

DESIGN AND SIMULATION OF A HIGH SPEED DOUBLE PRECISION FLOATING POINT UNIT USING VERILOG

B.ANIL KUMARI, K.SREENIVASA RAO2,

1Dept. of VLSI System design, 2Associate professor, Dept of E,C,E, Annamacharya Institute of Technology & Sciences, Rajampet, Andhra Pradesh, India

Abstract: Floating point format represents very large or small values, large range is required as the integer representation is no longer appropriate. These values can be represented using the IEEE-754 standard based floating point representation.

In existing system floating point ALU with universal logic gate we can perform addition, subtraction, multiplication and logical operation with less delay and less area using single precision. 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.

The proposed system presents high speed ASIC implementation of a floating point arithmetic unit which can perform addition, subtraction, multiplication, division functions on 64-bit operands that use the IEEE 754-2008 standard. Pre-normalization unit and post normalization units are also discussed along with exceptional handling. All the functions are built by feasible efficient algorithms with several changes incorporated that can improve overall latency, and if pipelined then higher throughput. The algorithms are modeled in Verilog HDL and the RTL code for adder, subtractor, multiplier, divider, square root are synthesized using Xilinx ISE tool.

Index Terms—floating point number, normalization, exceptions, latency, overflow, underflow, etc.

I. INTRODUCTION

An arithmetic circuit which performs digital arithmetic operations has many applications in digital coprocessors, application specific circuits, etc. Because of the advancements in the VLSI technology, many complex algorithms that appeared impractical to put into practice, have become easily realizable today with desired performance parameters so that new designs can be incorporated [2]. The standardized methods to represent floating point numbers have been instituted by the IEEE 754 standard through which the floating point operations can be carried out efficiently with modest storage requirements. The three basic components in IEEE 754 standard floating point numbers are the sign, the exponent, and the mantissa [3]. The sign bit is of 1 bit where 0 refers to positive number and 1 refers to negative number [3]. The mantissa, also called significand which is of 23 bits composes of the fraction and a leading digit which represents the precision bits of the number [3] [2]. The exponent with 8 bits represents both positive and negative exponents. A bias of 127 is added to the exponent to get the stored exponent [2]. Table 1 shows the bit ranges for single (32-bit) and double (64-bit) precision floating-point values [2]. A floating point number representation is shown in table 2. The value of binary floating point representation is as follows where S is sign bit, F is fraction bit and E is exponent field.

$$\text{Value of a floating point number} = (-1)^S \times \text{val}(F) \times 2^{\text{val}(E)}$$

Table 1: Bit Range For Single (32-Bit) And Double (64-Bit) Precision Floating-Point Values

	Sign	Exponent	Fraction	Bias
Single precision	1[31]	8[30-23]	23[22-00]	127
Double precision	1[63]	11[62-52]	52[51-00]	1023

Table 2: Floating Point Number Representation

64 bits		
Sign	Exponent	mantissa
1 bit	11 bits	52 bits

There are four types of exceptions that arise during floating point operations. The Overflow exception is raised whenever the result cannot be represented as a finite value in the precision format of the destination [13].

The Underflow exception occurs when an intermediate result is too small to be calculated accurately, or if the operation's result rounded to the destination precision is too small to be normalized [13]. The Division by zero exception arises when a finite nonzero number is divided by zero [13]. The Invalid operation exception is raised if the given operands are invalid for the operation to be performed [13].

In this paper, ASIC implementation of a high speed FPU has been carried out using efficient addition, subtraction, multiplication, division algorithms. Section II depicts the architecture of the floating point unit and methodology, to carry out the arithmetic operations. Section III presents the arithmetic operations that use efficient algorithms with some modifications to improve latency. Section IV presents the simulation results that have been simulated in Cadence RTL compiler using 180nm process. Section V presents the conclusion.

II. ARCHITECTURE AND METHODOLOGY

The FPU of a double precision floating point unit that performs add, subtract, multiply, divide functions is shown in figure 1 [1]. Two pre-normalization units for addition/subtraction and multiplication/division operations has been given [1].

Post normalization unit also has been given that normalizes the mantissa part [2]. The final result can be obtained after postnormalization. To carry out the arithmetic operations, two IEEE-754 format single precision operands are considered. Pre-normalization of the operands is done. Then the selected operation is performed followed by post-normalizing the output obtained. Finally the exceptions occurred are detected and handled using exceptional handling. The executed operation depends on a three bit control signal (z) which will determine the arithmetic operation is shown in table 3.

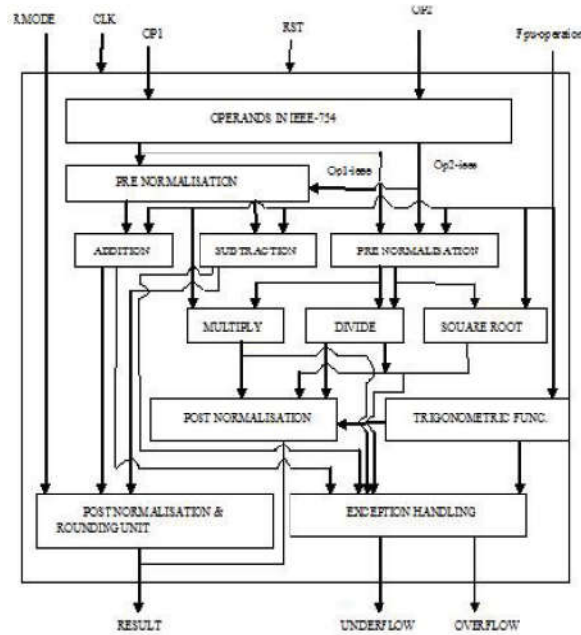


Fig.1: Block Diagram of floating point arithmetic unit [1]

Table 3: Floating Point Unit Operations

z(control signal)	Operation
2'b000	Addition
2'b001	Subtraction
2'b010	Multiplication
2'b011	Division
2'b100	Square root

III 64 BIT FLOATING POINT ARITHMETIC UNIT

A. Addition Unit:

One of the most complex operations in a floating point unit comparing to other functions which provides major delay and also considerable area. Many algorithms has been developed which focused to reduce the overall latency in order to improve performance. The floating point addition operation is carried out by first checking the zeros, then aligning the significand, followed by adding the two significands using an efficient architecture.

The obtained result is normalized and is checked for exceptions. To add the mantissas, a high speed carry lookahead has been used to obtain high speed. Traditional carry look ahead adder is constructed using AND, XOR and NOT gates.

The implemented modified carry look ahead adder uses only NAND and NOT gates which decreases the cost of carry lookahead adder and also enhances its speed also [4]. The 16 bit modified carry look ahead adder is shown in figure 2 and the metamorphosis of partial full adder is shown in figure 3 using which, a 24 bit carry look ahead adder has been constructed and performed the addition operation.

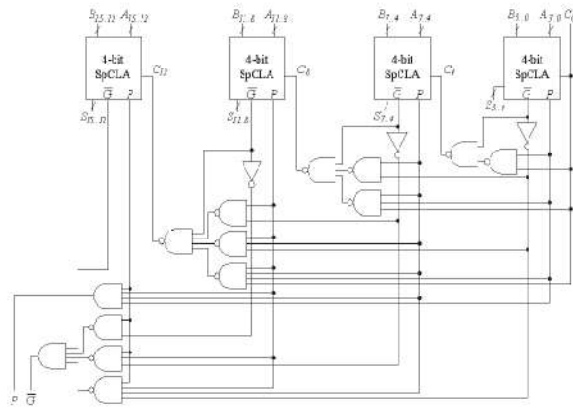


Fig.2: 16 bit modified carry look ahead adder [4]

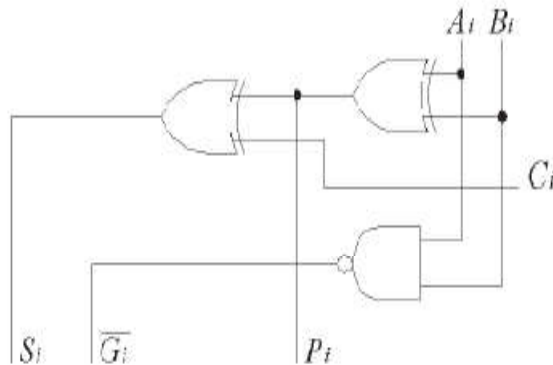


Fig.3: Metamorphosis of partial full adder [4]

B. Subtraction Unit:

Subtraction operation is implemented by taking 2's complement of second operand. Similar to addition operation, subtraction consists of three major tasks pre normalization, addition of mantissas, post normalization and exceptional handling. Addition of mantissas is carried out using the 24 bit MCLA shown in figure 2 and figure 3.

C. Multiplication Algorithm

Constructing an efficient multiplication module is an iterative process and 2n-digit product is obtained from the product of two n-digit operands. In IEEE 754 floating-point multiplication, the two mantissas are multiplied, and the two exponents are added. Here first the exponents are added from which the exponent bias (1023) is removed. Then mantissas have been multiplied using feasible algorithm and the output sign bit is determined by XORing the two input sign bits. The obtained result has been normalized and checked for exceptions. To multiply the mantissas Bit Pair Recoding (or Modified Booth Encoding) algorithm has been used, because of which the number of partial products gets reduced by about a factor of two, with no requirement of pre-addition to produce the partial products. It recodes the bits by considering three bits at a time.

Bit Pair Recoding algorithm increases the efficiency of multiplication by pairing. To further increase the efficiency of the algorithm and decrease the time complexity, Karatsuba algorithm can be paired with the bit pair recoding algorithm. One of the fastest multiplication algorithms is Karatsuba algorithm which reduces the multiplication of two n-digit numbers to $3n \log_2 3 \sim 3n \cdot 1.585$ single-digit multiplications and therefore faster than the classical algorithm, which requires n^2 single-digit products [11]. It allows to compute the product of two large numbers x and y using three multiplications of smaller numbers, each with about half as many digits as x or y, with some additions and digit shifts instead of four multiplications [11]. The steps are carried out as follows

Let x and y be represented as n -digit numbers with base B and $m < n$.

$$x = x_{1B}m + x_0$$

$$y = y_{1B}m + y_0$$

Where x_0 and y_0 are less than B^m [11]. The product is then

$$xy = (x_{1B}m + x_0)(y_{1B}m + y_0) = c_{1B}2m + b_{1B}m + a_1$$

$$\text{Where } c_1 = x_1y_1$$

$$b_1 = x_1y_0 + x_0y_1$$

$$a_1 = x_0y_0.$$

$$p_1 = p_1 - z_2 - z_0$$

$$p_1 = (x_1 + x_0)(y_1 + y_0)$$

Here c_1 , a_1 , p_1 has been calculated using bit pair recoding algorithm. Radix-4 modified booth encoding has been used which allows for the reduction of partial product array by half [n/2]. The bit pair recoding table is shown in table 3. In the implemented algorithm for each group of three bits (y_{2i+1} , y_{2i} , y_{2i-1}) of multiplier, one partial product row is generated according to the encoding in table 3.

Radix-4 modified booth encoding (MBE) signals and their respective partial products has been generated using the figures 4 and 5. For each partial product row, figure 4 generates the one, two, and neg signals. These values are then given to the logic in figure 5 with the bits of the multiplicand, to produce the whole partial product array. To prevent the sign extension the obtained partial products are extended as shown in figure 6 and the product has been calculated using carry save select adder.

Table 3: Bit-Pair Recoding [11]

BIT PATTERN		OPERATION
0 0 0	NO OPERATION	
0 0 1	1xa	prod=prod+a;
0 1 0	2xa-a	prod=prod+a;
0 1 1	2xa	prod=prod+2a;
1 0 0	-2xa	prod=prod-2a;
1 0 1	-2xa+a	prod=prod-a;
1 1 0	-1xa	prod=prod-a;
1 1 1	NO OPERATION	

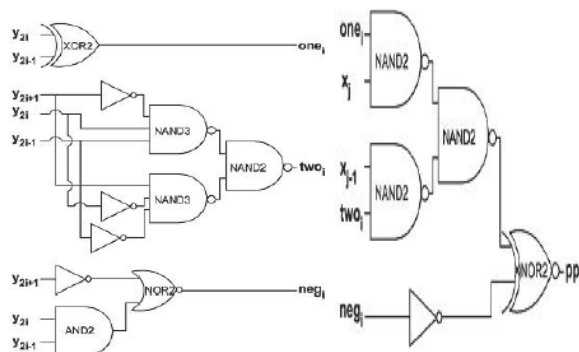


Fig.4: MBE signal generation [10] Fig.5: Partial Product Generation [10]

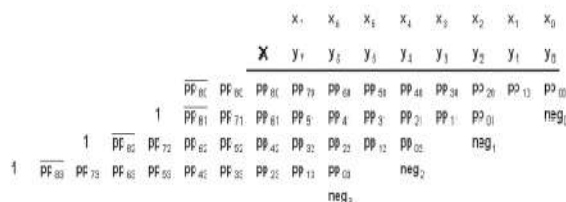


Fig.6: Sign prevention extension of partial products [10]

D. Division Algorithm

Division is the one of the complex and time-consuming operation of the four basic arithmetic operations. Division operation has two components as its result i.e. quotient and a remainder when two inputs, a dividend and a divisor are given. Here the exponent of result has been calculated by using the equation,

$$e_0 = e_A - e_B + \text{bias} (127) - z_A + z_B$$

followed by division of fractional bits [5] [6]. Sign of result has been calculated from XORing sign of two operands. Then the obtained quotient has been normalized [5] [6].

Division of the fractional bits has been performed by using non restoring division algorithm which is modified to improve the delay. The non-restoring division algorithm is the fastest among the digit recurrence division methods [5] [6]. Generally restoring division require two additions for each iteration if the temporary partial remainder is less than zero and this results in making the worst case delay longer [5] [6]. To decrease the delay during division, the non-restoring division algorithm was introduced which is shown in figure 7. Non-restoring division has a different quotient set i.e it has one and negative one, while restoring division has zero and one as the quotient set [5] [6].

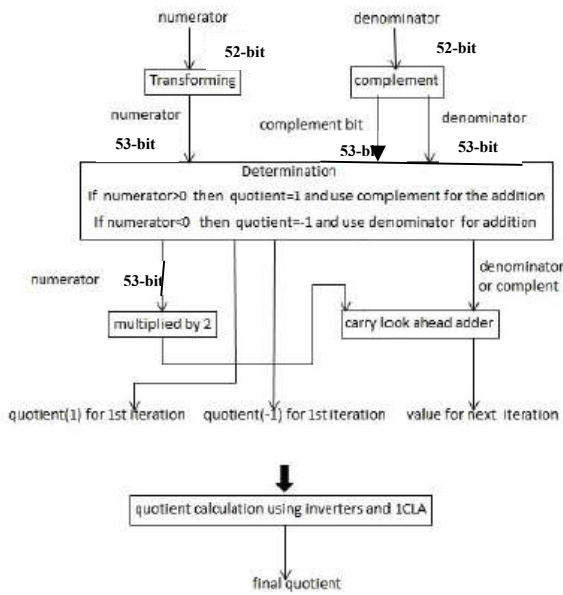


Fig.7: Non Restoring Division algorithm

Using the different quotient set, reduces the delay of non-restoring division compared to restoring division. It means, it only performs one addition per iteration which improves its arithmetic performance [6]. The delay of the multiplexer for selecting the quotient digit and determining the way to calculate the partial remainder can be reduced through rearranging the order of the computations. In the implemented design the adder for calculating the partial remainder and the multiplexer has been performed at the same time, so that the multiplexer delay can be ignored since the adder delay is generally longer than the multiplexer delay.

Second, one adder and one inverter are removed by using a new quotient digit converter. So, the delay from one adder and one inverter connected in series will be eliminated.

E. Square Root Unit

Square root operation is difficult to implement because of the complexity of the algorithms. Here a low cost iterative single precision non-restoring square root algorithm has been presented that uses a traditional adder/subtractor whose operation latency is 25 clock cycles and the issue rate is 24 clock cycles. If the biased exponent is even, the biased exponent is added to 126 and divided by two and mantissa is shifted to its left by 1 bit before computing its square root. Here before shifting the mantissa bits are stored in 52 bit register as 1.xx..xx.

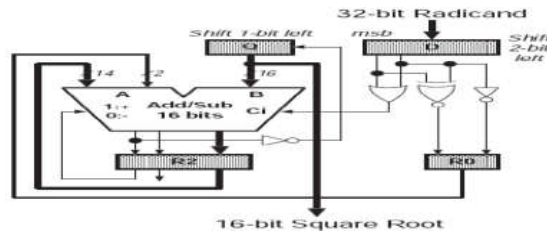


Fig.7: Non Restoring square root circuitry [15] [16]

After shifting it becomes $1x.xx...$. If the biased exponent is odd, the odd exponent is added to 127 and divided by two. Themantissa. The square root of floating point number has been calculated by using non restoring square root circuitry which is shown in figure 8 [15] [16].

IV SIMULATION RESULTS

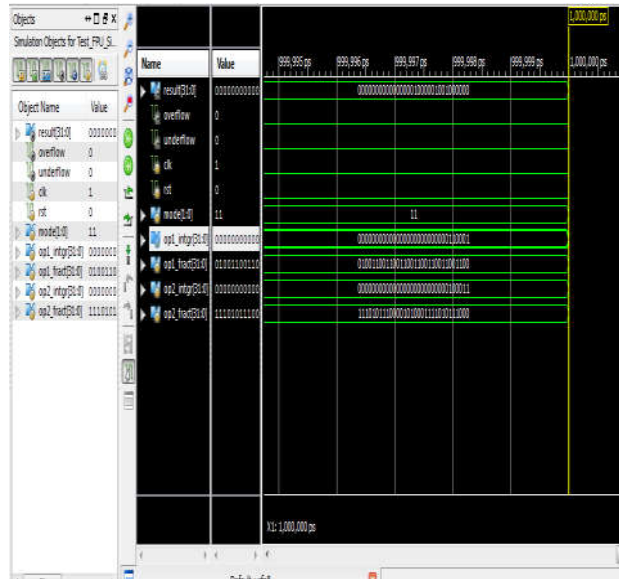


Figure : Implementation of 32 bit single precision operation

- In single precision floating point unit it executes only 32-bits.

Sotime consumption is more and speed alsoreduced. So for avoiding that problemwe propose a double precision floating point unit.

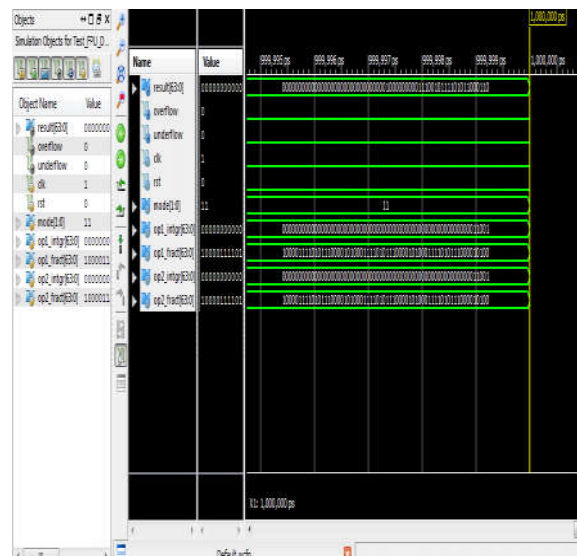


Figure : Implementation of 64 bit double precision

V CONCLUSION

The implementation of a high speed double precision FPU has been presented. The design has been synthesized with Xilinx tool. Strategies have been employed to realize optimal hardware and power efficient architecture. The layout generation of the presented architecture using the backend flow is an ongoing process and is being done using Cadence RTL compiler with 180nm process technology. Hence it can be concluded that this FPU can be effectively used for ASIC implementations which can show comparable efficiency and speed and if pipelined then higher throughput may be obtained.

REFERENCES

- [1] Rudolf Usselman, "Open Floating Point Unit, The Free IP Cores Projects".
- [2] Edvin Catovic, Revised by: Jan Andersson, "GRFPU – High Performance IEEE754 Floating Point Unit", Gaisler Research, Första Långatan 19, SE413 27 Göteborg, and Sweden.
- [3] David Goldberg, "What Every Computer Scientist Should Know About Floating-Point Arithmetic", *ACM Computing Surveys*, Vol 23, No 1, March 1991, Xerox Palo Alto Research Center, 3333 Coyote Hill Road, Palo Alto, California 94304.
- [4] Yu-Ting Pai and Yu-Kumg Chen, "The Fastest Carry Lookahead Adder", Department of Electronic Engineering, Huafan University.
- [5] Prof. Kris Gaj, Gaurav, Doshi, Hiren Shah, "Sine/Cosine using CORDIC Algorithm".
- [6] S. F. Oberman and M. J. Flynn, "Division algorithms and implementations," *IEEE Transactions on Computers*, vol. 46, pp. 833–854, 1997.
- [7] Milos D. Ercegovac and Tomas Lang, *Division and Square Root: Digit-Recurrence Algorithms and Implementations*, Boston: Kluwer Academic Publishers, 1994.
- [8] ANSI/IEEE Standard 754-1985, *IEEE Standard for Binary Floating-Point Arithmetic*, 1985.
- [9] Behrooz Parhami, *Computer Arithmetic - Algorithms and Hardware Designs*, Oxford: Oxford University Press, 2000.
- [10] Steven Smith, (2003), *Digital Signal Processing-A Practical guide for Engineers and Scientists*, 3rd Edition, Elsevier Science, USA.