

# Study of Single and Double Precision Floating Point Multiplier using Different Types of Adder

M.Tech. Scholar Marode Shubham Hemant, Asst. Prof. Dr. Neetesh Raghuwanshi

Department of Electronics & Communication Engineering  
 Sarvepalli Radhakrishnan University, Bhopal

**Abstract-** In high-performance computing applications such as signal processing and image processing, floating point multiplication is a critical and widely utilized arithmetic operation. The IEEE-754 standard floating point multiplier's high speed and memory requirements make its implementation more difficult in many applications that require quick processing. As a result, developing an efficient architecture to accomplish the floating point multiplication operation in less time and with less memory has become a major research focus in practically every discipline of electronics. This study review paper provides a detailed explanation of each of the proposed architectures for floating point multipliers, as well as their benefits and drawbacks. In this paper, we introduce an architecture that performs high speed IEEE 754 floating point multiplier using carry select adder (CSA). Here we are introduced two carry select based design. These designs are implementation Xilinx Vertex device family.

**Keywords-** IEEE754, Single Precision Floating Point (SP FP), Double Precision Floating Point (DP FP), Binary to Excess-1 ().

## I. INTRODUCTION

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. A floating-point unit (FPU) colloquially is a math coprocessor, which is a part of a computer system specially designed to carry out operations on floating point numbers. Typical operations that are handled by FPU are addition, subtraction, multiplication and division.

The aim was to build an efficient FPU that performs basic as well as transcendental functions with reduced complexity of the logic used reduced 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).

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

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

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

Sign Bit	Biased	Significand
1-bit	8/11-bit	23/52-bit

Fig 1. IEEE 754 Single Precision and Double Precision Floating Point Format.

## 1. IEEE 754 standard floating point multiplication algorithm:

A brief overview of floating point multiplication has been explained below [5-6].

- [1] Both sign bits  $S_1, S_2$  are need to be Xoring together, then the result will be sign bit of the final product.
- [2] Both the exponent bits  $E_1, E_2$  are added together, and then subtract bias value from it. So, we get exponent field of the final product.
- [3] Significand bits  $Sig_1$  and  $Sig_2$  of both the operands are multiply including their hidden bits.
- [4] 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.

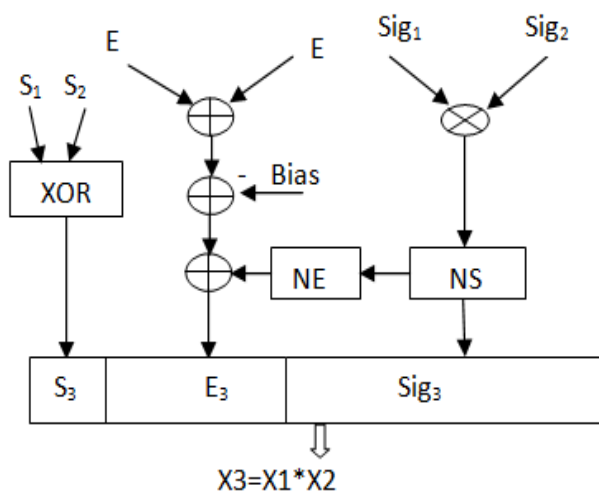


Fig 2. IEEE754 SP FP and DP FP Multiplier Structure,  
NE: Normalized exponent, NS: Normalized  
Significand

## 2. IEEE 754 Standards:

IEEE754 standard is a technical standard established by IEEE and the most widely used standard for floating-point computation, followed by many hardware (CPU and FPU) and software implementations [3]. Single-precision floating-point format is a computer number format that occupies 32 bits in a computer memory and represents a wide dynamic range of values by using a floating point. In IEEE 754-2008, the 32-bit with base 2 format is officially referred to as single precision or binary32. It was called single in IEEE 754-1985.

The IEEE 754 standard specifies a single precision number as having sign bit which is of 1 bit length, an exponent of width 8 bits and a significant precision of 24 bits out of which 23 bits are explicitly stored and 1 bit is implicit 1.

## II. LITERATURE REVIEW

### 1. Open Floating Point Unit:

This was the open source project done by Rudolf Usselmann. FPU described a single precision floating point unit which could perform add, subtract, multiply, divide, and conversion between FP number and integer. It consists of two pre-normalization units that can adjust the mantissa as well as the exponents of the given numbers. One unit is for addition and subtraction operation and the other one is for multiplication and division operations.

It also has different units for different operations that perform an actual addition subtraction, multiplication and division. It also has a shared post normalization unit that normalizes the fraction part. The final result after post-normalization is directed to a valid result which is in accordance to single precision FP format. The main drawback of this model was that most of the codes were written in MATLAB and due to this it is non-synthesizable.

### 2. GRFPU:

This high Performance IEEE754 FPU was designed at Gaisler Research for the improvement of FP operations of a LEON based systems. It supports both single precision and double precision operands. It implements all floating point arithmetic operations defined by the IEEE754 standard in hardware. All operations are dealt with the exception of denormalized numbers which are flushed to zero and supports all rounding modes. This advanced design combines low latency and high throughput.

The most common operations such as addition, subtraction and multiplication are fully pipelined which has throughput of one CC and a latency of three CC. More complex divide and square root operation takes between 1 to 24 CC to complete and execute in parallel with other FP operations. It can also perform operations like converse and compliment. It supports all SPARC V8 FP instructions. The main drawback of this model is that it is very expensive and complex to implement practically. In this section some authors papers are also discuss.

**Ms. Meenu et al.** computerized Signal handling turned into an application to make rapid information handling frameworks like 3D delivering, 4G portable web, and so forth, we really want best processors with elite execution information way units and there

is a developing requirement for research on elective techniques for signal handling equipment execution. In most frameworks utilizing computerized signal handling Multiply-Accumulate (MAC) is one of the fundamental capacities. The execution of the entire framework relies upon the exhibition of the MAC units set up.

**Paldurai et al.** augmentation of drifting point numbers observed broad use in DSP applications including tremendous reach. The basic part in drifting point augmentation is the increase of mantissas which involves  $24 \times 24$  piece whole number multiplier for single accuracy drifting point numbers. The speed of the framework can be upgraded by working on the speed of duplication. In this paper a 24 cycle Vedic multiplier has been proposed involving  $3 \times 3$  Vedic multiplier as its essential square.

The proposed and regular drifting point multipliers dependent on Vedic science are coded in Verilog, Synthesized and mimicked in ISE Simulator. Greatest combinational way postponement and number of cuts needed on FPGA are looked at for proposed and customary multipliers. The outcomes obviously show that proposed technique incredibly affect working on the speed and diminish the region needed on Spartan 6 FPGA.

**Ragini Parte et al.** drifting point number-crunching has a tremendous applications in DSP, computerized PCs, robots because of its capacity to address tiny numbers and enormous numbers just as marked numbers and unsigned numbers. Notwithstanding intricacy engaged with drifting point number juggling, its execution is expanding step by step. Here we examine the impacts of utilizing three distinct kinds of adders while computing the single accuracy and twofold accuracy drifting point increase. We likewise present the duplication of significand bits by disintegration of operands strategy for IEEE 754 norm.

**Irine Padma et al.** to address exceptionally enormous or little qualities, huge reach is needed as the number portrayal is no more suitable. These qualities can be addressed utilizing the IEEE 754 standard based drifting point portrayal. Duplicating drifting point numbers is a basic prerequisite for DSP applications including huge unique reach. The paper depicts the execution and plan of IEEE 754 Pipelined Floating Point Multiplier dependent on Vedic

Multiplication Technique. The contributions to the multiplier are given in IEEE 754, 32 bit design. The Urdhva Triyakbhyam sutra is utilized for the augmentation of mantissa. The sub-current and flood cases are taken care of.

**R. Sai Siva et al.** in this paper we portray an effective execution of an IEEE 754 single accuracy drifting point multiplier utilizing vedic math. The motivation behind utilizing vedic math is because of expansion in the quantity of fractional items in typical augmentation process, with utilizing vedic math incomplete items can be decreased so the region also power requirements of the drifting point multiplier can be decreased proficiently.

**Priyanka Koneru et al.** a quick and energy effective drifting point unit is constantly required in significant applications like computerized signal handling, picture handling, and ongoing information handling and media applications. As circuits get shrivel, the coordinated plan turns into a basic test as far as clock slant and clock dispersion. One alluring option is to utilize powerful offbeat circuits, which effortlessly oblige these planning disparities. In this paper, a solitary accuracy nonconcurrent drifting point multiplier is carried out utilizing VERILOG equipment depiction language.

**I. V. Vaibhav et al.** IEEE drifting point design was a standard arrangement utilized in all handling parts since Binary floating point numbers expansion is one of the principal limits used in cutting edge sign dealing with (DSP) application. In that work VHDL execution of Floating Point Multiplier utilizing old Vedic science is introduced. The thought for planning the multiplier unit is taken on from antiquated Indian science "Vedas". The Urdhvatriyakbhyam sutra will be utilized for the augmentation of Mantissa. The sub-current and over stream cases will be dealt with. The contributions to the multiplier in 32 digit design. The multiplier is planned in VHDL or VERILOG and reproduced utilizing Modelsim.

**D. Monniaux et al.** because of fast development in monetary, business, and Internet-based applications, there is an expanding want to permit PCs to work on both paired and decimal drifting point numbers. Thus, details for decimal drifting point support are being added to the IEEE-754 Standard for Floating-Point Arithmetic. In this paper, we present the plan and execution of a decimal drifting point viper that is

consistent with the current draft modification of this norm. The viper upholds procedure on 64-bit (16-digit) decimal drifting point operands. We give union outcomes showing the assessed region and deferral for our plan when it is pipelined to different profundities.

**Soumya Havaladar et al.** drifting point number can co-occurently foster a noticeable scope of numbers and an undeniable degree of accuracy. Duplication of drifting point numbers tracked down broad use in more extensive scope of innovative and business computations It is expected to execute quicker multipliers including restricted region and devouring decreased power. This paper proposes a drifting point multiplier which oversees flood, sub-current and adjusting. The proposed and traditional drifting point multipliers dependent on Vedic arithmetic would be coded in Verilog, Synthesized and Simulated utilizing ISE Simulator. Xilinx Virtex VI FPGA will be utilized for Hardware acknowledgment and Verification. It is proposed to analyze asset use and timing execution of the proposed multiplier with that of existing at this point.

### III. DIFFERENT TYPES OF ADDER

#### 1. 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.

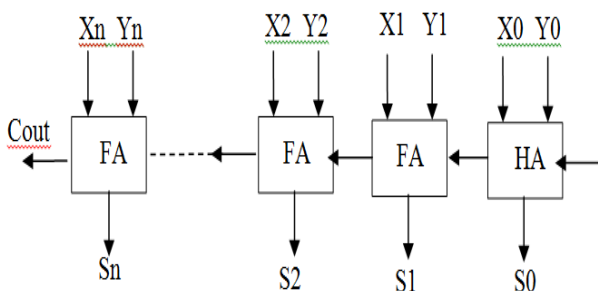


Fig 3. Parallel Adder (n=7 for SPFP and n=10 for DPFP)

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.

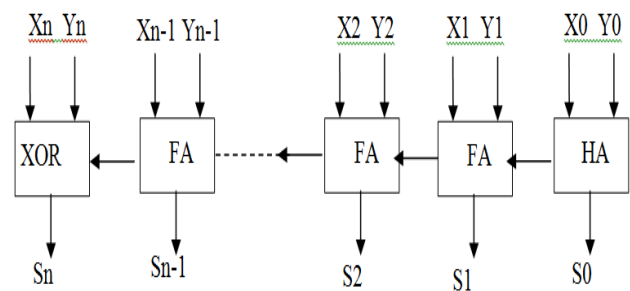


Fig 4. Modified Parallel Adder (n=7 for SPFP and n=10 for DPFP)

#### 2. 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.

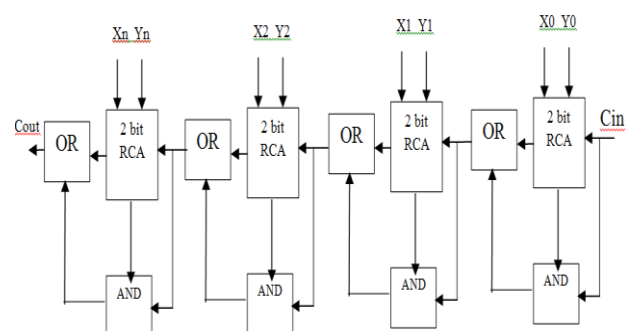


Fig 5. Carry Skip Adder.

#### 3. 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.

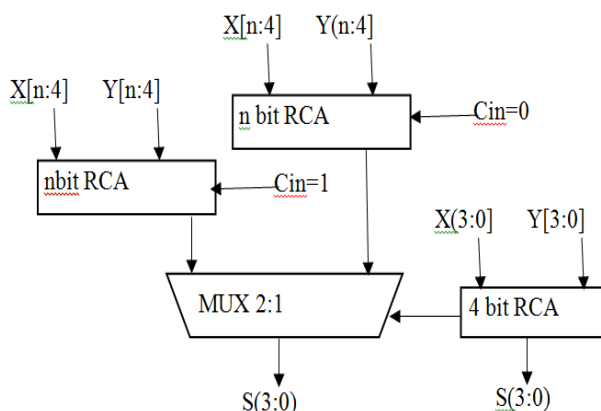


Fig 6. Carry Select Adder.

## 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.

### 1. 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.

### 2. 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 significand 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.

### 3. Significand Bit Calculation:

Significand 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 significand bits of both the operands in four groups. Multiply each group of one operand by each group of second operand. We get 16 product terms. Then we add all of them together very carefully by shifting the term to the left according to which groups of the operands are involved in the product term.

## 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] Paldurai.K and Dr.K.Hariharan "FPGA Implementation of Delay Optimized Single Precision Floating point Multiplier", 2015 International Conference on Advanced Computing and Communication Systems (ICACCS- 2015), Jan. 05-07-2015, Coimbatore, INDIA.
- [2] Irine Padma B.T and Suchitra. K, "Pipelined Floating Point Multiplier Based On Vedic Multiplication Technique," International Journal of Innovative Research in Science, Engineering and Technology (IJIRSET), ISSN: 2347-6710, Volume-3, Special Issue -5, July 2014.
- [3] R. Sai Siva Teja and A. Madhusudhan,"FPGA Implementation of Low- Area Floating Point Multiplier Using Vedic Mathematics", International Journal of Emerging Technology and Advanced Engineering (IJETAEE), Volume-3, Issue -12, December 2013, pp.362-366.
- [4] Priyanka Koneru, Tinnanti Sreenivasu, and Addanki Purna Ramesh, "Asynchronous Single Precision Floating Point Multiplier Using Verilog HDL," International Journal of Advanced Research in Electronics and Communication Engineering (IJARECE), ISSN:2278-909X, Volume-2, Issue - 11, November 2014, pp.885-887.
- [5] I. V. Vaibhav, K. V. Saicharan, B. Sravanthi and D. Srinivasulu, "VHDL Implementation of Floating Point Multiplier using Vedic Mathematics", International Conference on Electrical, Electronics and Communications (ICEEC) , ISBN-978-93-81693-66-03 , June 2014 pp.110-115.
- [6] Ms. Meenu S.Ravi and Mr. Ajit Saraf, "Analysis and study of different multipliers to design floating point MAC units for digital signal processing applications", International Journal of Research in Advent Technology, (IJRAT), ISSN:2321-9637, Volume-2, Issue-3, March 2014, pp.264-267.
- [7] D. Monniaux, "The pitfalls of verifying floating-point computations", ACM Transaction Programming Language System, Vol. 30, No. 3, pp. 1-12, May 2008.
- [8] 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.
- [9] Ragini Parte and Jitendra Jain, "Analysis of Effects of using Exponent Adders in IEEE- 754 Multiplier by VHDL", International Conference on Circuit, Power and Computing Technologies (ICCPCT), 2015 IEEE.
- [10] 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.
- [11] Purna Ramesh Addanki, Venkata Nagaratna Tilak Alapati and Mallikarjuna Prasad Avana, "An FPGA based High Speed IEEE-754 double precision floating point Adder/Subtractor and Multiplier using Verilog", in International Journal of Advance Science and Technology, vol. 52, March 2013.
- [12] Shashank Suresh, Spiridon F. Beldianu and Sotirios G. Ziavras "FPGA and ASIC square root designs for high performance and power efficiency", in 24th IEEE International conference on Application specific-systems, architecture and processors, June 2013.
- [13] 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, October 2011.