

# Review Paper on Matrix Multiplication using IEEE 754 Floating Point and Different Types of Adders

Manjusha Kumari<sup>1</sup>, Vijay Yadav<sup>2</sup>

<sup>1</sup>M. Tech. Scholar, <sup>2</sup>Assistant Professor

Department of Electronics and Communication, LNCT, Bhopal

**Abstract**-Due to advancement of new technology in the field of VLSI and Embedded system, there is an increasing demand of high speed and low power consumption processor. Speed of processor greatly depends on its multiplier as well as adder performance. In spite of complexity involved in floating point arithmetic, its implementation is increasing day by day. Due to which high speed adder architecture become important. Several adder architecture designs have been developed to increase the efficiency of the adder. In this paper, we introduce an architecture that studied of matrix multiplication using IEEE 754 floating point multiplier and different types of adder.

**Keywords**-IEEE 754, Single Precision Floating Point (SPFP), Double Precision Floating Point (DPFP)

## I. INTRODUCTION

Matrix multiplication is commonly used in most signal processing algorithms. It is also a frequently used kernel operation in a wide variety of graphics, image processing as well as robotic applications. The matrix multiplication operation involves a large number of multiplication as well as accumulation. Multipliers have large area, longer latency and consume considerable power compared to adders. Registers, which are required to store the intermediate product values, are also major power intensive component [1]. These components pose a major challenge for designing VLSI structures for large-order matrix multipliers with optimized speed and chip-area. However, area, speed and power are usually conflicting hardware constraints such that improving upon one factor degrades the other two [2].

The real numbers represented in binary format are known as floating point numbers. Based on IEEE-754 standard, floating point formats are classified into binary and decimal interchange formats. Floating point multipliers are very important in dsp applications. This paper focuses on double precision normalized binary interchange format. Figure 1 shows the IEEE-754 double precision binary format representation. Sign (s) is represented with one bit, exponent (e) and fraction (m or mantissa) are represented with eleven and fifty two bits respectively. For a number is said to be a normalized number, it must consist of 'one' in the MSB of the significant and exponent is greater than zero and smaller than 1023. The real number is represented by equations (i) & (2).

$$Z = (-1^s) \times 2^{(E-Bias)} \times (1.M) \dots\dots\dots(1)$$

$$Value = (-1^{signbit}) \times 2^{(Exponent-1023)} \times (1.Mantissa) \dots\dots\dots(2)$$

Biasing makes the values of exponents within an unsigned range suitable for high speed comparison.

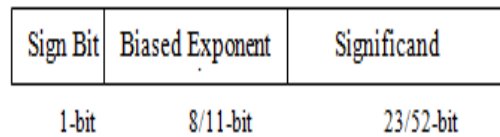
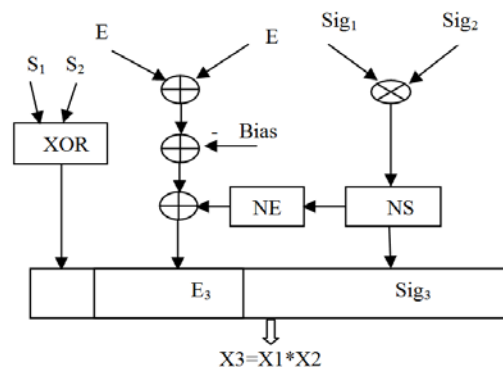


Figure 1 IEEE 754 Single precision and double precision floating point format.

- IEEE 754 standard floating point multiplication algorithm

A brief overview of floating point multiplication has been explained below [3-4].

- Both sign bits S1, S2 are need to be Xoring together, then the result will be sign bit of the final product.
- Both the exponent bits E1, E2 are added together, then subtract bias value from it. So, we get exponent field of the final product.



- Significant bits Sig1 and Sig2 of both the operands are multiply including their hidden bits.
- Normalize the product found in step 3 and change the exponent accordingly. After normalization, the leading "1" will become the hidden bit.

Above algorithm of multiplication algorithm is shown in Figure 2.

## II. LITERATURE REVIEW

Lakshmi kiran Mukkara et al. [1], for implementation of Low Power VLSI Architectures in the area of Digital Image Processing applications, Matrix Multiplication is a key arithmetic operation. To construct VLSI architectures with Low Power, High Speed and Low area, Matrix Multiplication design becomes complex. In this paper, a simple novel VLSI architecture for FPMM is presented. It is designed by considering Pseudo code for matrix multiplication, CSD multiplication algorithm for power reduction, Conventional floating point number format and Pipelining concept for improving speed. FPMM design is applicable for any arbitrary size of matrices by following matrix rules.

Soumya Havaldar et al. [2], gives an FPGA Based High Speed IEEE-754 Double Precision Floating Point Multiplier Using Verilog. This paper has implemented DFPF Multiplier using parallel Adder. A high speed floating point double precision multiplier is implemented on a Virtex-6 FPGA. In addition, the proposed design is compliant with IEEE-754 format and handles over flow, under flow, rounding and various exception conditions. The design achieved the operating frequency of 414.714 MHz with an area of 648 slices.

Ragini Parte et al. [3], IEEE point number-crunching has an immense application in DSP, advanced PCs, robots because of its capacity to speak to little numbers and huge numbers and in addition marked numbers and unsigned numbers. Disregarding unpredictability included in gliding point number juggling, its usage is expanding step by step. Here we break down the impacts of utilizing three unique sorts of adders while figuring the single accuracy and twofold exactness skimming point increase. We additionally exhibit the increase of significant bits by disintegration of operands strategy for IEEE 754 standard.

Ross Thompson et al. [4], IEEE-754 determines trade and number juggling positions also, routines for paired and decimal drifting point number juggling in PC programming world. The execution of a skimming point framework utilizing this standard should possible completely in programming, or in equipment, or in any blend of programming and equipment. This venture propose VHDL execution of IEEE - 754 Floating point unit .In proposed work the pack, unload and adjusting mode was actualized utilizing the VHDL dialect and reenactment was checked.

In this proposition work, DFPF Multiplier alongside SPFP Multiplier has been actualized with four ordinary Adders (PA, CSKA, CSA, and CSABEC). Think about their Results and CSA is known not the speediest snake among

every single customary viper. Be that as it may, CSA involves more territory as it has two parallel circuits to include the same bits yet with diverse convey data. As CSA figures the whole without sitting tight for the transitional conveys to spread stage by stage. Finally it is the obligation of multiplexer to pick and give the last right yield. CSABEC is an adjusted adaptation of CSA in which one of the parallel circuits is supplanted by the arrangement of Binary to Excess-1 Converters circuit (BECs). It is turned out to be an awesome way to deal with decrease the territory.

M. K. Jaiswal et al. [5], we display a nonintrusive simultaneous mistake identification (CED) technique for ensuring the control rationale of a contemporary coasting point unit (FPU). The proposed strategy depends on the perception that control rationale mistakes lead to broad information way defilement and influence, with high likelihood, the example some portion of the IEEE-754 coasting point representation. Consequently, type observing can be used to recognize blunders in the control rationale of the FPU. Anticipating the type includes generally basic operations; in this way, our system causes fundamentally bring down overhead than the traditional methodology of copying the control rationale of the FPU. In fact, trial results on the open SPARC T1 processor utilizing SPEC2006FP benchmarks demonstrate that when contrasted with control rationale duplication, which brings about a zone overhead of 17.9 percent of the FPU size, our technique acquires a region overhead of just 5.8 percent yet still accomplishes discovery of more than 93 percent of transient mistakes in the FPU control rationale. Besides, the proposed technique offers the subordinate advantage of additionally distinguishing 98.1 percent of the information way mistakes that influence the type, which can't be recognized by means of duplication of control rationale. At long last, when joined with a traditional deposit code-based strategy for the portion, our system prompts a complete CED answer for the whole FPU which gives scope of 94.1 percent of all blunders at a region expense of 16.32 percent of the FPU size.

Table 1: Summary of Literature Review

Entitle of paper	Approach ed used	Softwar e	Paramet er	Publish ed Year
Design of Vedic IEEE 754 Floating Point Multiplier	IEEE 754 floating point multiplier using Vedic Multiplier	Xilinx 12.1 i	Slice = 305, LUT= 598	IEEE 2016
Analysis of Effects of using	IEEE 754 floating point	Verilog using the	Slice = 353, LUT=	IEEE 2015

Exponent Adders in IEEE-754 Multiplier by VHDL	multiplier using Carry select adder	Model SIM SE 6.3	610	
An IEEE 754 Double-Precision Floating-Point Multiplier for Denormalized and Normalized Floating-Point Numbers	IEEE 754 floating point multiplier using floating number	Xilinx Spartan 3E xa3s250 e-4vqg100	Slice = 453, LUT= 678	IEEE 2015
High Performance FPGA Implementation of Double Precision Floating Point Adder/Subtractor	IEEE 754 floating point multiplier using adder and subtractor	Xilinx ISE 12.2	Slice = 543, LUT= 792	IEEE 2011

III. DIFFERENT TYPES OF ADDER

Parallel Adder:-

Parallel adder can add all bits in parallel manner i.e. simultaneously hence increased the addition speed. In this adder multiple full adders are used to add the two corresponding bits of two binary numbers and carry bit of the previous adder. It produces sum bits and carry bit for the next stage adder. In this adder multiple carry produced by multiple adders are rippled, i.e. carry bit produced from an adder works as one of the input for the adder in its succeeding stage. Hence sometimes it is also known as Ripple Carry Adder (RCA). Generalized diagram of parallel adder is shown in figure 3.

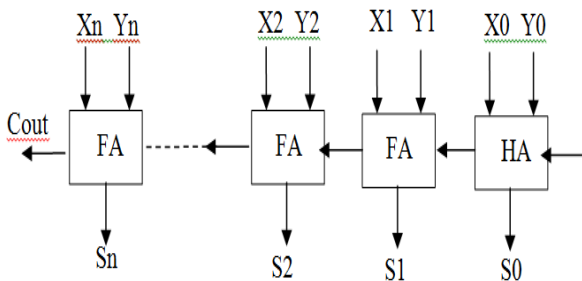


Figure 3: Parallel Adder (n=7 for SPFP and n=10 for DPFP)

An n-bit parallel adder has one half adder and n-1 full adders if the last carry bit required. But in 754 multiplier's exponent adder, last carry out does not required so we can

use XOR Gate instead of using the last full adder. It not only reduces the area occupied by the circuit but also reduces the delay involved in calculation. For SPFP and DPFP multiplier's exponent adder, here we Simulate 8 bit and 11 bit parallel adders respectively as show in figure 4.

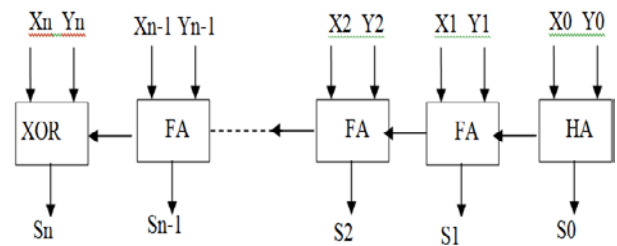


Figure 4: Modified Parallel Adder (n=7 for SPFP and n=10 for DPFP)

Carry Skip Adder:-

This adder gives the advantage of less delay over Ripple carry adder. It uses the logic of carry skip, i.e. any desired carry can skip any number of adder stages. Here carry skip logic circuitry uses two gates namely "and gate" and "or gate". Due to this fact that carry need not to ripple through each stage. It gives improved delay parameter. It is also known as Carry bypass adder. Generalized figure of Carry Skip Adder is shown in figure 5.

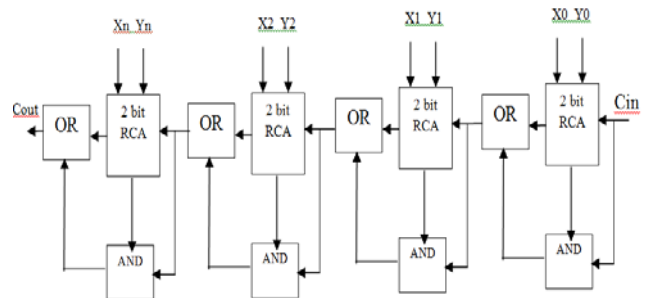


Figure 5: Carry Skip Adder

Carry Select Adder:-

Carry select adder uses multiplexer along with RCAs in which the carry is used as a select input to choose the correct output sum bits as well as carry bit. Due to this, it is called Carry select adder. In this adder two RCAs are used to calculate the sum bits simultaneously for the same bits assuming two different carry inputs i.e. '1' and '0'. It is the responsibility of multiplexer to choose correct output bits out of the two, once the correct carry input is known to it. Multiplexer delay is included in this adder. Generalized figure of Carry select adder is shown in figure 3.9. Adders are the basic building blocks of most of the ALUs (Arithmetic logic units) used in Digital signal processing and various other applications. Many types of adders are available in today's scenario and many more are developing day by day. Half adder and Full adder are the two basic types of adders. Almost all other adders are made with the different arrangements of these two basic

adders only. Half adder is used to add two bits and produce sum and carry bits whereas full adder can add three bits simultaneously and produces sum and carry bits.

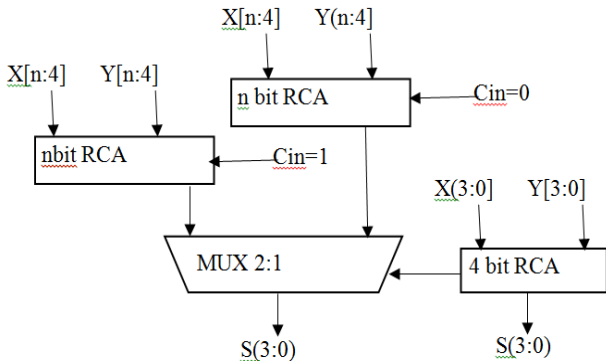


Figure 6: Carry Select Adder

IV. PROPOSED DESIGN

In IEEE754 standard floating point representation, 8 bit Exponent field in single precision floating point (SP FP) representation and 11 bit in double precision floating point (DP FP) representation are need to add with another 8 bit exponent and 11 bit exponent respectively, in order to multiply floating point numbers represented in IEEE 754 standard as explained earlier. Ragini et al. [10] has used parallel adder for adding exponent bits in floating point multiplication algorithm. We proposed the use of 8-bit modified CSA with dual RCA and 8-bit modified CSA with RCA and BEC for adding the exponent bits. We have found the improved area of 8-bit modified Carry select adder with RCA and BEC over the 8-bit modified CSA with dual RCA.

- Sign bit calculation

To calculate the sign bit of the resultant product for SP FP and DP FP multiplier, the same strategy will work. We just need to XOR together the sign bits of both the operands. If the resultant bit is '1', then the final product will be a negative number. If the resultant bit is '0', then the final product will be a positive number.

- Exponent bit calculation

Add the exponent bits of both the operands together, and then the bias value (127 for SPFP and 1023 for DPFP) is subtracted from the result of addition. This result may not be the exponent bits of the final product. After the significant multiplication, normalization has to be done for it. According to the normalized value, exponents need to be adjusted. The adjusted exponent will be the exponent bits of the final product.

- Significant bit calculation

Significant bits including the one hidden bit are need to be multiply, but the problem is the length of the operands. Number of bits of the operand will become 24 bits in case of SP FP representation and it will be 53 bits in case of DP

FP representation, which will result the 48 bits and 106 bits product value respectively. In this paper we use the technique of break up the operands into different groups then multiply them. We get many product terms, add them together carefully by shifting them according to which part of one operand is multiplied by which part of the other operand. We have decomposed the significant bits of both the operands in four groups. Multiply each group of one operand by each group of second operand.

Matrix Multiplication

In this design we have reduced the resource utilization in terms of number of multipliers and registers in lieu of the completion time.

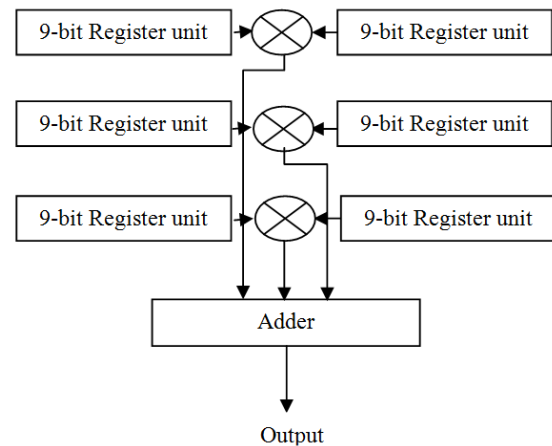


Figure 7 Proposed PPI-So design for n=3

This design is particularly useful where resources are limited and design can be compromised on basis of increased completion time. The basic working model for a 3 x 3 matrix-matrix multiplication is shown in figure 7 below.

Considering the matrix – matrix multiplication of two nxn matrices, the calculation is performed using n number of multipliers, n number of registers and n-1 number of adders. cycles are required to perform the matrix multiplication operation. Each multiplier has two input ports: one each from matrix A and B. In each cycle, n numbers of multiplications are performed and the products are fed to the adder block to give a single element of the output matrix, C. The data flow to the multipliers are such that, multiplier is fed from column of matrix A and row of matrix B, where 1 < k < n. At the multiplier, each element from matrix A is repeated for n consecutive cycles whereas the elements from matrix B are cycled back after n cycles. The partial products are then fed to the adder which computes the final result.

For a better understanding of the process, let us consider the matrix multiplication for n = 3 (as shown in figure 1). In this case, 3 multipliers and 3 registers are used to calculate and store the partial products respectively. These partial products are then fed to the adder block to compute

the final result. The first multiplier receives input from the first column of matrix A ( $a_{k1}$ ) and first row of matrix B ( $b_{1k}$ ), where. Each element of the matrix A at the first multiplier is repeated for 3 cycles, such that the data flow can be represented as . Similarly, at the first multiplier, the elements of B are repeated after 3 cycles, such that the input data-flow will be . The other two multipliers receive the component of A and B in the similar order as the first multiplier. After the multiplication, the partial products are fed to the adder which computes the elements of output matrix C in row major order given by . So the entire matrix multiplication operation is performed in  $=9$  cycles.

## V. CONCLUSION

IEEE754 standardize two basic formats for representing floating point numbers namely, single precision floating point and double precision floating point. Floating point arithmetic has vast applications in many areas like robotics and DSP. Delay provided and area required by hardware are the two key factors which are need to be consider Here we present single precision floating point multiplier by using two different adders namely modified CSA with dual RCA and modified CSA with RCA and BEC.

Among all two adders, modified CSA with RCA and BEC is the least amount of Maximum combinational path delay (MCDP). Also, it takes least number of slices i.e. occupy least area among all two adders.

## REFERENCES

- [1] Lakshmi kiran Mukkara and K.Venkata Ramanaiah, "A Simple Novel Floating Point Matrix Multiplier VLSI Architecture for Digital Image Compression Applications", 2nd International Conference on Inventive Communication and Computational Technologies (ICICCT 2018) IEEE.
- [2] Soumya Havaladar, K S Gurumurthy, "Design of Vedic IEEE 754 Floating Point Multiplier", IEEE International Conference On Recent Trends In Electronics Information Communication Technology, May 20-21, 2016, India.
- [3] Ragini Parte and Jitendra Jain, "Analysis of Effects of using Exponent Adders in IEEE- 754 Multiplier by VHDL", 2015 International Conference on Circuit, Power and Computing Technologies [ICCPCT] 978-1-4799-7074-2/15/\$31.00 ©2015 IEEE.
- [4] Ross Thompson and James E. Stine, "An IEEE 754 Double-Precision Floating-Point Multiplier for Denormalized and Normalized Floating-Point Numbers", International conference on IEEE 2015.
- [5] M. K. Jaiswal and R. C. C. Cheung, "High Performance FPGA Implementation of Double Precision Floating Point Adder/Subtractor", in International Journal of Hybrid Information Technology, vol. 4, no. 4, (2011) October.
- [6] B. Fagin and C. Renard, "Field Programmable Gate Arrays and Floating Point Arithmetic," IEEE Transactions on VLSI, vol. 2, no. 3, pp. 365-367, 1994.
- [7] N. Shirazi, A. Walters, and P. Athanas, "Quantitative Analysis of Floating Point Arithmetic on FPGA Based Custom Computing Machines," Proceedings of the IEEE Symposium on FPGAs for Custom Computing Machines (FCCM'95), pp.155-162, 1995.
- [8] Malik and S. -B. Ko, "A Study on the Floating-Point Adder in FPGAs", in Canadian Conference on Electrical and Computer Engineering (CCECE-06), (2006) May, pp. 86-89.
- [9] D. Sangwan and M. K. Yadav, "Design and Implementation of Adder/Subtractor and Multiplication Units for Floating-Point Arithmetic", in International Journal of Electronics Engineering, (2010), pp. 197-203.
- [10] L. Louca, T. A. Cook and W. H. Johnson, "Implementation of IEEE Single Precision Floating Point Addition and Multiplication on FPGAs", Proceedings of 83rd IEEE Symposium on FPGAs for Custom Computing Machines (FCCM'96), (1996), pp. 107-116.
- [11] Jaenicke and W. Luk, "Parameterized Floating-Point Arithmetic on FPGAs", Proc. of IEEE ICASSP, vol. 2, (2001), pp. 897-900.
- [12] Lee and N. Burgess, "Parameterisable Floating-point Operations on FPGA", Conference Record of the Thirty-Sixth Asilomar Conference on Signals, Systems, and Computers, (2002).
- [13] M. Al-Ashrafy, A. Salem, W. Anis, "An Efficient Implementation of Floating Point Multiplier", Saudi International Electronics, Communications and Photonics Conference (SIECP), (2011) April 24-26, pp. 1-5.