

Design of High Speed IEEE-754 Double Precision Floating Point Multiplier Using Dadda Algorithm

Vivek .K(M.Tech), Gundam Narahari, and S.Jagadeesh

ABSTRACT: Floating Point (FP) multiplication is widely used in large set of scientific and signal processing computation. Multiplication is one of the common arithmetic operations in these computations. A high speed floating point double precision multiplier is implemented in HDL. This paper presents a high speed binary double precision floating point multiplier based on Dadda Algorithm. To improve speed multiplication of mantissa is done using Dadda multiplier replacing Carry Save Multiplier. 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. **KEY WORDS:** Dadda Algorithm, Double precision, Floating point, Multiplier, IEEE-754, Verilog HDL.

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. 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 significand and exponent is greater than zero and smaller than 1023. The real number is represented by equations (1) & (2).

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

$$\text{Value} = (-1^{\text{Sign bit}}) * 2^{(\text{Exponent} - 1023)} * (1.\text{Mantissa}) \quad (2)$$

Floating point implementation has been the interest of many researchers. In an IEEE-754 single precision pipelined floating point multiplier is implemented with custom 16/18 bit three stage pipelined floating point multiplier, that doesn't support rounding modes [1]. L.Louca, T.A.Cook, W.H. Johnson [2] implemented a single precision floating point multiplier by using a digit-serial multiplier. The design achieved 2.3 MFlops and doesn't support rounding modes. The multiplier handles the overflow and underflow cases but rounding is not implemented. The design achieves 30 I MFLOPs with latency of three clock cycles. The multiplier was verified against Xilinx floating point multiplier core.

1.Vivek .K(M.Tech),2.Gundam Narahari,3.S.Jagadeesh, 1.M.TechStudent in SSJ Engineering College,Hyderabad, anudeep405@gmail.com.
2.Ass.prof In ECE Dept,SSJ Engg,Hyderabad,3. HOD and Ass.prof In ECE Dept,SSJ Engg,Hyderabad, jaaga.ssjecc@gmail.com.



Figure1. IEEE 754 Double Precision Floating Point Format. The double precision floating point multiplier presented here is based on IEEE-754 binary floating standard. We have designed a high speed double precision floating point multiplier using Verilog language. It operates at a very high frequency of 414.714 MFlops and occupies 648 slices. It handles the overflow, underflow cases and rounding mode.

I.FLOATING POINT MULTIPLICATION ALGORITHM

Multiplying two numbers in floating point format is done by

1. Adding the exponent of the two numbers then subtracting the bias from their result.
2. Multiplying the significand of the two numbers
3. 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 (leading one).

The following steps are necessary to multiply two floating point numbers.

1. Multiplying the significand i.e. (I.M1 * I.M2)
2. Placing the decimal point in the result
3. Adding the exponents i.e. (E1 + E2 - Bias)
4. Obtaining the sign i.e. s1 xor s2
5. Normalizing the result i.e. obtaining 1 at the MSB of the results "significand"
6. Rounding the result to fit in the available bits
7. Checking for underflow/overflow occurrence

II.IMPLEMENTATION OF DOUBLE PRECISION FLOATING POINT MULTIPLIER

In this paper we implemented a double precision floating point multiplier with exceptions and rounding. Figure 2 shows the multiplier structure that includes exponents addition, significand multiplication, and sign calculation. Figure 3 shows the multiplier, exceptions and rounding that are independent and are done in parallel.

A_exponent B_exponent A_mantissa B_mantissa

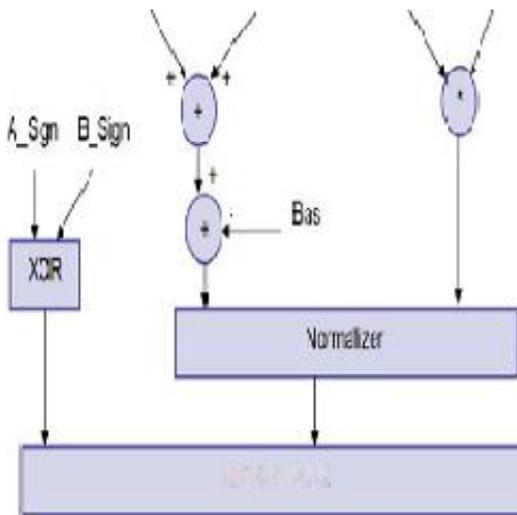


Figure 2. Multiplier structure

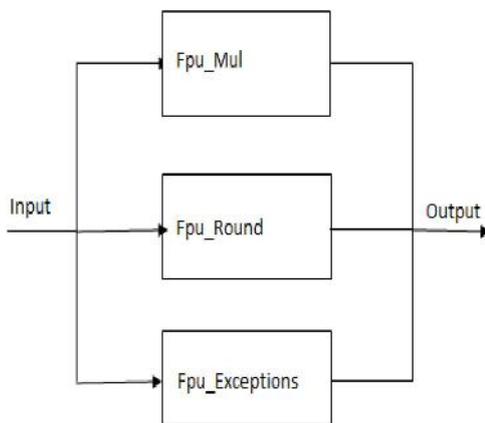


Figure 3. Multiplier structure with rounding and exceptions

III. MULTIPLIER

Existing Multiplier:

Carry Save Multiplier:

This unit is used to multiply the two unsigned significant numbers and it places the decimal point in the multiplied product. The unsigned significant multiplication is done on 24 bit. The result of this significant multiplication will be called the IR. Multiplication is to be carried out so as not to affect the whole multiplier's performance. In this carry save multiplier architecture is used for 24X24 bit as it has a moderate speed with a simple architecture. In the carry save multiplier, the carry bits are passed diagonally downwards (i.e. the carry bit is propagated to the next stage). Partial products are generated by ANDing the inputs of two numbers and passing them to the appropriate adder. Carry save multiplier has three main stages:

1. The first stage is an array of half adders.
2. The middle stages are arrays of full adders. The number of middle stages is equal to the significant size minus two.
3. The last stage is an array of ripple carry adders. This stage is called the vector merging stage.

The count of adders (Half adders and Full adders) in each stage is equal to the significant size minus one. For example, a 4x4 carry save multiplier is shown in Figure 8 and it has the following stages:

1. The first stage consists of three half adders.

2. Two middle stages; each consists of three full adders.
3. The vector merging stage consists of one half adder and two full adders.

The decimal point is placed between bits 45 and 46 in the significand multiplier result. The multiplication time taken by the carry save multiplier is determined by its critical path. The critical path starts at the AND gate of the first partial products (i.e. a_1b_0 and a_0b_1), passes through the carry logic of the first half adder and the carry logic of the first full adder of the middle stages, then passes through all the vector merging adders. The critical path is marked in bold in Figure 4.

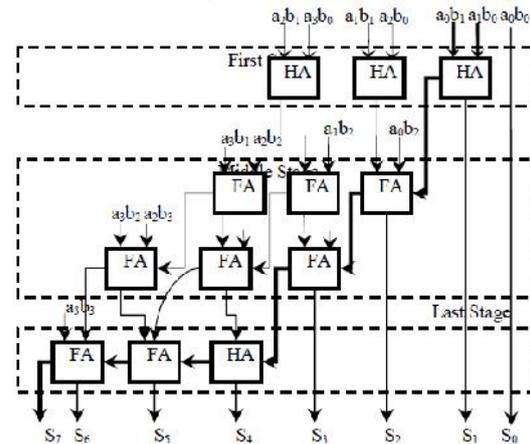


Fig. 4. 4x4 bit Carry Save multiplier

In Figure 4

1. Partial product: $a_i b_j$ a_i and b_j
2. HA: half adder.
3. FA: full adder.

Proposed multiplier

Dadda Multiplier:

Dadda proposed a sequence of matrix heights that are predetermined to give the minimum number of reduction stages. To reduce the N by N partial product matrix, dadda multiplier develops a sequence of matrix heights that are found by working back from the final two-row matrix. In order to realize the minimum number of reduction stages, the height of each intermediate matrix is limited to the least integer that is no more than 1.5 times the height of its successor.

The process of reduction for a dadda multiplier [3] is developed using the following recursive algorithm

1. Let $d_1=2$ and $d_{j+1} = \lceil 1.5*d_j \rceil$, where d_j is the matrix height for the j th stage from the end. Find the smallest j such that at least one column of the original partial product matrix has more than d_j bits.
2. In the j th stage from the end, employ (3, 2) and (2, 2) counter to obtain a reduced matrix with no more than d_j bits in any column.
3. Let $j = j-1$ and repeat step 2 until a matrix with only two rows is generated.

This method of reduction, because it attempts to compress each column, is called a column compression technique.

Another advantage of utilizing Dadda multipliers is that it utilizes the minimum number of (3, 2) counters. Therefore, the number of intermediate stages is set in terms of lower bounds: 2, 3, 4, 6, 9 . . .

Fig. 6. Simulation Result of Double Precision Floating Point Multiplier

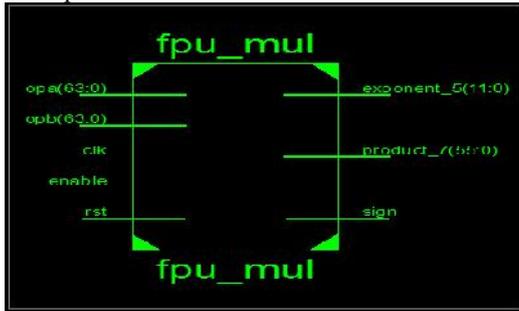


Fig.7. RTL Schematic for top level module

CONCLUSION

The double precision floating point multiplier supports the IEEE-754 binary interchange format. The design achieved the operating frequency of 414.714 MFLOOPS with area of 648 slices. The implemented design is verified with single precision floating point multiplier [4] and Xilinx core, it provides high speed and supports double precision, which gives more accuracy compared to single precision. This design handles the overflow, underflow, and truncation rounding mode.

REFERENCES

- [1] 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.
- [2] 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), pp. 107-116,1996.
- [3] Whytney J. Townsend, Earl E. Swartz, "A Comparison of Dadda and Wallace multiplier delays". Computer Engineering Research Center, The University of Texas.
- [4] Mohamed AI-Ashraf), Ashraf Salem, Wagdy Anis., "An Efficient Implementation of Floating Point Multiplier ", Saudi International Electronics, Communications and Photonics Conference (SIEPCPC), pp. 1-5,24-26 April 2011.
- [5] B. Lee and N. Burgess, "Parameterisable Floating-point Operations on FPG A," Conference Record of the ThirtySixth Asilomar Conference on Signals, Systems, and Computers, 2002.
- [6] Xilinx13.4, Synthesis and Simulation Design Guide", UG626 (v13.4) January 19, 2012.
- [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] L. Louca, T. A. Cook, and W. H. Johnson, "Implementation of IEEE Single Precision Floating Point Addition and Multiplication on FPGAs," Proceedings of 83 the IEEE Symposium on FPGAs for Custom Computing Machines (FCCM'96), pp. 107-116, 1996.