

An Efficient Implementation of Floating Point Multiplier

ASIYA THAPASWIN PATTAN
M. Tech Student, Dept of ECE ,
Vaagdevi Institute of Technology & Science,
Proddatur, Kadapa (DT), AP-516361

C.Md. ASLAM
HOD, Dept. of ECE,
Vaagdevi Institute of Technology & Science,
Proddatur, Kadapa (DT), AP-516361

V.RAMESH
Assistant Professor, Dept. of ECE,
Vaagdevi Institute of Technology & Science,
Kadapa, Proddatur, Kadapa Dt. A.P.-516361

Abstract—This paper describes an efficient implementation of an IEEE 754 single precision floating point multiplier targeted for Xilinx Virtex-5 FPGA. VHDL is used to implement a technology-independent pipelined design. The multiplier implementation handles the overflow and underflow cases. Rounding is not implemented to give more precision when using the multiplier in a multiply and Accumulate (MAC) unit. With latency of three clock cycles the design achieves 301 MFLOPs. The multiplier was verified against Xilinx floating point multiplier core.

Keywords -floating point; multiplication; FPGA; CAD design flow.

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 only on single precision normalized binary interchange format. 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; in this case the real number is represented by (1)

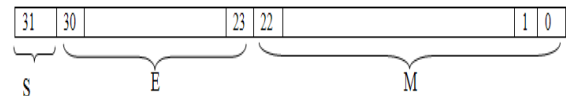


Figure 1. IEEE single precision floating point format

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

Where

$$M = m_{22} 2^{-1} + m_{21} 2^{-2} + m_{20} 2^{-3} + \dots + m_1 2^{-22} + m_0 2^{-23};$$

$$Bias = 127.$$

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, and 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). Floating-point implementation on FPGAs has been the interest of many researchers.

II. FLOATING POINT MULTIPLICATION ALGORITHM

As stated in the introduction, normalized floating point numbers have the form of $Z = (-1^S) * 2^{(E - Bias)} * (1.M)$. To multiply two floating point numbers the following is done:

1. Multiplying the significand; i.e. $(1.M1 * 1.M2)$
2. Placing the decimal point in the result
3. Adding the exponents; i.e. $(E1 + E2 - Bias)$
4. Obtaining the sign; i.e. $s1 \text{ xor } s2$
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

Consider a floating point representation similar to the IEEE 754 single precision floating point format, but with a reduced number of mantissabits (only 4) while

still retaining the hidden '1' bit for normalized numbers:

$$A = 0\ 1000100\ 0100 = 40,$$

$$B = 1\ 1000001\ 1110 = -7.5$$

To multiply A and B

$$\begin{array}{r}
 1.0100 \\
 \times 1.1110 \\
 \hline
 00000 \\
 10100 \\
 10100 \\
 10100 \\
 \hline
 10100 \\
 1001011000 \\
 \hline
 10.01011000 \\
 \hline
 \begin{array}{r}
 10000100 \\
 + 10000001 \\
 \hline
 10000101
 \end{array}
 \end{array}$$

The exponent representing the two numbers is already shifted/biased by the bias value (127) and is not the true exponent; i.e. $E_A = E_{A\text{-true}} + \text{bias}$ and $E_B = E_{B\text{-true}} + \text{bias}$

And

$$E_A + E_B = E_{A\text{-true}} + E_{B\text{-true}} + 2\ \text{bias}$$

So we should subtract the bias from the resultant exponent otherwise the bias will be added twice.

$$\begin{array}{r}
 10000101 \\
 - 01111111 \\
 \hline
 10000110
 \end{array}$$

4. Obtain the sign bit and put the result together:

$$1\ 10000110\ 10.01011000$$

5. Normalize the result so that there is a 1 just before the radix point (decimal point). Moving the radix point one place to the left increments the exponent by 1; moving one place to the right decrements the

Exponent by 1.

$$1\ 10000110\ 10.01011000\ (\text{before normalizing})$$

$$1\ 10000111\ 1.001011000\ (\text{normalized})$$

The result is (without the hidden bit):

$$1\ 10000111\ 00101100$$

6. The mantissa bits are more than 4 bits (mantissa available bits); rounding is needed.

If we applied the Truncation rounding mode then the stored value is:

$$1\ 10000111\ 0010.$$

In this paper we present a floating point multiplier in which rounding support isn't implemented. Rounding support can be added as a separate unit that can be accessed by the multiplier or by a floating point adder, thus accommodating for more precision if the multiplier is connected directly to an adder in a MAC unit. Fig.2 shows the multiplier structure; Exponents addition, Significand multiplication, and Result's sign calculation are independent and are done in parallel. The significand multiplication is done on two 24 bit numbers and results in a 48 bit product, which we will call the intermediate product (IP). The IP is represented as (47 downto 0) and the decimal point is located between bits 46 and 45 in the IP. The following sections detail each block of the floating point multiplier.

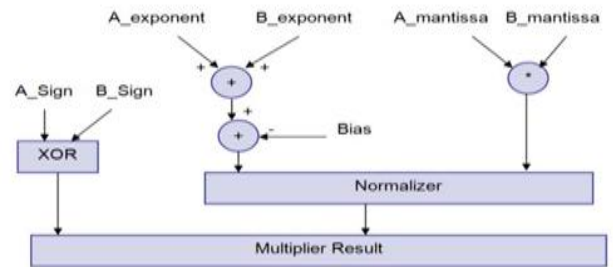


Figure 2. Floating point multiplier block diagram

III. HARDWARE OF FLOATING POINT MULTIPLIER

A. Sign bit calculation

Multiplying two numbers results in a negative sign number if one of the multiplied numbers is of a negative value. By the aid of a truth table we find that this can be obtained by XORing the sign of two inputs.

B. Unsigned Adder (for exponent addition)

This unsigned adder is responsible for adding the exponent of the first input to the exponent of the second input and subtracting the Bias (127) from the addition result (i.e. $A_exponent + B_exponent - \text{Bias}$). The result of this stage is called the intermediate exponent. The add operation is done on 8 bits, and there is no need for a quick result because most of the calculation time is spent in the significand multiplication process (multiplying 24 bits by 24 bits); thus we need a moderate exponent adder and a fast significand multiplier. An 8-bit ripple carry adder is used to add the two input exponents. As shown in Fig. 3 a ripple carry adder is a chain of cascaded full adders and one

half adder; each full adder has three inputs (A, B, Ci) and two outputs (S, Co). The carry out (Co) of each adder is fed to the next full adder (i.e. each carry bit "ripples" to the next full adder).

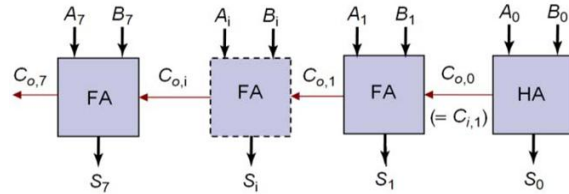
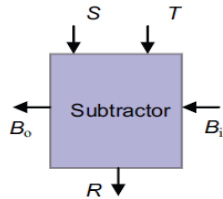


Figure 3. Ripple Carry Adder

The addition process produces an 8 bit sum (S7 to S0) and a carry bit (Co, 7). These bits are concatenated to form a 9 bit addition result (S8 to S0) from which the Bias is subtracted. The Bias is subtracted using an array of ripple borrow subtractors.

A normal subtractor has three inputs (minuend (S), subtrahend (T), Borrow in (Bi)) and two outputs (Difference (R), Borrow out (Bo)).



The subtractor logic can be optimized if one of its inputs is a constant value which is our case, where the Bias is constant (127|10 = 00111111|2).

Table I shows the truth table for a 1-bit subtractor with the input T equal to 1 which we will call "one subtractor (OS)"

TABLE I. 1-BIT SUBTRACTOR WITH THE INPUT T=1

S	T	Bi	Difference(R)	Bo
0	1	0	1	1
1	1	0	0	0
0	1	1	0	1
1	1	1	1	1

The Boolean Equation (2) and (3) represent this subtractor:

$$Difference(R) = \overline{S} \oplus B_i \quad (2)$$

$$Borrowout (Bo) = \overline{S} + B_i \quad (3)$$

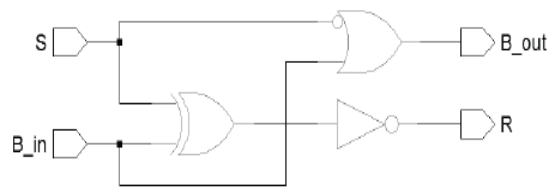


Figure 4. 1-bit subtractor with the input T = 1

Table II shows the truth table for a 1-bit subtractor with the input T equal to 0 which we will call "zero subtractor (ZS)"

TABLE II. 1-BIT SUBTRACTOR WITH THE INPUT T=0

S	T	Bi	Difference(R)	Bo
0	0	0	0	0
1	0	0	1	0
0	0	1	1	1
1	0	1	0	0

The Boolean Equation (4) and (5) represent this subtractor:

$$Difference(R) = S \oplus B_i \quad (4)$$

$$Borrowout (Bo) = \overline{S} \cdot B_i \quad (5)$$

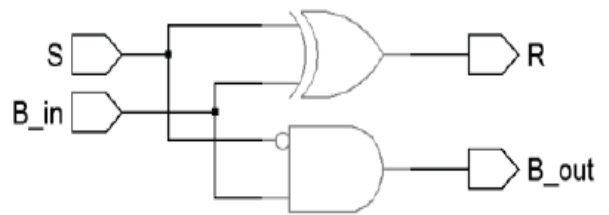


Figure 5. 1-bit subtractor with the input T = 0

Fig. 6 shows the Bias Subtractor which is a chain of 7 one subtractors (OS) followed by 2 zero subtractors (ZS); the borrow output of each Subtractor is fed to the next Subtractor. If an underflow occurs then $E_{result} < 0$ and the number is out of the IEEE 754 single precision normalized numbers range; in this case the output is signaled to 0 and an underflow flag is asserted.

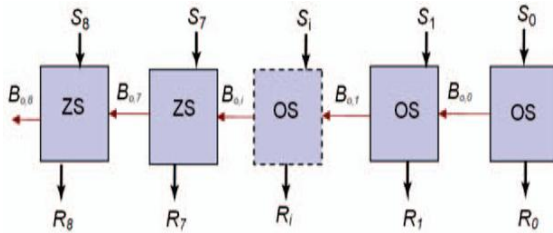


Figure6. Ripple Borrow Subtractor

C. Unsigned Multiplier (for significand multiplication)

This unit is responsible for multiplying the unsigned significand and placing the decimal point in the multiplication product. The result of significand multiplication will be called the intermediate product (IP).

The unsigned significand multiplication is done on 24 bit. Multiplier performance should be taken into consideration so as not to affect the whole multiplier's performance. A 24x24 bit carry save multiplier architecture is used 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 made by ANDing the inputs together 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.
- 3- The last stage is an array of ripple carry adders. This stage is called the vector merging stage.

The number of adders (Half adders and Full adders) in each stage is equal to the significand size minus one. For example, a 4x4 carry save multiplier is shown in Fig. 7 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 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 Fig. 7

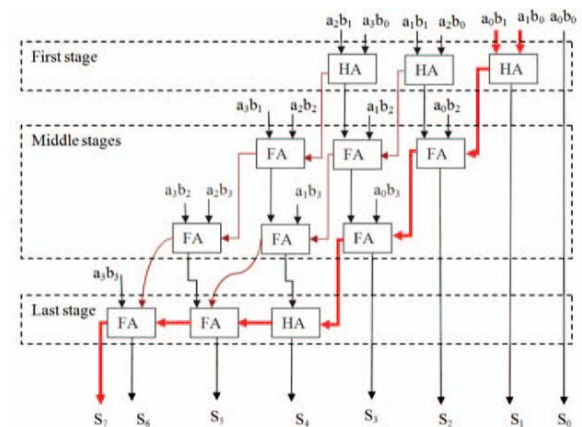


Figure7. 4x4 bit Carry Save multiplier

In Fig. 7:

- 1- Partial product: $a_ib_j = a_i$ and b_j
- 2- HA: half adder
- 3- FA: full adder

D. Normalizer

The result of the significand multiplication (intermediate product) must be normalized to have a leading '1' just to the left of the decimal point (i.e. in the bit 46 in the intermediate product). Since the inputs are normalized numbers then the intermediate product has the leading one at bit 46 or 47

- 1- If the leading one is at bit 46 (i.e. to the left of the decimal point) then the intermediate product is already a normalized number and no shift is needed.
- 2- If the leading one is at bit 47 then the intermediate product is shifted to the right and the exponent is incremented by 1.

The shift operation is done using combinational shift logic made by multiplexers.

Fig. 8 shows a simplified logic of a Normalizer that has an 8 bit intermediate product input and a 6 bit intermediate exponent input.

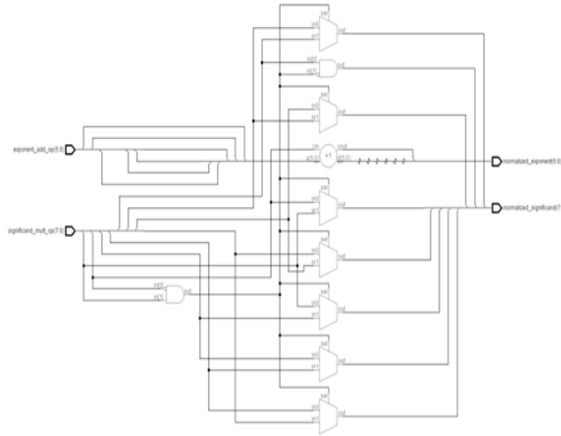


Figure8. Simplified Normalizer logic

IV. UNDERFLOW/OVERFLOW DETECTION

Overflow/underflow means that the result’s exponent is too large/small to be represented in the exponent field. The Exponent of the result must be 8 bits in size, and must be between 1 and 254 otherwise the value is not a normalized one. Between 1 and 254 otherwise the value is not a normalized one.

An overflow may occur while adding the two exponents or during normalization. Overflow due to exponent addition may be compensated during subtraction of the bias; resulting in a normal output value (normal operation). An underflow may occur while subtracting the bias to form the intermediate exponent. If the intermediate exponent < 0 then it’s an underflow that can never be compensated; if the intermediate exponent = 0 then it’s an underflow that may be compensated during normalization by adding 1 to it.

Assume that E1 and E2 are the exponents of the two numbers A and B respectively; The result’s exponent is calculated by (6)

$$E_{result} = E_1 + E_2 - 127 \tag{6}$$

E1 and E2 can have the values from 1 to 254; resulting in Eresult having values from -125 (2-127) to 381 (508-127); but for normalized numbers, Eresult can only have the values from 1 to 254. Table III summarizes the

Eresult different values and the effect of normalization on it.

TABLE III. NORMALIZATION EFFECT ON RESULT’S EXPONENT AND OVERFLOW/UNDERFLOW DETECTION

E_{result}	Category	Comments
$-125 \leq E_{result} < 0$	Underflow	Can't be compensated during normalization
$E_{result} = 0$	Zero	May turn to normalized number during normalization (by adding 1 to it)
$1 < E_{result} < 254$	Normalized number	May result in overflow during normalization
$255 \leq E_{result}$	Overflow	Can't be compensated

V. PIPELINING THE MULTIPLIER

In order to enhance the performance of the multiplier, three pipelining stages are used to divide the critical path thus increasing the maximum operating frequency of the multiplier. The pipelining stages are imbedded at the following locations:

1. In the middle of the significand multiplier, and in the middle of the exponent adder (before the bias subtraction).
2. after the significand multiplier, and after the exponent adder.
3. At the floating point multiplier outputs (sign, exponent and mantissa bits).

Fig. 9 shows the pipelining stages as dotted lines.

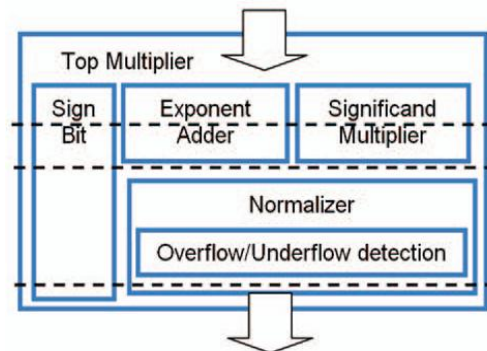


Figure 9. Floating point multiplier with pipelined stages

Three pipelining stages mean that there is latency in the output by three clocks. The synthesis tool “retiming” option was used so that the synthesizer uses its optimization logic to better place the pipelining registers across the critical path

VI. RESULT

Fig.10 shows the simulation result of top multiplier.

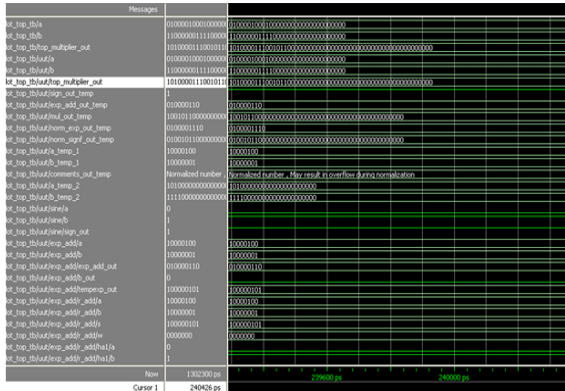


Figure 10. Simulation result of top multiplier

Fig.11 shows the RTL diagram of top multiplier which is having inputs A, B of 32 bit and output of 56bit.

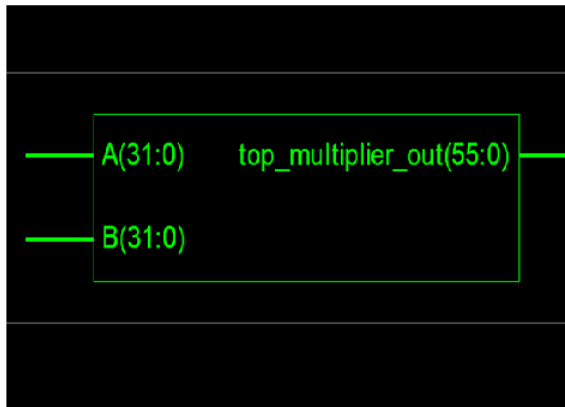


Figure 11. RTL diagram of top multiplier

VII. IMPLEMENTATION AND TESTING

The whole multiplier (top unit) was tested against the Xilinx floating point multiplier core generated by Xilinx coregen. Xilinx core was customized to have two flags to indicate overflow and underflow, and to have a maximum latency of three cycles. Xilinx core implements the “round to nearest” rounding mode.

The area of Xilinx core is less than the implemented floating point multiplier because the latter doesn’t truncate/round the 48 bits result of the mantissa multiplier which is reflected in the amount of function generators and registers used to perform operations on the extra bits; also the speed of Xilinx core is

affected by the fact that it implements the round to nearest rounding mode.

VIII. CONCLUSIONS AND FUTURE WORK

This paper presents an implementation of a floating point multiplier that supports the IEEE 754-2008 binary interchange format; the multiplier doesn’t implement rounding and just presents the significant multiplication result as is (48 bits); this gives better precision if the whole 48 bits are utilized in another unit; i.e. a floating point adder to form a MAC unit. The design has three pipelining stages and after implementation on a Xilinx Virtex5 FPGA it achieves 301 MFLOPs.

REFERENCES

- [1] IEEE 754-2008, IEEE Standard for Floating-Point Arithmetic, 2008.
- [2] 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.
- [3] 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.
- [4] 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.
- [5] A. Jaenicke and W. Luk, "Parameterized Floating-Point Arithmetic on FPGAs", Proc. of IEEE ICASSP, 2001, vol. 2, pp. 897-900.
- [6] B. Lee and N. Burgess, “Parameterisable Floating-point Operations on FPGA,” Conference Record of the Thirty-Sixth Asilomar Conference on Signals, Systems, and Computers, 2002
- [7] “DesignChecker User Guide”, HDL Designer Series 2010.2a, Mentor Graphics, 2010
- [8] “PrecisionR Synthesis User’s Manual”, Precision RTL plus 2010a update 2, Mentor Graphics, 2010.



Mrs. Asiya Thapaswin Pattan, is currently doing post graduation in Vaagdevi institute of Tech & Science, Proddatur, Kadapa (Dt), A.P with the specialization of VLSI



Mr. V. Ramesh, is an Assistant Professor in ECE dept. at Vaagdevi institute of Tech & Science, Proddatur. He obtained his B.Tech degree in ECE from MeRITS, JNTUA, Udayagiri in 2007, M.Tech degree in Electronic Instrumentation and Communication Systems from S.V. University, Tirupati in 2009. He has 4 years of teaching experience and his area of interest includes Electronics instrumentation and Antennas. He has published 5 papers in referred international journals and also 3 papers at national level conferences.



Mr. C. Mahammed Aslam is the HOD, Dept. of ECE at Vagdevi institute of Tech & Science, Proddatur. He obtained his B.Tech degree in ECE from Dr. Babasaheb Ambedkar Maratwada University, Aurangabad in 1997, M.Tech degree in Digital Electronics and Computer Science from JNTU, Hyderabad in 2007. He has 11 years of teaching experience. He registered his Ph.D from JNTU, Anantapur in digital Image Processing. His area of interest is Microprocessor and EDC circuits. He has published 3 papers in international journals.