Artificial Intelligence formulated this projection for compatibility purposes from the original article published at Global Journals. However, this technology is currently in beta. *Therefore, kindly ignore odd layouts, missed formulae, text, tables, or figures.*

¹ Energy Efficient FMA for Embedded Multimedia Application

Mandala Rakesh Raj¹ and Ms S. Sujana²

¹ VARDHAMAN COLLEGE OF ENGINEERING

Received: 6 December 2013 Accepted: 2 January 2014 Published: 15 January 2014

6 Abstract

2

3

This article presents energy efficient fused multiplyadd for multimedia applications. Low cost, 7 low power and high performance factors diddle the design of many microprocessors directed to 8 the low-power figuring market. The floating point unit occupies a significant percentage of the a silicon area in a microprocessor due to its wide data bandwidth and the area occupied by the 10 multiply array. The fused floating-point multiply-add unit is utilitarian for digital signal 11 processing (DSP) applications such as fast Fourier transform (FFT) and discrete cosine 12 transform (DCT). The proposed designs are implemented for single precision and synthesized 13 with a 45-nm standard cell library. To improve the performance of the fused floating point 14 multiply-add unit, we are supervening upon leading zero anticipation with the novel leading 15 zero detection, as the novel leading one detection algorithm allowing us to significantly reduce 16 the anticipation failure rates. 17

18

19 Index terms—floating point, binary128, fused multiply add, simd, implementation, computer arithmetic.

20 1 Introduction

his paper presents energy efficient fused multiplyadd unit for multimedia applications. In this, floating point 21 can be implemented by using different precisions. As we have SP (Single Precision), DP (Double Precision), 22 QP (Quadruple Precision). IEEE Binary 32 which is pertained as single precision and binary64 referred to as 23 24 double precision. Floating-point operations have both gained widespread popularity in versatile multimedia and 25 scientific applications, resulting in modern processors patronize both the precisions. Due to accumulation errors in computations they are becoming deficient for today large-scale applications. This precision problem can be 26 overcome by one promising approach that is by using the binary 128 which is referred to as quadruple precision 27 or QP format. The accuracy and numerical stability of many applications can be improved by introducing 28 this format and is already specified in the new IEEE-754-2008 standard. Another approach to which we have 29 to improve the performance is using the fused multiply add (FMA) operation which yields one rounding error 30 for two operations [3]. The first FMA is introduced in the year 1990 by IBM RS/6000 [6], [7]. After FMA is 31 implemented by several companies like HP, MIPS, ARM and Intel. 32

Many algorithms are developed on floating-point fused multiply add unit to decrease its latency [2], [4]. As we can say it is a key feature of the floating-point unit because it greatly increases the floating-point performance and accuracy since rounding is performed only once for the result.

A Field Programmable Gate Array, FPGA provides a versatile and inexpensive way to implement and test VLSI designs. It is mostly used in low volume applications that cannot afford silicon fabrication or designs which require frequent changes or upgrades.

In FPGA's, the bottleneck for designing efficient floating-point units has mostly been area with advancement in FPGA architecture [3], there is a significant increase in FPGA densities so latency has been the main focus of attention in order to improve performance.

The main contribution and objective of our work is to implement the architecture which is proposed by Lang/Bruguera but with little change to facilitate the implementation.

In reminder of this paper is organized as follows. Proposed FMA unit section 2 backgrounds, section 3 proposed methods and section 4 describes its general architecture. Section 5 provides the evaluation results and Section 6 concludes this paper.

47 **2 II.**

48 3 Background

In this paper, the floating-point fused multiplyadd operation A x B + C is implemented for the IEEE floatingpoint format. In this format, a floating-point number X represents the valuep e s f X ? $\times \times$? = ?) 1 (

52 , Where s , f , ? ,e ,and p are integers with s being the sign bit; f the normalized mantissa; ? the radix, 2 for 53 binary; e the biased exponent; and p the biased number.

The fused multiply-add unit gets the input operands A, B, and C with valuesea a sa f A 2 . .) 1 (? = , eb b s5 sb f B 2 . .) 1 (? = and ec c sc f C 2 . .) 1 (? =

and performs the fused multiply-add operation:) 2 . .) 1 (2 . .) 1 ((ec c sc eb ea b a sb sa f f f rnd C B A 57 $? + ? = + \times + ?$

Where the computed fused multiply-add result is rounded and normalized. The FMA architecture proposed before implemented in several floating-point ? Inorder to reduce the latency, the bit inversion and alignment of the significand of A is done in parallel with the multiplication [2]. The bit inversion provides the one's implement of A for an effective subtraction.

? The shift amount of the alignment depends on the value of d = E a - (E b + E c), where E a, E b , E c are the exponents of the A, B, and C operands, respectively.

? When d? 0 (i.e. E a > (E b + E c)), in a conventional alignment, B x C would have to be aligned with a right shift of d bits. vectors that are reduced together with the aligned A using 3:2 CSA, because the product has only 106 bits. The 55 most-significant bits will be sign extension bits, for theses 55 most significant bits, we use two multiplexers, one to select between A and inverted A as a sum vector and the second one to select between zeros and A as a carry vector by Xor-ing sign extension bits then the outputs of the two multiplexers are concatenated at the output of the CSA to obtain the 161-bit sum and carry words. c) The next step is the addition of the carry and sum words and the determination of the shift amount for normalization.

? The carry and sum words, obtained at the output of the CSA, are added in a 161-bit one's complement adder
(with end around carry adjustment for effective subtraction). As the result can be negative, a complementer is
required at the output of the adder.

? In parallel with the significands addition, the normalization shift is determined. The LZA (Leading Zero 74 Anticipator) [12] produces the amount of the shift directly from the operands. d) Once the result of the addition 75 is obtained, the normalization shift can be performed since the shift amount has been already determined. 76 A normalization shift is required to place the most significant bit of the result at bit 0; as a consequence, 77 78 normalization is performed to compensate for the cancellation produced in subtraction as well as to compensate 79 for the way the alignment is performed. With this scheme, the delay of the FMA operation is determined 80 by the sum of the delays of the following components: multiplier, 3-2 CSA, 161-bit adder plus complementer, normalization, and rounding. On the other hand, the main hardware components are: multiplier, alignment 81 shifter, 3:2 CSA, LZA, 161-bit adder, normalization shifter, and rounder. 82

83 **4** III.

⁸⁴ 5 Proposed Architecture

We now describe the proposed FMA architecture. Since the unit is quite complex, we present this description in a
single step. In this section, we give an overview of the scheme, with just enough detail to make it understandable
and believable. Here we use Fig. 2 to illustrate the description.

The objective of the proposed FMA architecture is to reduce the overall delay, and Power. Since, in floatingpoint addition and multiplication, one of the approaches to reduce latency has been to combine addition with rounding [5], [10], [12], [13], we follow the same approach. For this approach, in floating-point addition and multiplication, the order of normalization and rounding is interchanged. This seems impractical to do for FMA because, before the normalization, the rounding position is not known. The solution we explore is to perform the normalization before the addition.

94 6 Architecture

To improve the performance of the fused floating point multiply-add unit, we are supervening upon leading zero anticipation with the novel leading zero detection, as the novel leading one detection algorithm allowing us to

significantly reduce the anticipation failure rates. The proposed leading digit is worked using tree structure, where inputs of n bits are divided into n/2 pairs of bits.

For each pair of bits a two bit count is generated and another counter is triggered to calculate the depth of the tree. For example a four digit can be paired into two pairs and a counter is used to find the one/ zero in pair i and i+1 and the second counter is used to find the value of required bit in which pair and at which level. This method is continued of log 2 (n) levels.

To boost this tree structure we use a structure method which speed's up by a 4 bit or even 8bit to reduce the hierarchy of the tree structure.

105 7 IV.

Detailed Description of some Modules of the Architecture a) 3:2 CSA The multiplier produce 106-bit sum and carry vector that are reduced together with the aligned A using 3:2 CSA. Although the output of the multiplier must be positive number because we multiply two positive numbers (sign and magnitude representation), one of the two output vectors of the multiplier may be negative because of using booth algorithm which use negative sets {-1,-2} which convert a positive number with sign and magnitude representation to a negative number with two's complement representation. The addition of sum and carry vectors must be a positive number but one of them, not both, may be negative.

Instead of using 161-bit CSA, only the 106 leastsignificant bits of the aligned A are needed as input to the 3:2 CSA, because the product of sum and carry vectors has only 106 bits and The 55 most-significant bits will be sign extension bits which have two cases {0, 0} if both sum and carry vectors are positive or {0, 1} if one of them is negative. For the 55 mostsignificant bits, we use two multiplexers, one to select between A and inverted A as a sum vector and the second one to select between zeros and A as a carry vector by Xor-ing sign extension bits then the outputs of the two multiplexers are concatenated at the output of the CSA to obtain the 161-bit

119 sum and carry words.

¹²⁰ 8 b) Leading zero Anticipator

The leading zero anticipator (LZA) has two main parts: the encoding of the leading-one position i.e. detection 121 module and the correction of this position i.e. correction module. The detection module are divided into two 122 parts the first one is called LZA and it determines the number of leading zeros i.e. the position of the leading one. 123 By producing a string of zeros and ones where the position of the most significant 1 is the Year 2014 J position 124 of the leading one. The second part, called leading zero detectors (LZD), counts the number of zero digits from 125 the left-most position until the first nonzero digit is reached i.e. leading one position. Since the detection is done 126 from most significant bit to least significant bit regardless of the carry that may come from the least significant 127 bit, the detection of leading one position may be off by one bit. The LZA logic takes two input strings and uses 128 a set of logical equations given below. 129

After LZA logic LZD is used to drive the normalization shifter [11] by encoding the position of leading one to its weighted binary representation.0), . . (). . (111111>+?+×=+++++iggzztgzgtfi iiiiiiiiii100.ttf = Where iiibat? = iiibag. = iiibaz. =

The LZD unit assumes n bits as input and produces log 2 n bits of the leading one position Pattern Position Valid The IEEE 754 binary floating-point standard defines a set of normalized numbers and four sets of special numbers. Of the four types of special numbers, three do not require computation for arithmetic operations. These include Not-a-Numbers (NaN), infinities, and zeros. De-normalized numbers, also known as subnormal or denormals are the fourth type of special number and do require computation.

Normalized numbers can be described by the following: bias xe xs f X ? $\times \times \times$? = 2 1) 1 (

Where X is the value of the normalized number, X s is the sign bit, X f is the fractional part of the significand, Xe is the exponent, and bias is the bias of the format which corresponds to 127, 1,023, and 16,383, for single, double, and quad, respectively. Denormalized numbers can be described by the following:0, 2 0) 1 (1 ? $\times =$ $\times \times \times ? = ?$ f i X f X e bias xs

The denormal format differs from a normal number in that there is no implied bit and the exponent is not equal to X e -bias, but , instead, is forced up by 1 to E min , which is equal to -126, -1,022, and -16,382, depending on the format.

Using the results from the LZD, the result from the adder is shifted left to normalize the result. That means now the first bit is 1.

The normalize is mostly a large shifter. The shifter has more than one stage. The stages are organized from the coarsest to the finest. The last stage performs a shift by one or two due to correction signal. This should have a negligible effect on the delay of the last stage.

¹⁵¹ 9 d) Rounding

The IEEE 754 floating-point standard has been widely adopted since its introduction in 1985. The standard 152 153 requires that all arithmetic operations are rounded so as to maintain as high a degree of accuracy and uniformity across different computing platforms as possible. The rounding decision is taken by knowing also sticky and round 154 bits. The sticky bit is calculated from the result by OR-ing all least significant bits after the round bit. Rounding 155 operations were originally viewed as a final separate step in most arithmetic circuit implementations. This has 156 been merged with the carrypropagate addition in floating-point adders by delaying normalization until after 157 rounding. Four different rounding modes are laid down by the IEEE floating-point standard [8], [9]: rounding 158 toward 0, rounding to nearest (even), and rounding to \pm . Rounding to nearest (even) is the standard's default 159

mode; rounding toward zero is helpful in many DSP applications; and rounding to \pm is used in interval arithmetic, which affords bounds to be specified on the accuracy of a number.

¹⁶² 10 J e XIV Issue VI Version

163 I V. Simulations Results

¹⁶⁴ **11 b)** Power

The power is the important aspect of the any architecture as the power decreases the power consumption of the entire processor decreases. In this project the power of the both proposed and the previous architecture are calculated using Cadence RC complier in different TSMC standard libraries. The proposed architecture is efficient in terms of power.

169 12 Conclutions

Architecture for a floating-point Multiply-Add-Fused (FMA) unit that reduces the latency of the traditional FMA units has been proposed. The proposed FMA is based on the combination of the final addition and the rounding, by using proposed LZD. This novel leading one detection algorithm allowing us to significantly

reduce the anticipation failure rates. We embedded the proposed technique in Fused floating point multiply and Accumulation unit and its silicon area and performance with other existing solution. This approach has been

used previously to reduce the latency of the floating-point addition and multiplication. However, it can be used

only if the normalization is performed after the rounding and this is not possible for the FMA operation because

the rounding position is not known until the normalization has been performed. To overcome this difficulty, we

¹⁷⁸ propose that the normalization be carried out before the addition. This required a careful design of some other parts of the FMA, the leading-zeros-detector (LZD). ¹



Figure 1: T © 2014

179

 $^{^1 \}ensuremath{\textcircled{O}}$ 2014 Global Journals Inc. (US)



Figure 2: Fig. 1 :



Figure 3: Je



669

3

Figure 4: Fig. 2:



Figure 6: Fig. 3:

"Curso-Beeline"= 100,000ps											TreA = 100 KObs
Nare +	Casa •	R.M.p.	St,ZTIps	198,400,03	a ma	se allips	Sejillips	93,Z0ps	59,411 ps	39,600 ps	59,811pt
म् 🚚 म्हण	1.111.19	1111,110000									
e di ka	10,000 £	1001_0011111									
ទ 🚺 ត្ រព្	1 289, 99	CHARLEN AND									
P. 10 0 8	2.00000	1001_0000000									
E 🗣 g [[0.47]	1.000,00	000,000000									
8 1 10.07	a un no	1111_111111									
pqt 🖗 B	A (18_81) F	481,103848									
5 10 x 10 x 10	1 288,85	Desa Annanan									
0.01.08	J 112 12)	uar intans									
1411	3.03.01	10,000									

Figure 7: Fig. 4 : Fig. 5 : Fig. 6 :

 \mathbf{II}

Methods	Area	Power
Using-nm	(μm^2)	(mW)
Existing	87,908	19
Method(90nm)		
Proposed	$6,\!984$	6.35
Method (45nm)		
VI.		

Figure 8: Table II :

- 180 [Computers (2004)], Computers Aug. 2004. 53 (8) p. .
- 181 [Even and Seidel (2000)] 'A Comparison of Three Rounding Algorithms for IEEE Floating-Point Multiplication'.
- 182 G Even , P M Seidel . *IEEE. Trans. Computers* July 2000. 49 (7) p. .
- [Hinds ()] 'An Enhanced Floating Point Coprocessor for Embedded Signal Processing and Graphics Applications'.
 C Hinds . Proc. 33rd Asilomar Conf. Signals, Systems, and Computers, (33rd Asilomar Conf. Signals, Systems,
- and Computers) 1999. p. .
- 186 [Chen and Cheng ()] 'Architectural Design of a Fast Floating-Point Multiplication-Add Fused Unit Using Signed-
- Digit Addition'. L Chen, J Cheng. Proc. Euromicro Symp. Digital Systems Design, (Euromicro Symp. Digital
 Systems Design) 2001. p. 346.
- 189 [Yeh (1999)] Fast method of floating point multiplication and accumulation, H Yeh. February 1999.
- 190 [Montoye (1990)] Floating Point Unit for Calculating A=XY+Z Having Simultaneous Multiply and Add, E 191 Montoye. November 1990. (United States patent)
- [Lang and Bruguera (2004)] 'Floating-Point Fused Multiply-Add with Reduced Latency'. T Lang , J Bruguera .
 IEEE Transactions on Computers August 2004. 53 (8) p. .
- [Lang and Bruguera] 'Floating-Point Multiply-Add-Fused with Reduced Latency'. T Lang , J D Bruguera . *IEEE Trans*
- [Seidel and Even ()] 'How Many Logic Levels Does Floating-Point Addition Require?'. P M Seidel , G Even .
 Proc. Int'l Conf. Computer Design (ICCD 98), (Int'l Conf. Computer Design (ICCD 98)) 1998. p. .
- [Schmookler and Nowka (2001)] 'Leading Zero Anticipation and Detection –A Comparison of Methods'. M
 Schmookler, K Nowka . Proceedings of the 15th IEEE Symposium on Computer Arithmetic, (the 15th
 IEEE Symposium on Computer ArithmeticVail, Colorado) June 2001. p. .
- [Santoro et al. ()] 'Rounding Algorithms for IEEE Multipliers'. M R Santoro , G Bewick , M A Horowitz . Proc.
 IEEE Ninth Symp. Computer Arithmetic, (IEEE Ninth Symp. Computer Arithmetic) 1989. p. .
- [Santoro et al. (1989)] 'Rounding Algorithms for IEEE Multipliers'. M Santoro , G Bewick , M A Horowitz .
 Proceedings of the 9th IEEE Symposium on Computer Arithmetic, (the 9th IEEE Symposium on Computer ArithmeticSanta Monica, California, USA) September 1989. p. .
- [Montoye et al. (1990)] 'Second -Generation RISC Floating Point with Multiply-Add Fused'. R Montoye , E
 Hokenek , S Runyon . *IEEE journal of solid-state circuits* October 1990. 25 (5) p. .
- 208 [Oberman et al. ()] 'The SNAP Project: Design of Floating-Point Arithmetic Units'. S F Oberman, H Al-Twaijry
- , M J Flynn . Proc. IEEE 13th Symp. Computer Arithmetic, (IEEE 13th Symp. Computer Arithmetic) 1997.
 p. .