

CSE 306 Assignment 1
A1
Group 5

Md. Zarif Ul Alam : 1705010
Jamilus Sheium : 1705012
Naeem Ahmed : 1705014
Md. Shadmim Hasan Sifat : 1705021
Solaiman Ahmed : 1705022

May 28, 2021

1 Introduction

Floating-Point Adder is a combinational circuit which takes two floating points as inputs and provides their sum which is another floating point as output. Its implementation requires some basic n bits adder, subtractor, shifter, multiplexer, comparator and some other basic gates.

Floating-Point Adder is designed to perform “Floating Point arithmetic” which is by far the most used way of approximating real number arithmetic for performing numerical calculations on modern computers.

The floating-point numbers representation is based on the scientific notation: the decimal point is not set in a fixed position in the bit sequence, but its position is indicated as a base power.

All the floating-point numbers are composed by four components:

- Sign: it indicates the sign of the number (0 positive and 1 negative)
- Significant: it sets the value of the number
- Exponent: it contains the value of the base power (biased)
- Base: the base (or radix) is implied and it is common to all the numbers (2 for binary numbers)

The steps involved in the design of a Floating-Point Adder are as follows:

1. Extracting signs, exponents and fractions of both A and B numbers.
2. Treating the special cases:
 - Operations with A or B equal to zero
 - Operations with $\pm\infty$
 - Operations with NaN

(For simplicity of our design, we are only dealing with the first case)

3. Finding out what type of numbers are given:
 - Normalized
 - Unnormalized

(For simplicity of our design, we are assuming that numbers are given in normalized form)

4. By comparing the exponent of the numbers, finding out their difference which is actually the required shifting amount and also the smaller & larger number.
5. Shifting the lower exponent number fraction to the right $[\text{Exp1} - \text{Exp2}]$ bits. Setting the output exponent as the highest exponent.

6. Working with the operation symbol and both signs to calculate the output sign and determine the operation to do.
7. Addition of the numbers and detection of overflow (carry bit)
8. Standardizing fraction shifting it to the left up the first one will be at the first position and updating the value of the exponent according with the carry bit and the shifting over the fraction.

2 Problem Specification

Design a floating-point adder circuit which takes two floating points as inputs and provides their sum, another floating point as output.

Each floating point will be 16 bits long with following representation:

Sign	Exponent	Fraction
1 Bit	4 Bits	11 Bits

3 Flowchart

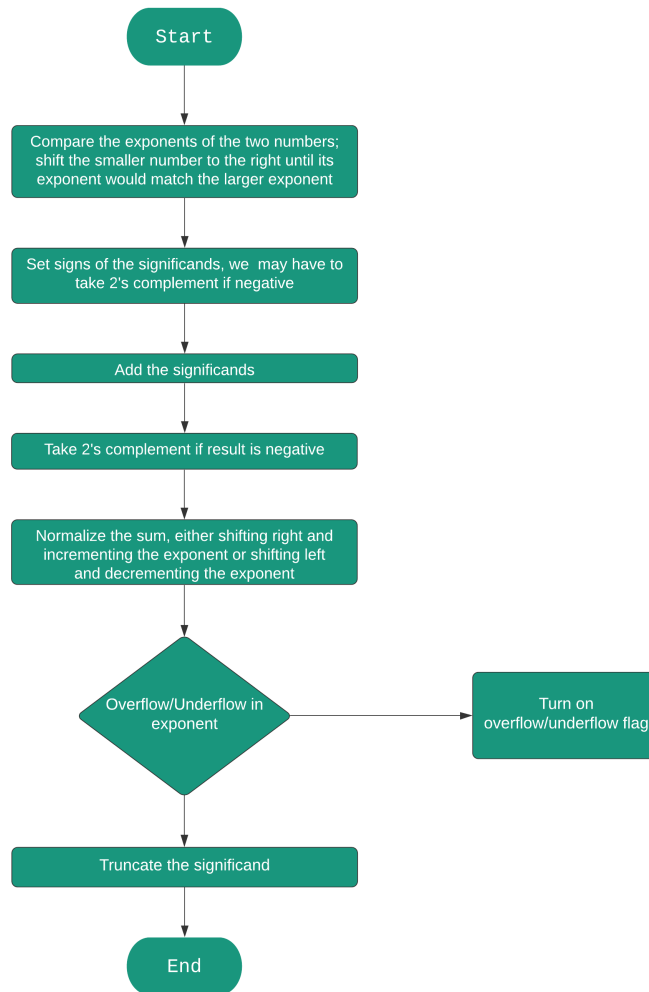


Figure 1: Flowchart of Floating Point Adder

4 ICs used with count as a chart

5 Simulator

Logisim Version 2.7.1

6 Discussion

While implementing the circuit, we had to change our design several times. Some designs required a lot of ICs. To optimize number of ICs in our design, we had to discard those. Again, we invested enough of our time in cross-checking corner cases, overflow, underflow conditions for each of sample test cases cautiously. However, because of truncating and rounding up the sum, sometimes we encountered precision loss of minimum one or maximum two bits. Also, in some cases we reused some logic gate outputs to minimize the number of ICs used. Considering all these aspects, we finally implemented the most optimized design we could find.