

IEEE 754 arithmetic using lists

Kapil Aggarwal (16305R010)
Hareesh Kumar (16305R013)

November 22, 2017

Contents

1	Problem Statement	2
2	Motivation	3
3	Approach	4
4	Implementation Details	5
5	Algorithms	7
5.1	Addition	7
5.2	Subtraction	7
5.3	Multiplication	8
6	Working Example	9
7	Conclusion	10

Chapter 1

Problem Statement

IEEE 754 standard specifies interchange and arithmetic formats and methods for binary and decimal floating-point arithmetic in computer programming environments. An implementation of a floating-point system conforming to this standard may be realized entirely in software, entirely in hardware, or in any combination of software and hardware ([1]). The goal of this project is to develop a application capable of doing the floating point arithmetic in Haskell using Lists.

Chapter 2

Motivation

IEEE 754 standard provides a discipline for performing floating-point computation that yields results independent of whether the processing is done in hardware, software, or a combination of the two. This standard defines a ways for systems to perform binary and decimal floating-point arithmetic.

Chapter 3

Approach

The IEEE 754 format contains the representations of floating point number in 2 formats as shown in Figure1.

Figure 3.1: IEEE 754 Format - Number Representations

Sign S	8 bit - biased Exponent E	23 bits - unsigned fraction P
-------------	--------------------------------	---------------------------------

(a) IEEE single precision data format

Sign S	11 bit - biased Exponent E	52 bits - unsigned fraction p
-------------	---------------------------------	---------------------------------

(b) IEEE double precision data format

Floating points number are represented in normalized format i.e.

$$-1^S \times (1.0 + 0.M) \times 2_{E-bias}$$

where M is mantissa and E is exponent and S is sign bit.

We begin the implementation of the application by first converting the input numbers that are in decimal format into the following format

(sign, exponent, mantissa).

After that we start implementing the arithmetic operations. We have implemented the following 3 arithmetic operations as functions in our project.

1. Addition
2. Subtraction
3. Multiplication

Above functions take 2 operand and return the results in (sign, exponent, mantissa) format. We have to convert it back to the decimal format.

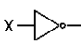
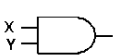
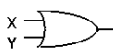
Chapter 4

Implementation Details

In the microprocessor the operation are implemented using gates. So we started the implementation by first designing following basic gates in boolean algebra as these are basic building blocks of any boolean circuit.

1. AND
2. OR
3. XOR
4. NOT

Figure 4.1: Boolean Gates

NOT	X'	\bar{X}	$\sim X$		$X \rightarrow Y$	<table> <tr> <th>X</th> <th>Y</th> </tr> <tr> <td>0</td> <td>1</td> </tr> <tr> <td>1</td> <td>0</td> </tr> </table>	X	Y	0	1	1	0									
X	Y																				
0	1																				
1	0																				
AND	$X \cdot Y$	XY	$X \wedge Y$		$X, Y \rightarrow Z$	<table> <tr> <th>X</th> <th>Y</th> <th>Z</th> </tr> <tr> <td>0</td> <td>0</td> <td>0</td> </tr> <tr> <td>0</td> <td>1</td> <td>0</td> </tr> <tr> <td>1</td> <td>0</td> <td>0</td> </tr> <tr> <td>1</td> <td>1</td> <td>1</td> </tr> </table>	X	Y	Z	0	0	0	0	1	0	1	0	0	1	1	1
X	Y	Z																			
0	0	0																			
0	1	0																			
1	0	0																			
1	1	1																			
OR	$X + Y$		$X \vee Y$		$X, Y \rightarrow Z$	<table> <tr> <th>X</th> <th>Y</th> <th>Z</th> </tr> <tr> <td>0</td> <td>0</td> <td>0</td> </tr> <tr> <td>0</td> <td>1</td> <td>1</td> </tr> <tr> <td>1</td> <td>0</td> <td>1</td> </tr> <tr> <td>1</td> <td>1</td> <td>1</td> </tr> </table>	X	Y	Z	0	0	0	0	1	1	1	0	1	1	1	1
X	Y	Z																			
0	0	0																			
0	1	1																			
1	0	1																			
1	1	1																			

After implementing the basic boolean circuits , we have to implement the following arithmetic circuit using the basic gates.

1. Half ADDER
2. Full ADDER
3. Half SUBTRACTOR
4. Full SUBTRACTOR
5. Multiplier

Since the above implemented circuits are capable of performing only single bit operations , we have to implement the circuits that can perform n-bit calculations using single bit operations.

Lets have a look at how above circuits are implemented using boolean get Half adder circuit are implemented.

Figure 4.2: Half Adder

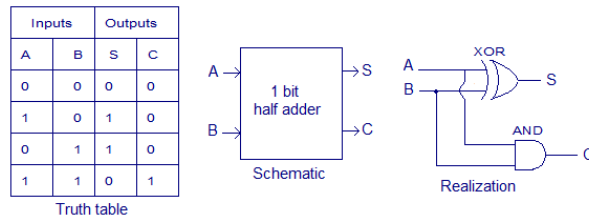
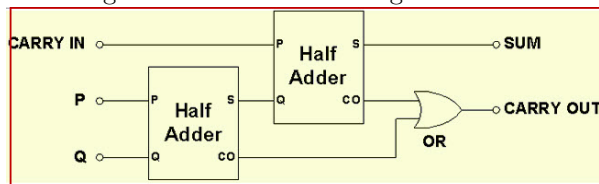


Figure 4.3: Full adder using half adders



Chapter 5

Algorithms

5.1 Addition

1. First, convert the two representations to scientific notation. Thus, we explicitly represent the hidden 1.
2. In order to add, we need the exponents of the two numbers to be the same. We do this by rewriting Y. This will result in Y being not normalized, but value is equivalent to the normalized Y.
3. Add $x - y$ to Y's exponent. Shift the radix point of the mantissa Y left by $x - y$ to compensate for the change in exponent.
4. Add the two mantissas of X and the adjusted Y together.
5. If the sum in the previous step does not have a single bit of value 1, left of the radix point, then adjust the radix point and exponent until it does.
6. Convert back to the normalized floating point representation.

5.2 Substraction

1. First, convert the two representations to scientific notation. Thus, we explicitly represent the hidden 1.
2. In order to add, we need the exponents of the two numbers to be the same. We do this by rewriting Y. This will result in Y being not normalized, but value is equivalent to the normalized Y.
3. Add $x - y$ to Y's exponent. Shift the radix point of the mantissa Y left by $x - y$ to compensate for the change in exponent.
4. Subtracts the two mantissas of X and the adjusted Y together.
5. If the difference in the previous step does not have a single bit of value 1, left of the radix point, then adjust the radix point and exponent until it does.
6. Convert back to the normalized floating point representation.

5.3 Multiplication

1. First, convert the two representations to scientific notation. Thus, we explicitly represent the hidden 1.
2. Let x be the exponent of X . Let y be the exponent of Y . The resulting exponent (call it z) is the sum of the two exponents. z may need to be adjusted after the next step.
3. Multiply the mantissa of X to the mantissa of Y . Call this result m .
4. If m does not have a single 1 left of the radix point, then adjust the radix point so it does, and adjust the exponent z to compensate.
5. Add the sign bits, mod 2, to get the sign of the resulting multiplication.
6. Convert back to the one byte floating point representation, truncating bits if needed.

Chapter 6

Working Example

Figure 6.1: IEEE 754 - Arithmetic Operation

```
*Main> (convert_to_ieee 123.45)
(0,[1,0,0,0,0,1,0,1],[1,1,1,0,1,1,0,1,1,1,0,0,1,1,0,0,1,1,0,0,1,1,0])
*Main>
*Main> (convert_to_ieee 67.89)
(0,[1,0,0,0,0,1,0,1],[0,0,0,0,1,1,1,1,1,0,0,0,1,1,1,1,0,1,0,1,1,1,0])
*Main>
*Main> addition (convert_to_ieee 12.345) (convert_to_ieee 67.89)
(0,[1,0,0,0,0,1,0,1],[0,1,0,0,0,0,0,0,1,1,1,1,0,0,0,0,1,0,1,0,0,0,1])
*Main>
*Main> convert_from_ieee $ addition (convert_to_ieee 12.345) (convert_to_ieee 67.89)
80.23499298095703
*Main>
*Main> convert_from_ieee $ subtract (convert_to_ieee 12.345) (convert_to_ieee 67.89)
-55.545005798339844
*Main>
*Main> convert_from_ieee $ multiply (convert_to_ieee 12.345) (convert_to_ieee 67.89)
838.1019897460938
*Main> █
```

Chapter 7

Conclusion

By using the algorithm of various operations on the input parameters , we can get the results that are similar to actual results of the various operations. So the application is able to do arithmetic operations on the inputs operands as per IEEE standards for floating points operations.

Bibliography

- [1] 754-2008 - ieee standard for floating-point arithmetic.