROM . . . . .	14
6.8	Testbenches . . . . .	17
<b>7</b>	<b>Bibliography</b>	<b>23</b>

# 1 Project Overview

## 1.1 Specifications

Design an arithmetic logic unit on a 32-bit architecture capable of performing addition and multiplication operations on floating point numbers, encoded in the IEEE 754 standard. The ALU will be implemented on an FPGA, and must run and display various example operations. The ALU must be able to perform the following operations:

- 32-bit floating point addition
- 32-bit floating point multiplication

## 1.2 Design

The implementation of the arithmetic logic unit will be designed in VHDL, using Vivado. The ALU will be contained within a master project which will provide the input and output display functionalities according to the FPGA provided at the laboratory. The project will be broken down into multiple smaller sources, each of which will be implemented using a behavioural style. Finally, all the project sources will be combined structurally to form a cohesive unit.

## 2 Bibliographic Research

### 2.1 Arithmetic Logic Unit

In computing, an arithmetic logic unit (ALU) is a combinational digital circuit that is able to perform arithmetic on binary numbers. It is a fundamental building block of many other circuits, such as CPU or GPU. The ALU receives two input operands, an operation code, to decide which operation to perform, and optionally, a carry-in bit. Then, the result of the operation will be sent on the output signal.

### 2.2 Floating Point

In computing, numbers are represented in binary notation, as a series of bits, which can take values of either 1 or 0. In order to represent fractional parts, we must place a comma somewhere inside the number's representation. However, this greatly reduces the range of numbers which can be represented on a limited amount of bits. In order to increase the range, we can imagine the comma moving (floating) left or right across the binary representation, according to what number we want to represent.

Thus, we will reserve a few bits for storing the position of the comma, creating a floating point number.

### 2.3 IEEE 754

The IEEE 754 defines the most widely used standard for floating point numbers. Instead of storing the position of the comma, and the number itself, we will take advantage of scientific notation. In decimal scientific notation, the magnitude of the number is represented as a power of 10, and the "mantissa," essentially the core of the number, is stored as a fractional number between 1 and 10.

Moving this format to binary, we can observe that no matter the number, the mantissa will always start with a leading "1." Thus, we can assume this part of the number and save 1 bit. The magnitude represents the exponent of a power of 2. Furthermore, there is one more bit for the sign of the number, 0 for positive, and 1 for negative.

Overall, for a 32 bit number, the IEEE 754 standard formats the bits as such:

- 1 bit for the sign of the number
- 8 bits for the exponent
- 23 bits for the mantissa

The IEEE 754 standard also mentions the format for special values. There are five such numbers:

- Zero — All exponent bits and all mantissa bits are set to 0. Can be either positive or negative depending on the sign bit.
- Infinity — All exponent bits are set to 1, and all mantissa bits are set to 0. Can be either positive or negative depending on the sign bit.

- NaN (Not a Number) — All exponent bits are set to 1, and at least one mantissa bit is set to 1. The sign bit is not taken into consideration.

## **3 Development Plan**

### **3.1 Development Overview**

It is proposed that by each project meeting a different subsection of the Floating Point ALU will be completed, to be tested at the laboratory, and any irregularities in behaviour to be eliminated. By the end of the semester, the whole project must be completed and must be presented by the last week.

### **3.2 Project Lab 1**

The project overview and development plan will be finalised and presented.

### **3.3 Project Lab 2**

The supporting environment for the project must be implemented and tested on the FPGA. It must be capable of displaying numbers stored in a Read-Only Memory using a Seven-Segment Display, and must be capable of switching between multiple display modes using switches or buttons as input.

### **3.4 Project Lab 3**

The floating point addition subsection of the ALU must be implemented as a separate circuit and must be fully functional.

### **3.5 Project Lab 4**

The floating point multiplication subsection of the ALU must be implemented as a separate circuit and must be fully functional.

### **3.6 Project Lab 5**

Any remaining part of the project will be implemented, and the board's functionality will be fully tested. Any errors in the circuit's operations will be removed. The documentation will be updated as well, should the need arise.

## 4 Floating Point ALU Structure

The Operation of the Floating Point ALU is composed of three steps:

1. Check the validity of the inputs and the operation to be performed, and either proceed to the next step, or write the appropriate value and end the algorithm.
2. Perform the operation, with all required steps.
3. Realign the number to the IEEE 754 standard, and check for exponent overflow or underflow.

In order to implement the Floating Point ALU, we will use Registers, Comparators, a smaller ALU capable of integer operations, a Wallace Tree Multiplier, and the Control Unit, implemented as a Finite State Machine.

We can speed up the circuit by utilizing multiple comparators. For this implementation, we will use two comparators.

### 4.1 Algorithm

#### 4.1.1 Floating Point Adder

The Floating point Adder takes as inputs the two numbers, A and B, which are divided into their respective sign, exponent, and mantissa (with a 1 implicitly placed on the 24th bit).

First, we compare the exponents, and the mantissa of the number with the smaller exponent is shifted to the right and its exponent is incremented until the exponents are equal. This way, the numbers are aligned to the same exponent. Next, we perform the addition of the mantissa. If the signs are equal, we simply add the two mantissas together, otherwise, we subtract the smaller mantissa from the larger one, and use the sign of the number with the larger mantissa for the result.

Finally, we realign the resulting number to the IEEE-754 standard. If the mantissa has become too large, and there is a '1' on the 25th bit, we shift the result right and increment the exponent. If the mantissa is too small, we shift left and decrement the exponent until there is a '1' on exactly the 24th bit.

#### 4.1.2 Floating Point Multiplier

The Floating Point Multiplier takes as inputs the two numbers, A and B, which are divided into their respective sign, exponent, and mantissa (with a 1 implicitly placed on the 24th bit).

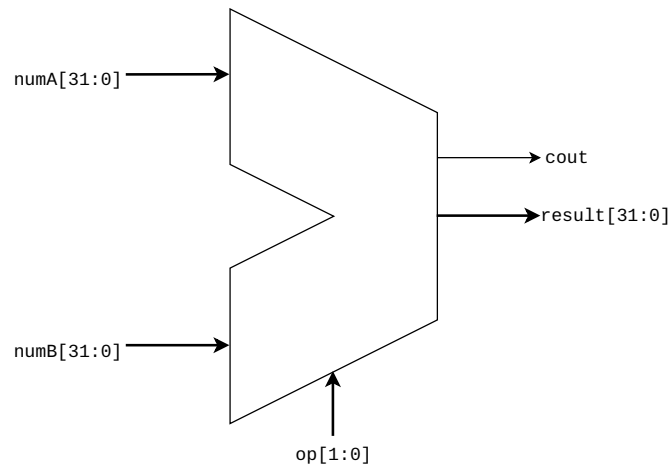
First, we add the two exponents together and subtract 127 from the resulting sum. This is because the exponents are implicitly stored as  $(127 + \text{exponent})$ . Next, we perform the multiplication on the mantissas. We can use a Wallace Tree Multiplier to perform the multiplication in a single clock cycle, which simplifies the circuit significantly.

Since we use a standard 32 bit Wallace Tree Multiplier and the mantissa has only 24 bits, we extend with zero to the right. The two mantissas will always start with a 1 on the MSB, and by multiplying them, the resulting number will also have a 1 on the 64th bit. We simply select the next most significant 23 bits for the result mantissa, with no need to shift and align the number.

## 4.2 Design

### 4.2.1 ALU

Figure 1: Simple Integer ALU



The ALU takes as input the two numbers, and a selection for which operation to be performed. Since it performs elementary operations, this is a combinational logic circuit, and requires no clock.

The outputs are the result of the operation, and an additional carry out bit, marked as *cout*.

The ALU is capable of performing four operations:

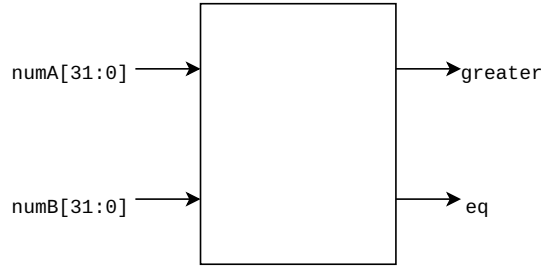
- Addition:  $result = numA + numB$
- Subtraction:  $result = numA - numB$
- Increment  $result = numA + 1$
- Decrement  $result = numA - 1$

In case of any overflow or underflow, the *cout* bit is set to 1, otherwise it remains 0.

### 4.2.2 Comparator

The comparator takes as input two numbers, and the outputs *eq* and *greater* are set to 1 if  $numA > numB$  and if  $numA = numB$  respectively.

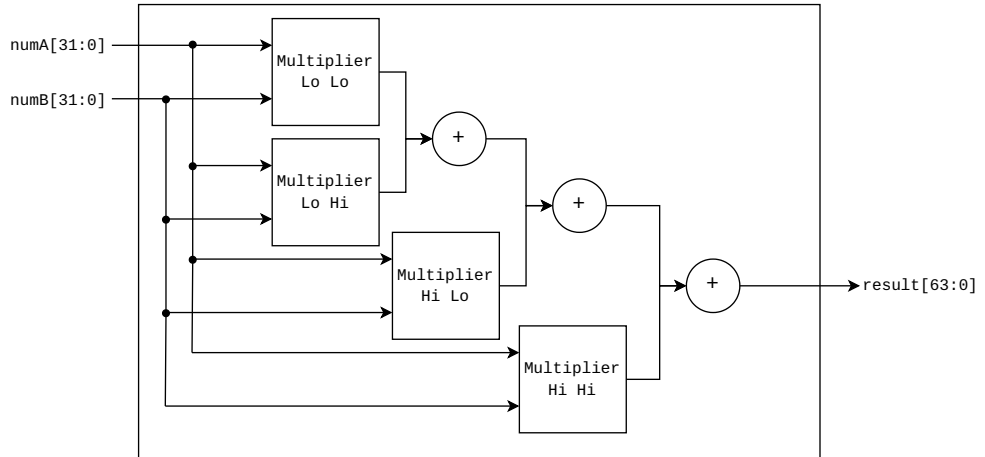
Figure 2: Integer Comparator



### 4.2.3 Multiplier

The mantissa multiplier was implemented as a wallace tree, capable of performing 32 bit multiplication. The result is on 64 bits. The desing was created recursively, starting from the 4 bit wallace tree multiplier. Each number is split into *Low* and *High*, and then multiplied, shifted, and added accordingly.

Figure 3: Structure of Wallace Tree Multiplier



The entity for the Components used is available [here](#).

### 4.2.4 Control Unit

The Control Unit is the most complex part of the floating point adder. It is a finite state machine which controls all the select signals for the various components of the design. It switches states synchronously, and based on the state and

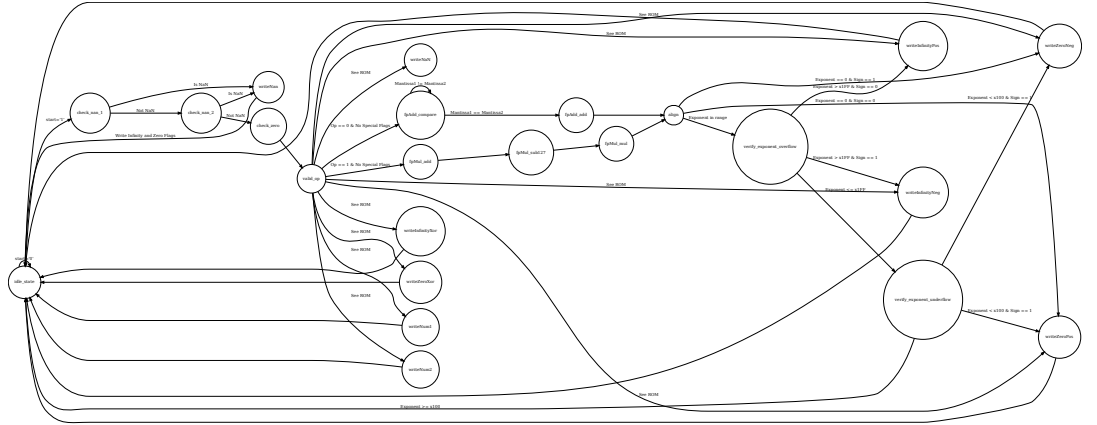


results from differing components, sets the select signals for the multiplexers, the write signals for the registers, and the operation for the ALU.

The Control Unit is a finite state machine with 22 states, divided into 6 groups:

- Idle: *idle\_state*  
Machine is idle and not busy performing any operation.
- Check Input: *check\_nan\_1*, *check\_nan\_2*, *check\_zero*, *valid\_op*  
Check the validity of the inputs and either proceed with the algorithm or write the result directly.
- Perform FP Addition: *fpAdd\_compare*, *fpAdd\_add*  
Shift mantissa and increment exponents until they are equal, then add the mantissas.
- Perform FP Multiplication: *fpMul\_add*, *fpMul\_sub127*, *fpMul\_mul*  
Add the two exponents together, subtract 127 from the result, and multiply the mantissas.
- Realign to Standard: *align*, *verify\_exponent\_overflow*, *verify\_exponent\_underflow*  
Realign mantissa to standard and check for overflow or underflow in the exponent.
- Write Special: *writeNan*, *writeInfinityXor*, *writeInfinityPos*, *writeInfinityNeg*, *writeZeroXor*, *writeZeroPos*, *writeZeroNeg*, *writeNum1*, *writeNum2*  
Used when operation is invalid.

Figure 4: The Control Unit transition diagram

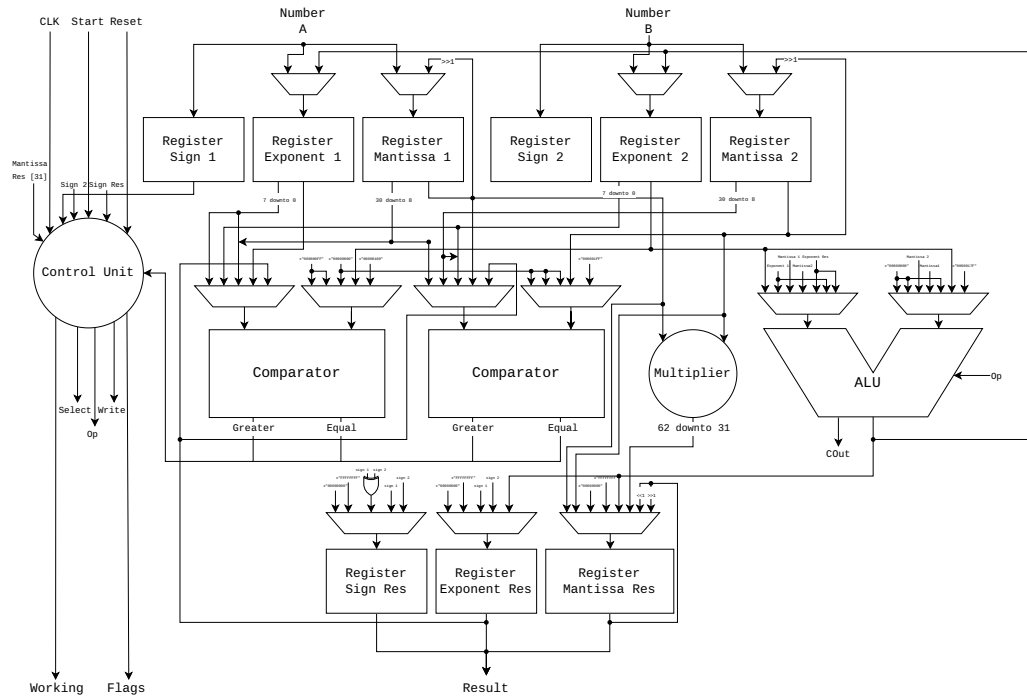


The entity for the Control Unit is available here.

#### 4.2.5 Top Level Design

Following the design specifications, all the components were connected. A separate register was used for each part of the inputs and outputs, and two comparators were used to speed up the process. Multiplexers were used to switch between the inputs to the different components. Since there is more space in the register than is needed to store the mantissa, it is stored left aligned, and padded with 0. Only the most significant bits from the mantissa register and multiplier result are used, ensuring there is never a mantissa overflow.

Figure 5: The top level design



### 4.3 Implementation

The code for the Comparator, ALU, Control Unit and Registers was written in a behavioural style. The code for the Multiplier, and Top Level was written in a structural style.

For ease of implementation, separate register components were written with a 4 : 1, and an 8 : 1 multiplexer included.

The Top Level Design is mostly composed of port maps, with a couple of multiplexers for the ALU and Comparators.

The Control Unit is a Finite State Machine composed of two processes:

- The first process calculates the next state of the machine. The switching between states is done synchronously. When checking the inputs, internal flags are set to 1 or 0. These flags and the operation input then form an address in a ROM which gives the next state.
- The second process represents a multitude of multiplexers and simple logic gates which calculate the output signals based on the current state and the inputs given. This process is asynchronous.

The code for the ROM is available [here](#).

## 5 Testing

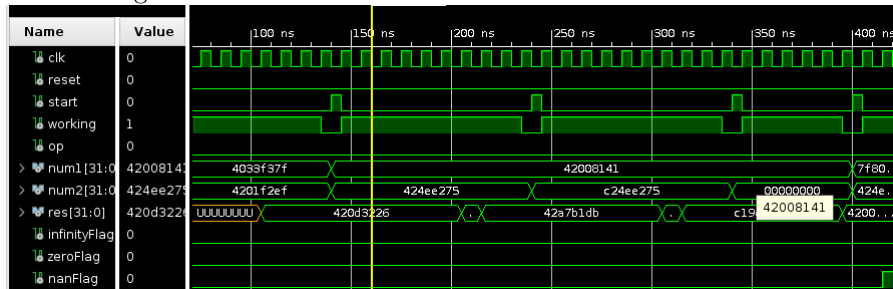
Testbenches were written for the Addition Operation and the Multiplication Operation. Multiple combinations of inputs were used, and the machine's ability to detect wrong inputs was also tested.

The inputs are captured immediately as the start value is read as 1 on a rising edge of the clock. A change in inputs afterwards will not affect the performance. While the machine is busy, the working bit is set to 1, otherwise, while it is idle it remains 0. Once the algorithm is finished and the working bit is set to 0, the output will remain stable. If the machine is started again, the output will still remain stable for at least *one* clock cycle. Afterwards, it cannot be guaranteed what will appear on the output while the working flag is active.

### 5.1 Addition

#### 5.1.1 Regular Behaviour

The FP ALU performs as expected under regular circumstances. It performs as expected when adding positive numbers, when adding negative numbers, and when adding zero to a number.

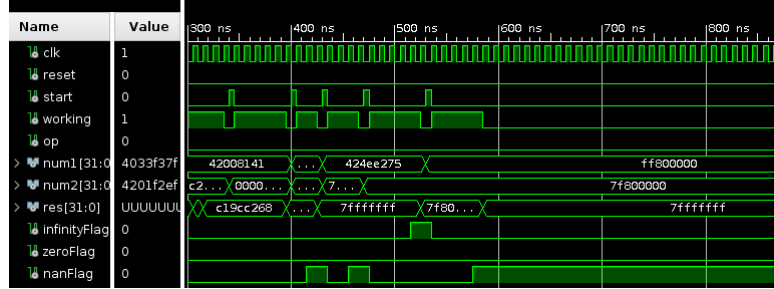


### 5.1.2 Special Cases

Based on the ROM in the Control Unit, the special cases for addition are as follows:

- $(+\infty) + (+\infty) = +\infty$
- $(+\infty) + (-\infty) = NaN$
- $(+\infty) + (\pm 0) = +\infty$
- $(\pm\infty) + x = \pm\infty$
- $(-\infty) + (-\infty) = -\infty$
- $(-\infty) + (\pm 0) = -\infty$
- $(\pm 0) + (\pm 0) = +0$
- $(\pm 0) + (\mp 0) = +0$
- $(\pm 0) + x = +0$

Also including their reciprocals. These special cases were tested, and the correct result was given.

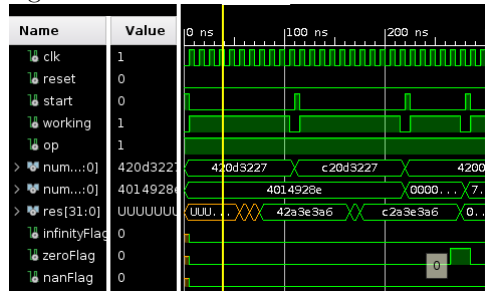


The code for the Add Testbench is available [here](#).

## 5.2 Multiplication

### 5.2.1 Regular Behaviour

Under regular circumstances the multiplication is performed correctly. When multiplying numbers with the same sign, the resulting sign is positive, and when multiplying with 0, the result will also be 0, with the same logic applied for the sign.

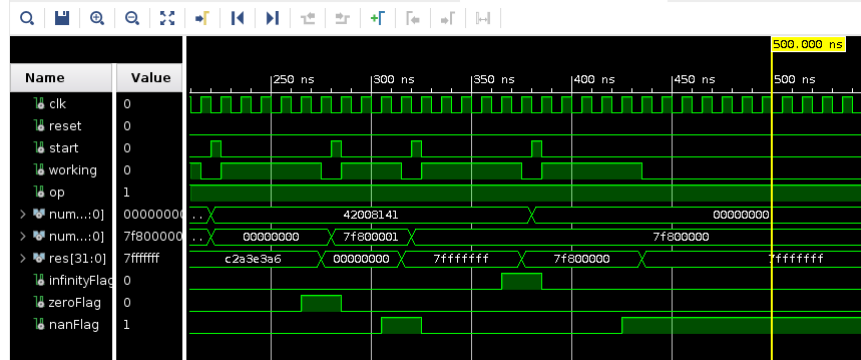


### 5.2.2 Special Cases

Based on the ROM in the Control Unit, the special cases for addition are as follows:

- $(\pm\infty) * (\pm\infty) = +\infty$
- $(\pm\infty) * (\mp\infty) = -\infty$
- $(\pm\infty) * (\pm 0) = NaN$
- $(\pm\infty) * (\mp 0) = NaN$
- $(\pm\infty) * x = (XOR)\infty$
- $(-\infty) + (\pm 0) = -\infty$
- $(\pm 0) * (\pm 0) = +0$
- $(\pm 0) * (\mp 0) = -0$
- $(\pm 0) * x = (XOR)0$

Also including their reciprocals. These special cases were tested, and the correct result was given.



The code for the Mul Testbench is available [here](#).

## 6 Code

### 6.1 ALU

```
1  entity alu32bit is
2  Port (
3      numA : in STD_LOGIC_VECTOR (31 downto 0);
4      numB : in STD_LOGIC_VECTOR (31 downto 0);
5      op   : in STD_LOGIC_VECTOR (1 downto 0);
6      output : out STD_LOGIC_VECTOR (31 downto 0);
7      cout : out std_logic;
8  end alu32bit;
```

### 6.2 Comparator

```
1  entity comparator32bit is
2  Port ( numA : in STD_LOGIC_VECTOR (31 downto 0);
3        numB : in STD_LOGIC_VECTOR (31 downto 0);
4        eq   : out STD_LOGIC;
5        greater : out STD_LOGIC);
6  end comparator32bit;
```

### 6.3 Register with multiplexer

```
1  entity register32bitW4x1Mux is
2  Port ( A : in STD_LOGIC_VECTOR (31 downto 0);
3        B : in STD_LOGIC_VECTOR (31 downto 0);
4        C : in STD_LOGIC_VECTOR (31 downto 0);
5        D : in STD_LOGIC_VECTOR (31 downto 0);
6        S : in std_logic_vector (1 downto 0);
7        write : in STD_LOGIC;
8        clk : in STD_LOGIC;
9        output : out STD_LOGIC_VECTOR (31 downto 0));
10 end register32bitW4x1Mux;
```

### 6.4 Wallace Tree Multiplier

```
1  entity multiplier32bit is
2  Port ( numA : in STD_LOGIC_VECTOR (31 downto 0);
3        numB : in STD_LOGIC_VECTOR (31 downto 0);
4        result : out STD_LOGIC_VECTOR (63 downto 0));
5  end multiplier32bit;
```

## 6.5 Control Unit

```
1  entity controlUnit is
2  Port (
3      clk : in STD_LOGIC;
4      reset: in std_logic;
5      start:in std_logic;
6      op:in std_logic;
7
8      writeSign1:out std_logic;
9      writeSign2:out std_logic;
10     writeSignRes:out std_logic;
11
12     writeExponent1:out std_logic;
13     writeExponent2:out std_logic;
14     writeExponentRes:out std_logic;
15
16     writeMantissa1:out std_logic;
17     writeMantissa2:out std_logic;
18     writeMantissaRes:out std_logic;
19
20     selectSign1:out std_logic_vector(1 downto 0);
21     selectExponent1:out std_logic_vector(1 downto 0);
22     selectMantissa1:out std_logic_vector(1 downto 0);
23
24     selectSign2:out std_logic_vector(1 downto 0);
25     selectExponent2:out std_logic_vector(1 downto 0);
26     selectMantissa2:out std_logic_vector(1 downto 0);
27
28     selectSignRes:out std_logic_vector(2 downto 0);
29     selectExponentRes:out std_logic_vector(2 downto 0);
30     selectMantissaRes:out std_logic_vector(2 downto 0);
31
32     outSign1:in std_logic;
33     outSign2:in std_logic;
34     outSignRes:in std_logic;
35     outMantissaRes31:in std_logic;
36
37     comparator1Select:out std_logic_vector(2 downto 0);
38     comparator1Eq:in std_logic;
39     comparator1Greater:in std_logic;
40
41     comparator2Select:out std_logic_vector(2 downto 0);
42     comparator2Eq:in std_logic;
43     comparator2Greater:in std_logic;
44
45     alu32bitOp:out std_logic_vector(1 downto 0);
46     alu32bitSelect:out std_logic_vector(2 downto 0);
47     alu32bitCout: in std_logic;
48
49     infinityFlag:out std_logic;
50     zeroFlag:out std_logic;
51     nanFlag:out std_logic;
52
53     working:out std_logic
54 );
55 end controlUnit;
```

## 6.6 ALU Top level

```
1  entity aluFloat is
2  Port (
3      clk : in STD_LOGIC;
4      start:in std_logic;
5      reset:in std_logic;
6      op : in STD_LOGIC;
7
8      num1 : in STD_LOGIC_VECTOR (31 downto 0);
9      num2 : in STD_LOGIC_VECTOR (31 downto 0);
10
11     res : out STD_LOGIC_VECTOR (31 downto 0);
12     working:out std_logic;
13     infinityFlag:out std_logic;
14     zeroFlag:out std_logic;
15     nanFlag:out std_logic
16 );
17 end aluFloat;
```

## 6.7 ROM

```
1  case internalFlags is
2  --      internalFlags<=op & infinityPos1 & infinityNeg1 &
3  --      infinityPos2 & infinityNeg2 &
4  --      zeroPos1 & zeroNeg1 & zeroPos2 & zeroNeg2;
5      when "00000000"=>
6          state<=fpAdd_compare;
7
8      -- +infinity + x
9      when b"0_10_10_00_00"=> -- +inf + +inf
10         state<=writeInfinityPos;
11     when b"0_10_01_00_00"=> -- +inf + -inf
12         state<=writeNan;
13     when b"0_10_00_00_10"=> -- +inf + +0
14         state<=writeInfinityPos;
15     when b"0_10_00_00_01"=> -- +inf + -0
16         state<=writeInfinityPos;
17     when b"0_10_00_00_00"=> -- +inf + x
18         state<=writeInfinityPos;
19
20     -- -infinity + x
21     when b"0_01_10_00_00"=> -- -inf + +inf
22         state<=writeNan;
23     when b"0_01_01_00_00"=> -- -inf + -inf
24         state<=writeInfinityNeg;
25     when b"0_01_00_00_10"=> -- -inf + +0
26         state<=writeInfinityNeg;
27     when b"0_01_00_00_01"=> -- -inf + -0
28         state<=writeInfinityNeg;
29     when b"0_01_00_00_00"=> -- -inf + x
30         state<=writeInfinityNeg;
31
32     -- +0 + x
33     when b"0_00_10_10_00"=> -- +0 + +inf
34         state<=writeInfinityPos;
35     when b"0_00_01_10_00"=> -- +0 + -inf
36         state<=writeInfinityNeg;
37     when b"0_00_00_10_10"=> -- +0 + +0
38         state<=writeZeroPos;
```



```

38     when b"0_00_00_10_01"=> -- +0 + -0
39         state<=writeZeroPos;
40     when b"0_00_00_10_00"=> -- +0 + x
41         state<=writeNum2;
42
43     -- -0 + x
44     when b"0_00_10_01_00"=> -- -0 + +inf
45         state<=writeInfinityPos;
46     when b"0_00_01_01_00"=> -- -0 + -inf
47         state<=writeInfinityNeg;
48     when b"0_00_00_01_10"=> -- -0 + +0
49         state<=writeZeroPos;
50     when b"0_00_00_01_01"=> -- -0 + -0
51         state<=writeZeroNeg;
52     when b"0_00_00_01_00"=> -- -0 + x
53         state<=writeNum2;
54
55     -- x + special
56     when b"0_00_10_00_00"=> -- x + +inf
57         state<=writeInfinityPos;
58     when b"0_00_01_00_00"=> -- x + -inf
59         state<=writeInfinityNeg;
60     when b"0_00_00_00_10"=> -- x + +0
61         state<=writeNum1;
62     when b"0_00_00_00_01"=> -- x + -0
63         state<=writeNum1;
64
65     -- internalFlags<=op & infinityPos1 & infinityNeg1 &
66     infinityPos2 & infinityNeg2 &
67     zeroPos1 & zeroNeg1 & zeroPos2 & zeroNeg2;
68     when "100000000"=>
69         state<=fpMul_add;
70
71     -- +infinity * x
72     when b"1_10_10_00_00"=> -- +inf * +inf
73         state<=writeInfinityPos;
74     when b"1_10_01_00_00"=> -- +inf * -inf
75         state<=writeInfinityNeg;
76     when b"1_10_00_00_10"=> -- +inf * +0
77         state<=writeNan;
78     when b"1_10_00_00_01"=> -- +inf * -0
79         state<=writeNan;
80     when b"1_10_00_00_00"=> -- +inf * x
81         state<=writeInfinityXor;
82
83     -- -infinity * x
84     when b"1_01_10_00_00"=> -- -inf * +inf
85         state<=writeInfinityNeg;
86     when b"1_01_01_00_00"=> -- -inf * -inf
87         state<=writeInfinityPos;
88     when b"1_01_00_00_10"=> -- -inf * +0
89         state<=writeNan;
90     when b"1_01_00_00_01"=> -- -inf * -0
91         state<=writeNan;
92     when b"1_01_00_00_00"=> -- -inf * x
93         state<=writeInfinityXor;
94
95     -- +0 * x
96     when b"1_00_10_10_00"=> -- +0 * +inf
97         state<=writeNan;
98     when b"1_00_01_10_00"=> -- +0 * -inf
99         state<=writeNan;

```

```

99      when b"1_00_00_10_10"=> -- +0 * +0
100         state<=writeZeroPos;
101      when b"1_00_00_10_01"=> -- +0 * -0
102         state<=writeZeroNeg;
103      when b"1_00_00_10_00"=> -- +0 * x
104         state<=writeZeroXor;
105
106      -- -0 * x
107      when b"1_00_10_01_00"=> -- -0 * +inf
108         state<=writeNan;
109      when b"1_00_01_01_00"=> -- -0 * -inf
110         state<=writeNan;
111      when b"1_00_00_01_10"=> -- -0 * +0
112         state<=writeZeroNeg;
113      when b"1_00_00_01_01"=> -- -0 * -0
114         state<=writeZeroPos;
115      when b"1_00_00_01_00"=> -- -0 * x
116         state<=writeZeroXor;
117
118      -- x * special
119      when b"1_00_10_00_00"=> -- x * +inf
120         state<=writeInfinityXor;
121      when b"1_00_01_00_00"=> -- x * -inf
122         state<=writeInfinityXor;
123      when b"1_00_00_00_10"=> -- x * +0
124         state<=writeZeroXor;
125      when b"1_00_00_00_01"=> -- x * -0
126         state<=writeZeroXor;
127
128      when others=>
129         state<=writeNan;
130  end case;

```

## 6.8 Testbenches

```
1  library IEEE;
2  use IEEE.STD_LOGIC_1164.ALL;
3
4  entity tbAddFloat is
5  -- Port ( );
6  end tbAddFloat;
7
8  architecture Behavioral of tbAddFloat is
9
10 component aluFloat is
11 Port (
12     clk : in STD_LOGIC;
13     start:in std_logic;
14     reset:in std_logic;
15     op : in STD_LOGIC;
16
17     num1 : in STD_LOGIC_VECTOR (31 downto 0);
18     num2 : in STD_LOGIC_VECTOR (31 downto 0);
19
20     res : out STD_LOGIC_VECTOR (31 downto 0);
21     working:out std_logic;
22     infinityFlag:out std_logic;
23     zeroFlag:out std_logic;
24     nanFlag:out std_logic
25 );
26 end component;
27
28 signal clk,reset,start,working,op:std_logic:= '0';
29 signal num1,num2,res:std_logic_vector(31 downto 0);
30
31 signal infinityFlag,zeroFlag,nanFlag:std_logic:= '0';
32
33 function to_string ( a: std_logic_vector) return string is
34 variable b : string (1 to a'length) := (others => NUL);
35 variable stri : integer := 1;
36 begin
37     for i in a'range loop
38         b(stri) := std_logic'image(a((i)))(2);
39         stri := stri+1;
40     end loop;
41 return b;
42 end function;
43
44
45 begin
46
47 clk<=not clk after 5ns;
48
49 tb:aluFloat port map(
50     clk=>clk,
51     reset=>reset,
52     start=>start,
53     op=>op,
54
55     num1=>num1,
56     num2=>num2,
57     res=>res,
58     working=>working,
59
60     infinityFlag=>infinityFlag,
```

```

61     zeroFlag=>zeroFlag,
62     nanFlag=>nanFlag
63 );
64
65 process
66 begin
67
68     num1<=b"0_10000000_01100111111001101111111"; --
69         2.8117368221282958984375
70     num2<=b"0_10000100_00000011111001011101111"; --
71         32.487239837646484375
72     start<='1';
73
74     --expecting 35.29897665977478027344
75     -- 420d3226
76     wait on clk;
77     start<='0';
78     wait until working='0';
79
80     report to_string(res);
81     wait on clk;
82     =====
83     num1<=b"01000010000000001000000101000001"; -- 32.126225
84     num2<=b"01000010010011101110001001110101"; -- 51.721149
85     start<='1';
86
87     -- expecting 83.847374
88     -- 42a7b1db
89     wait on clk;
90     start<='0';
91     wait until working='0';
92     report to_string(res);
93     wait on clk;
94
95     =====
96     num1<=b"01000010000000001000000101000001"; -- 32.126225
97     num2<=b"11000010010011101110001001110101"; -- -51.721149
98     start<='1';
99
100    -- expecting -19.594924
101    -- c19cc268
102    wait on clk;
103    start<='0';
104    wait until working='0';
105    report to_string(res);
106    wait on clk;
107
108    =====
109    num1<=b"01000010000000001000000101000001"; --32.126225
110    num2<=b"00000000000000000000000000000000"; --0
111    start<='1';
112
113    -- expecting 32.126225
114    -- 42008141
115    wait on clk;
116    start<='0';
117    wait until working='0';
118    report to_string(res);
119    wait on clk;
120

```

```

121 -----
122 num1<=b"0_11111111_000000000000000000000001"; --NaN
123 num2<=b"01000010010011101110001001110101"; --51.721149
124 start<='1';
125
126 -- expecting NaN
127 -- 7fffffff
128 wait on clk;
129 start<='0';
130 wait until working='0';
131 report to_string(res);
132 wait on clk;
133 -----
134 num1<=b"01000010010011101110001001110101"; --51.721149
135 num2<=b"0_11111111_000000000000000000000001"; --NaN
136 start<='1';
137
138 -- expecting NaN
139 -- 7fffffff
140 wait on clk;
141 start<='0';
142 wait until working='0';
143 report to_string(res);
144 wait on clk;
145 -----
146 num1<=b"01000010010011101110001001110101"; --51.721149
147 num2<=b"0_11111111_000000000000000000000000"; -- +inf
148 start<='1';
149
150 -- expecting +inf
151 -- 7f800000
152 wait on clk;
153 start<='0';
154 wait until working='0';
155 report to_string(res);
156 wait on clk;
157
158 -----
159 num1<=b"1_11111111_000000000000000000000000"; -- -inf
160 num2<=b"0_11111111_000000000000000000000000"; -- +inf
161 start<='1';
162
163 -- expecting +NaN
164 -- 7fffffff
165 wait on clk;
166 start<='0';
167 wait until working='0';
168 report to_string(res);
169 wait on clk;
170
171 wait for 1000ns;
172 end process;
173
174 end Behavioral;

```

Listing 1: Add Testbench

```

1  library IEEE;
2  use IEEE.STD_LOGIC_1164.ALL;
3
4  entity tbMulFloat is
5  -- Port ( );
6  end tbMulFloat;
7
8  architecture Behavioral of tbMulFloat is
9
10 component aluFloat is
11 Port (
12     clk : in STD_LOGIC;
13     start:in std_logic;
14     reset:in std_logic;
15     op  : in STD_LOGIC;
16
17     num1 : in STD_LOGIC_VECTOR (31 downto 0);
18     num2 : in STD_LOGIC_VECTOR (31 downto 0);
19
20     res : out STD_LOGIC_VECTOR (31 downto 0);
21     working:out std_logic;
22     infinityFlag:out std_logic;
23     zeroFlag:out std_logic;
24     nanFlag:out std_logic
25 );
26 end component;
27
28 signal clk,reset,start,working,op:std_logic:='0';
29 signal num1,num2,res:std_logic_vector(31 downto 0);
30
31 signal infinityFlag,zeroFlag,nanFlag:std_logic:='0';
32
33 function to_string ( a: std_logic_vector) return string is
34 variable b : string (1 to a'length) := (others => NUL);
35 variable stri : integer := 1;
36 begin
37     for i in a'range loop
38         b(stri) := std_logic'image(a((i)))(2);
39         stri := stri+1;
40     end loop;
41 return b;
42 end function;
43
44
45 begin
46
47 clk<=not clk after 5ns;
48
49 tb:aluFloat port map(
50     clk=>clk,
51     reset=>reset,
52     start=>start,
53     op=>op,
54
55     num1=>num1,
56     num2=>num2,
57     res=>res,
58     working=>working,
59
60     infinityFlag=>infinityFlag,
61     zeroFlag=>zeroFlag,
62     nanFlag=>nanFlag

```

```

63 );
64
65 process
66 begin
67     op<='1';
68     start<='1';
69     num1<="01000010000011010011001000100111"; --
        35.298976898193359375
70     num2<="01000000000101001001001010001110"; --
        2.321444988250732421875
71
72     --expecting 81.94463245721374278219
73     --42a3e3a6
74     wait on clk;
75     start<='0';
76     wait until working='0';
77
78     report to_string(res);
79     wait on clk;
80     =====
81
82     op<='1';
83     start<='1';
84     num1<="11000010000011010011001000100111"; --
        -35.298976898193359375
85     num2<="01000000000101001001001010001110"; --
        2.321444988250732421875
86
87     --expecting -81.94463245721374278219
88     --c2a3e3a6
89     wait on clk;
90     start<='0';
91     wait until working='0';
92
93     report to_string(res);
94     wait on clk;
95     =====
96     num1<=b"01000010000000001000000101000001"; --32.126225
97     num2<=b"00000000000000000000000000000000"; --0
98     start<='1';
99
100     -- expecting 0
101     -- 00000000
102     wait on clk;
103     start<='0';
104     wait until working='0';
105     report to_string(res);
106     wait on clk;
107     =====
108
109
110     num1<=b"01000010000000001000000101000001"; --32.126225
111     num2<=b"0_11111111_000000000000000000000001"; --NaN
112     start<='1';
113
114     -- expecting NaN
115     -- 7fffffff
116     wait on clk;
117     start<='0';
118     wait until working='0';
119     report to_string(res);
120     wait on clk;

```

```

121
122
123
124 num1<=b"010000100000000001000000101000001"; --32.126225
125 num2<=b"0_11111111_000000000000000000000000"; --+inf
126 start<='1';
127
128 -- expecting +inf
129 -- 7f800000
130 wait on clk;
131 start<='0';
132 wait until working='0';
133 report to_string(res);
134 wait on clk;
135
136
137
138 num1<=b"00000000000000000000000000000000"; --0
139 num2<=b"0_11111111_000000000000000000000000"; --+inf
140 start<='1';
141
142 -- expecting NaN
143 -- 7fffffff
144 wait on clk;
145 start<='0';
146 wait until working='0';
147 report to_string(res);
148 wait on clk;
149
150
151 wait for 1000ns;
152
153 end process;
154 end Behavioral;

```

Listing 2: Mul Testbench



## 7 Bibliography

### References

- [1] Steven Petryk, Adding IEEE-754 Floating Point Numbers, Oct 2015
- [2] Wikipedia, IEEE 754, Oct 2024
- [3] Jan Misali, how floating point works, May 2022