



Performance Analysis of Floating Point Arithmetic

P.Pushpalatha¹, Pinapothu Srilakshmi²

Assistant professor, Department of Electronics and Communication Engineering, UCEK, JNTUK
Kakinada, Andhra Pradesh, India¹

P.G. Student, Department of Electronics and Communication Engineering, UCEK, JNTUK
Kakinada, Andhra Pradesh, India²

pushpalatha86@gmail.com¹, Srilakshmi.hyma@gmail.com²

ABSTRACT: This paper mainly presents performance analysis of floating point arithmetic. Floating point arithmetic performs addition, subtraction, multiplication and division. This paper briefly explains about the floating addition, subtraction and multiplication operations. These operations are used in numerical calculations, FFT and butterfly operations. Floating point numbers are one possible way of representing real numbers in binary format. The arithmetic units in modern microprocessors execute advanced applications such as 3D graphics, multimedia, signal processing, and a variety of scientific computations that require complex mathematic computations. The fixed-point number system is not sufficient to handle such complex computations. In contrast, the floating-point notation, which is specified in the IEEE-754 Standard for floating-point arithmetic, represents a wide range of numbers from tiny fractions to extremely large numbers. Floating point multiplier is simple in terms of overall structure it requires more logic area and delay compared to the floating-point adder and floating point subtraction.

KEYWORDS: Floating point, Floating point multiplication, Floating point arithmetic.

I. INTRODUCTION

Floating point numbers are one possible way of representing real numbers in binary format. The IEEE 754 standard presents two different floating point formats, binary interchange format and decimal interchange format. Multiplying floating point numbers is a critical requirement for DSP applications involving large dynamic range. This paper focuses on the floating point arithmetic in single precision. Fig. 1 shows the IEEE 754 single precision binary format representation. It consists of a one bit sign (S), an eight bit exponent (E), and a twenty three bit fraction (M or Mantissa). An extra bit is added to the fraction to form what is called the significand. If the exponent is greater than 0 and smaller than 255, and there is 1 in the MSB of the significand then the number is said to be a normalized number. The arithmetic units in modern microprocessors execute advanced applications such

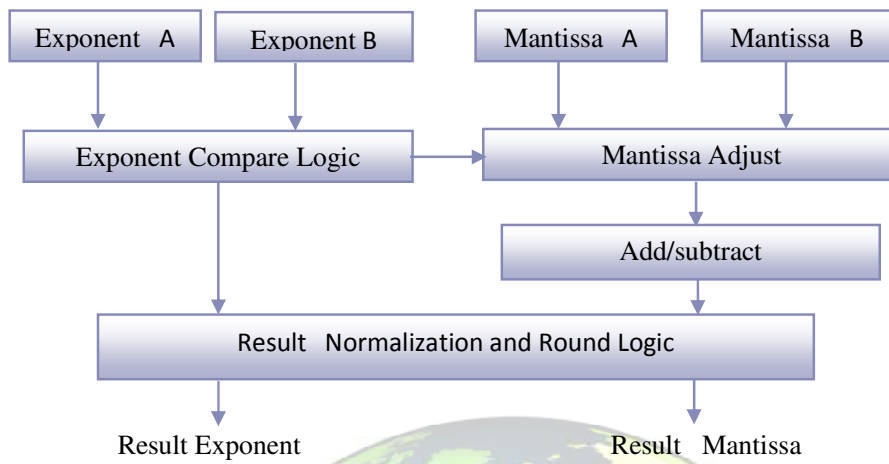


Fig 2: Addition/Subtraction of floating point numbers.

IV. FLOATING POINT SUBTRACTION

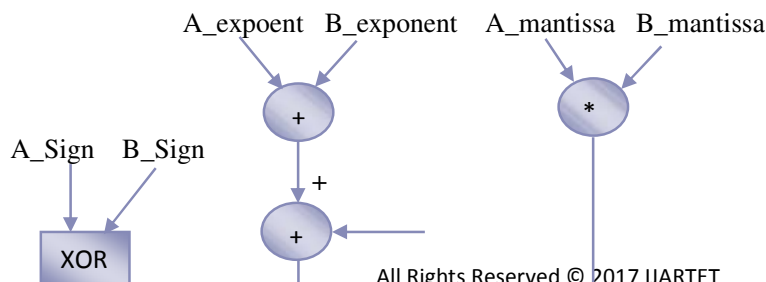
A simple method to subtract floating point numbers is to first represent them with the same exponent. Subtraction of two numbers in floating point format is done by

1. Compare the exponents of the two numbers and shift smaller number to right until its exponent would match the larger exponent.
2. Subtraction the significands.
3. Normalization the sum, either shifting right and incrementing exponent or shifting left and decrementing the exponent.
4. Check the result overflow or underflow.
5. Round the significand to the appropriate number of bits and finally obtain floating point subtraction result.

V. FLOATING POINT MULTIPLICATION

The floating-point multiplier takes two input operands and produces a rounded product result. Although the floating-point multiplier is simple in terms of overall structure, it requires more logic area and power consumption compared to the floating-point adder.

Multiplying two numbers in floating point format is done by adding the exponent of the two numbers then subtracting the bias from their result, multiplying the significand of the two numbers and calculating the sign by xoring the sign of the two numbers. In order to represent the multiplication result as a normalized number there should be 1 in the MSB of the result.





- Bias

Fig 3: floating point multiplication

To multiply two floating point numbers the following is done:

1. Multiplying the significand ($1.M1 * 1.M2$).
2. Placing the decimal point in the result.
3. Adding the exponents ($E1 + E2 - Bias$).
4. Obtaining the sign $s1 \text{ xor } s2$.
5. Normalizing the result obtaining 1 at the MSB of the results significand.
6. Rounding the result and finally obtain floating point multiplier result.

VI. RESULT COMPARSION TABLE

S.NO	TYPE OF OPERATION	AREA	DELAY
1.	ADDITION	489 slices	23.690 ns
2.	SUBTRACTION	509 slices	22.996 ns
3.	MULTIPLING	694 slices	31.07 ns

VII. SIMULATION RESULT

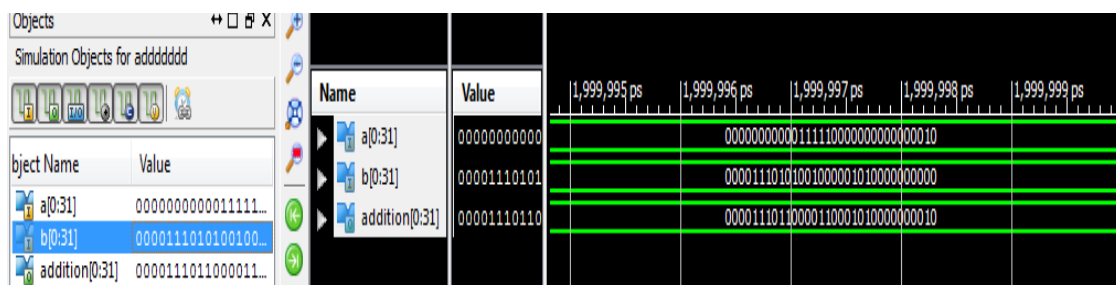




Fig 4: floating point addition

Project Name	Value	Name	Value	1,999,995 ps	1,999,996 ps	1,999,997 ps	1,999,998 ps	1,999,999 ps
a[31:0]	000011010101011000	a[31:0]	00001101010		000011010101010000000000000000			
b[31:0]	0000111111101010010	b[31:0]	000011111110		000011111110100101000000000000			
mul[31:0]	110000000000000000	mul[31:0]	110000000000		110000000000000000000000000000			

Fig 5: floating point multiplication

VI. CONCLUSION

Floating arithmetic operations are most commonly used mathematical operations in numerical calculations. In this mainly discuss about the performance analysis of the floating point addition, floating point subtraction and floating point multiplication. Floating point multiplication requires more area and delay when compare to addition and subtraction.

REFERENCES

- [1]. IEEE Standard for floating point arithmetic, IEEE Standard 754-2008, New York, Inc., Aug.29, 2008.
- [2]. IEEE Standard 754 Floating Point Numbers, by Steve Hollasch.
- [3]. Floating point arithmetic unit using verilog, by lalitha gongavar and rajan chaudhary.
- [4]. Christo Ananth, Muthamil Jothi.M, M.Priya, V.Manjula, "Parallel RC4 Key Searching System Based on FPGA", International Journal of Advanced Research in Management, Architecture, Technology and Engineering (IJARMATE), Volume 2, Special Issue 13, March 2016, pp: 5-12
- [5]. 32 bit Single Precision floating point Multiplier, Ms.Radhika Jumde, AVBIT, Pawnar, Wardha.
- [6]. Normalization on floating point Multiplication using Verilog HDL, V.Narasimha, V. Swathi.
- [7]. Floating Point Adder and Multiplier, Eduardo Sanchez EPFL- HEIG- VD. Aug-2013.