

# Implementation of Efficient Architecture for Dual-Mode Floating Point Division

K.Raja Kameswari<sup>1</sup>, Dhanbal.R<sup>2</sup>

<sup>1</sup> M.Tech VLSI, Vellore Institute of Technology, Vellore.

<sup>2</sup>Vellore Institute of Technology, School of Electronics Engineering, Vellore

**Abstract:** - This paper proposed an architecture for the double precision floating point division it works for both single and dual mode. Floating point division is used in various applications like engineering and scientific applications. This proposed architecture works for dual mode which can calculate on two sets of single precision operands and one set of double precision operands in parallel. This architecture works on the principle of series expansion multiplicative methodology for mantissa computation. In this a Radix- 4 Modified Booth Multiplier is used, which is used in dual mode mantissa computation in iterative fashion and also in floating point division the key elements used are leading zero counter, dynamic left and right shifters, rounding etc. The dual mode architecture which is proposed in this paper is synthesized in Synopsys. Compared to the other double precision division architectures. The proposed architecture performs better in factors like area, time period and throughput.

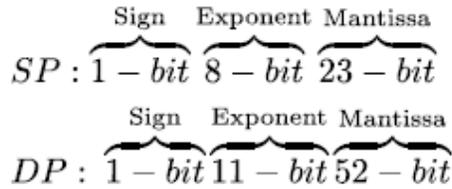
**Index terms-** Floating point Division, Dual –Mode functionality, Modified Booth Algorithm.

## I. INTRODUCTION

Floating point arithmetic is mostly used in many scientific, engineering and signal processing applications its large vital range and suitable platform for designers to register their algorithms. Moreover executing floating point arithmetic numbers is very difficult. The arithmetic operations like addition, subtraction, multiplication and division. Among these division is generally more challenging but division is mostly used in scientific and engineering applications. So, it is the need to implement efficient architecture for floating point. Generally dividers are designed to use in the iteration and latency is more compared to other arithmetic operations. Floating point division [1] is the arithmetic core used in scientific and engineering applications. The hardware complexity of the division in floating point is high compared to other operations like addition, subtraction and multiplication and also division architecture require more area and compared to be low performance. So, based on the area parameter we aim to unified and configurable architecture is proposed. In this paper we going to propose an architecture which computes both single and double precision vectors used for vector processing as an alternative different vectors of single and double precision, it can form a configurable floating point block, each configurable floating point arithmetic blocks can be used as double precision or two parallel single precision. Because of this block we can improve area and gives the best performance. The proposed architecture is based on multiplicative division algorithm series expansion methodology [2]algorithm. this one of the efficient methods in multiplicative division method like Goldschmidt and Newton raphson methods[10]. The implementation may do in different methods like digital recurrence, multiplication based and other

techniques. Coming to the digital recurrence algorithm SRT is the iterative form it takes more implementation time and coming to multiplication based upon approximation and inverse of division in iteration fashion. series expansion algorithm is based on computing initial quotient and then computing the quotient by using remainder and add the quotients. this expansion method will give the best performance, less area, lower latency and also compared to above methods this method require less memory requirement and faster compared to the digital recurrence method. In this architecture dual mode, Modified booth multiplier Radix 4 in the use of mantissa division, which has less area and performance is high over single-mode multiplier. since the underlying integer multiplier in mantissa division unit. an architecture iterative proposed by using 1-stage integer multiplier to get area efficiency. This paper proposed an architecture for the division which can be either a parallel two single precision and double precision and this architecture is designed based on the series expansion method. this architecture DpdSP division architecture supports subnormal and normal operands with round-to-nearest rounding method. a design only normal support has also implemented and compared with proposed architecture implemented designs like infinity, divide-by-zero, zero. proposed architecture is designed for both normal and subnormal which is helpful for rounding with round-to-nearest rounding for both single and double precision results. All the major blocks like dynamic right and left shifters, leading one detection and mantissa division are designed for efficient purpose. a single mode double precision based on the computational flow. A Floating point arithmetic computation involves separating the sign, exponent and mantissa part of the operands and finally combine them after normalization and rounding. the IEEE

standard format for floating point numbers[6] in the proposed architecture are single and double precision numbers are as follows



### II.BACKGROUND

The algorithm for floating point division architecture is shown below

- 1: ( $IN1$  (Dividend),  $IN2$  (Divisor)) Input Operands;
- 2: **Data Extraction & Exceptional Check-up:**  
 $\{S1(\text{Sign}1), E1(\text{Exponent}1), M1(\text{Mantissa}1)\} \leftarrow IN1$   
 $\{S2, E2, M2\} \leftarrow IN2$   
 Check for Infinity, Sub-Normal, Zero, Divide-By-Zero
- 3: **Process both Mantissa for Sub-Normal:**  
 Leading One Detection of both Mantissa ( $\rightarrow L\_Shift1, L\_Shift2$ )  
 Dynamic Left Shifting of both Mantissa
- 4: **Sign, Exponent & Right-Shift-Amount Computation:**  
 $S \leftarrow S1 \oplus S2$   
 $E \leftarrow (E1 - L\_Shift1) - (E2 - L\_Shift2) + BIAS$   
 $R\_Shift\_Amount \leftarrow (E2 - L\_Shift2) - (E1 - L\_Shift1) - BIAS$
- 5: **Mantissa Computation:**  $M \leftarrow M1/M2$
- 6: **Dynamic Right Shifting of Quotient Mantissa**
- 7: **Normalization & Rounding:**  
 Determine Correct Rounding Position  
 Compute ULP using Guard, Round & Sticky Bit  
 $M \leftarrow M + ULP$   
 1-bit Right Shift Mantissa in Case of Mantissa Overflow  
 Update Exponent
- 8: **Finalizing Output:**  
 Determine STATUS signal & Resolve Exceptional Cases  
 Determine Final Output

This algorithm is helpful to implement for both subnormal and normal processing .algorithm includes case handling and processing. Floating point arithmetic implementation computing the sign, exponent and mantissa parts of the operands and combined them after normalization and rounding.

### III. DDPSP DIVISION ARCHITECTURE

The proposed architecture is shown in Fig1.it consists of three pipelining stages. In the architecture each stage is discussed coming to the operands there are two 64 –bit operands ,one dividend(in1) and divisor (in2)are the initial inputs along with a control signal is available dp\_sp(dual precision or dual single precision).both the operands consists of two parallel SP

operands (two sets of 32-bit)and DP operands(64-bit) is displayed in Fig2. Algorithm is divided into two stages. Coming to first stage of architecture is contains stage 2 and stage 3 of algorithm. It includes data extraction, sub normal processing and exceptional case handling[3] and also includes mantissa division.

The data extraction computation is shown in Fig 3 .primary operands and extract the signs ( $sp1\_s1,sp1\_s2,sp2\_s1,sp2\_s2,dp\_s1,dp\_s2$ ) and exponents ( $sp1\_e1,sp1\_e2,sp2\_e1,sp2\_e2,dp\_e1$ and  $dp\_e2$ )and mantissa( $sp1\_m1,sp\_m2,sp2\_m1,sp2\_m2,dp\_e1$ and $dp\_e2$ ). sizes of the above components is shown in Fig2

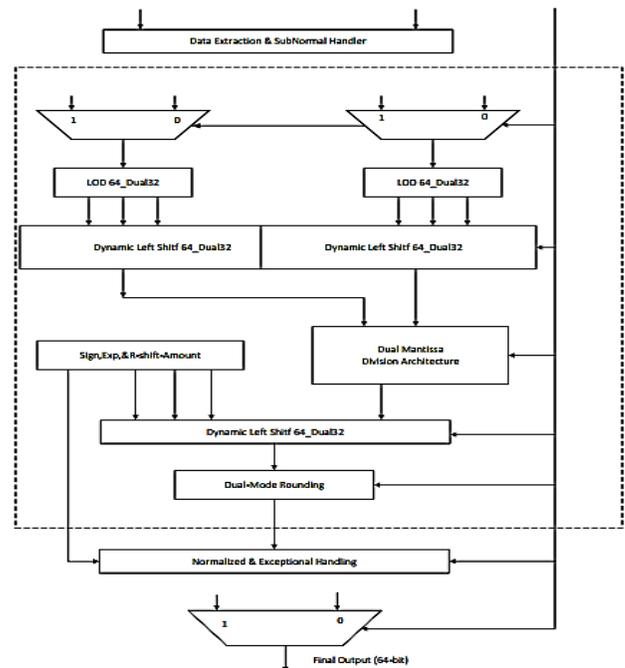


Fig1:Double precision with dual single precision Division architecture

s:sign,\_e:exponent,\_m:mantissa,dp\_sp:mode(DP/Dual-sp)dp:double precision,sp:single precision,\_ls:leftshift,\_rs:right shift.

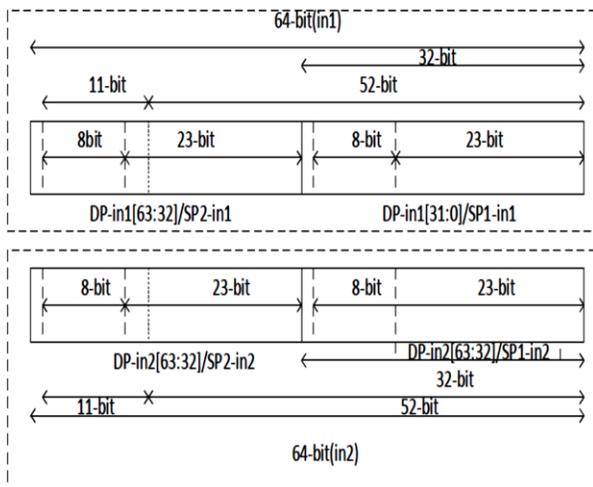


Fig2:Dpdsp Input Output Format

The input and output encoding is taken under the format of IEEE standard format. The subnormal handling and exceptional check computation checks subnormal, NaN (Not a Number) and infinity. It also checks for division by zero and has to be shared among both SP and DP. After data extraction and subnormal and exceptional handling, the uni-mantissa (M1 and M2) is obtained by using a mux as shown in the architecture. Because of this uni-mantissa is helpful in further blocks which is helpful in efficient resource sharing. The next two blocks are leading one detector and dynamic shifter. The dynamic shifter helps to convert the subnormal format to normal format. Firstly, it provides the amount of left shift amount by using the block leading zero counter and the amount given by the LOD and then dynamic left shifter shifts the mantissa. The architecture of the LOD is designed in a fashion that firstly it starts with a basic block of 2:1 LOD and finally 64:6 LOD which has 32:5 LOD. The initial block consists of basic gates like AND, OR, and NOT gates which makes for both SP's and DP's. The final result is cost-minimal compared to DP LOD. The uni-mantissas are given to the Dynamic Left Shifter (dual mode). The shifter shifts the amount of left shifting given by the LOD (as in Fig 1). In the processing of DP, the SP's left shift amounts are kept to be zero. Similarly, during SPs DP left shift amounts are kept to zero. The left shifter dual mode is shown in [8]. This shifter consists of 6 stage barrel shifter in which 5 are dual mode and the first shifter is single mode. In the  $n$ th stage, it is only used to perform double precision coming to dual mode. The stage consists of multiplexers, which shift the bits corresponding to both SP and DP. The main block in the DPdSP division architecture is the Mantissa division architecture for dual mode. After shifting the mantissa is converted into normalized form  $m1, m2$ . The normalized inputs are given to two LUT tables available in the mantissa architecture. In these Look Up Tables [7], one is used for both DP and SP2 and the other look-up table is used for SP1 only, as shown in Fig. 3.

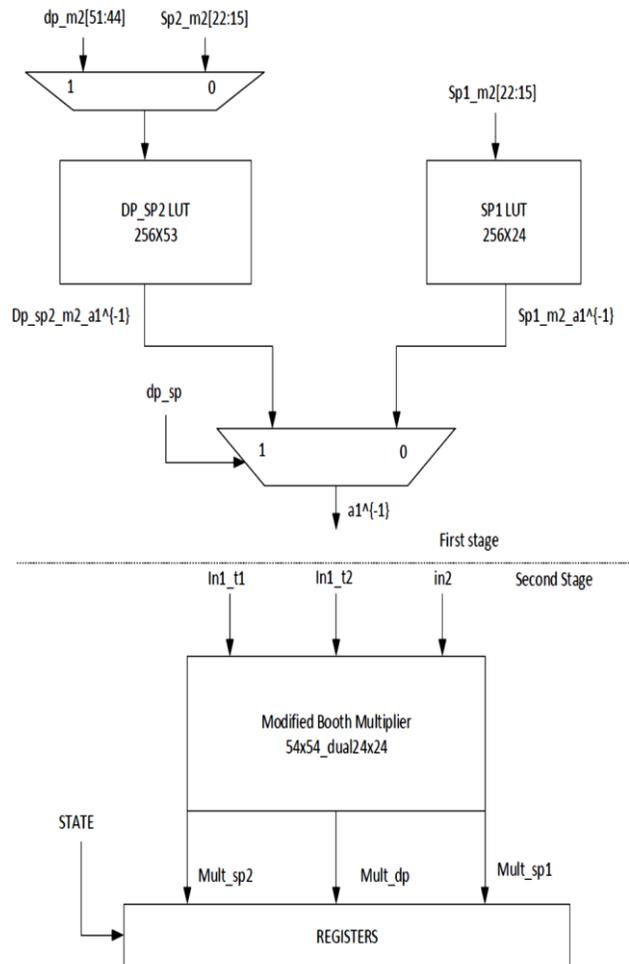


Fig.3.DpdSP Dualmode Mantissa architecture for division

In the second pipelining stage (4 & 5 stages) it consists of sign, exponent, and right shift amount, and the result is helpful for the remaining stages. The computed sign bits are a simple XOR operation. Mantissa generation is a very difficult part in division arithmetic. The mantissa obtained from the first stage of the mantissa division architecture is given to the modified booth multiplier. For the effective inputs to the multiplier, a Dual mode FSM is used. The booth multiplier is shown in the below figure, which is used to compute in dual mode. The booth multiplier is helpful to reduce partial products and for DP processing. Here, it is a 54-bit integer multiplier, and the DADDA tree has 8 levels, which compress the entire partial products into two operands and further is given to a parallel prefix adder. The final product contains both dual SP and DP results in Fig. 4.

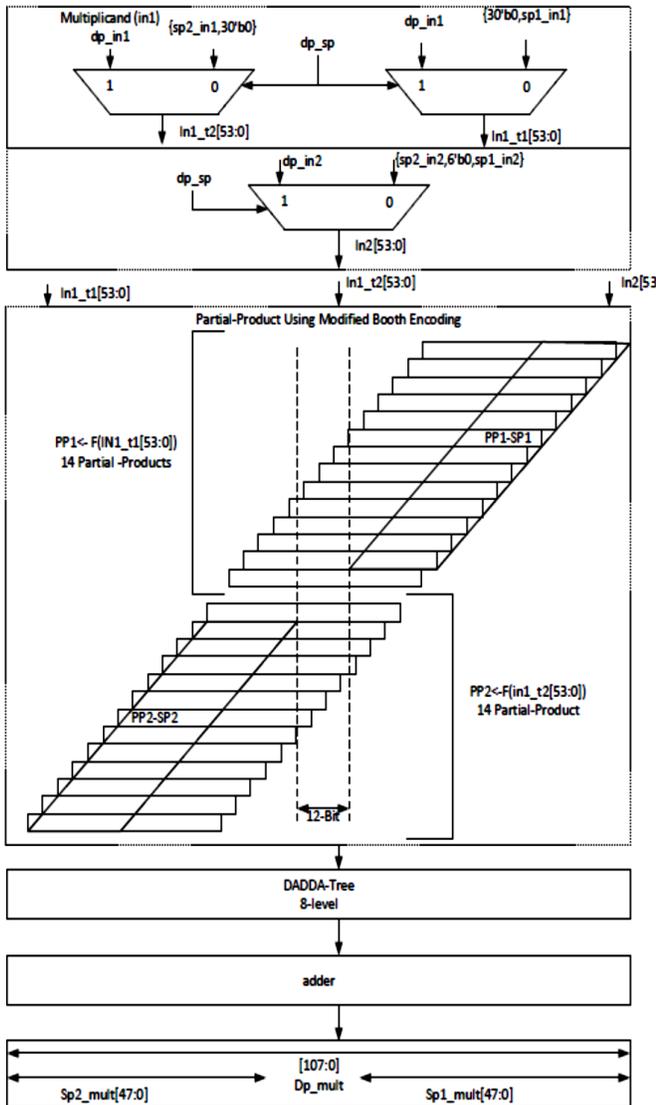


Fig.4. Dual-Mode Modified Booth Multiplier Architecture

Coming to the third stage architecture consists of 6, 7, 8 stages in algorithm in this case the exponent underflow and mantissa quotient is given as the input to the dynamic right shifter or dual mode and then rounding is done which undergoes normalisation and exceptional processing. To the dynamic right shifter mantissa quotient and right shift amount is given as the inputs. right shifter is similar to dual mode dynamic left shifter also has 6 stages 5 stages are dual mode and initial stage is single mode. top muxs provide SP quotient and third mux combined and give DP quotient to produce right shifting. The proposed dual mode rounding [3] produces round to the nearest method and it contains two steps ULP and ULP addition (unit-at-last-place) with mantissa quotients. ULP is based on the guard, round and sticky bits. ULP is done for dual modes for DP and both SPs. mantissa quotient contains both SPs and DP quotients. ULP additions consist of 32-bit

incrementer. the rounded quotient obtained from dual mode rounding is further normalized into SPs and DP. if there is any overflow due to the unit-at-last place addition we require 1 bit right shift for overflow. Finally, the normalized results for double precision and two parallel single precision are multiplexed by the MUX(2:1) to produce 64-bit output which contains quotients of single and double precision.

### III. RESULTS



Fig.5. Simulation results of DPdSP architecture

	[4]	[5]	Proposed Multiplier	
Sub normal	X	✓	✓	X
Technology	90nm	90nm	90nm	90nm
Gate Count (based on minimum size inverter)	212854	118085	66416	62649
Period (in ns)	1.41	19	1.72	1.67
Power (mW)	-	11.92	30.9	33.45
Throughput (in Clock Cycles)	29/15	1/1	10/8	

Tabel.1. Comparison table of different divider architectures

The simulation results of DPdSP architecture is shown in Fig.5. and Area, Power results of the proposed and related architecture are given below in Tabel.1. By the results we can understand that the DpdSP architecture is area efficient and produce less power and through put compared to the remaining architectures

### IV. CONCLUSION

This paper proposed dual mode double precision architecture for floating point division arithmetic. It is used for both dual single and double precision operands. This architecture is made based on series expansion mantissa division. A dual mode modified radix 4 booth multiplier is used which will help in low overhead. Entire architecture is used to perform dual mode computation and less hardware. By the DPdSP architecture for dual mode division we can use other multiplicative algorithms like Goldschmidt and Newton Raphson. This is the future work.

#### REFERENCES

- [1] J.-C. Jeong, W.-C. Park, W. Jeong, T.-D. Han, and M.-K. Lee, "A cost-effective pipelined divider with a small lookup table," *IEEE Trans. Comput.*, vol. 53, no. 4, pp. 489–495, Apr. 2004.
- [2] M.K. Jaiswal, M. Balakrishnan and Paul, "Series expansion based efficient architectures for double precision floating point division," *Circuits, Signal Processing*, vol. 33, no. 11, pp. 3499-3526, 2014.
- [3] Manish Kumar Jaiswal and Hayden K. "Area Efficient architecture for Dual-Mode Double Precision Floating Point Division," *IEEE Transactions on Circuits and Systems*, VOL. 64, Feb 2017
- [4] A. Isseven and A. Akka, "A dual mode quadruple floating point divider," in *proc. 2006*, pp. 1697-1701
- [5] M. Jaiswal, R. Cheung, "configurable architecture for double/two parallel single precision floating point division," in *proc. IEEE Computers, VLSI (ISVLSI)*, Jul. 2004, pp. 332-337
- [6] IEEE Standard for Floating point Arithmetic, IEEE Standard 754-2008, Aug. 2008, pp. 1-70.
- [7] B. Pasca, "Correctly rounded floating point division for DSP enabled FPGAs," in *proc conf Field program*, Aug. 2012, pp. 249-254.
- [8] M. Jaiswal, B. Varma, M. Balakrishnan, K. Paul and "Configurable architectures for multi mode floating point adders" VOL. 62, Aug 2015.
- [9] A. Akka, S, "Dual mode quadruple precision floating point adder" in *Proc. Euromicro symposium, Digital Systems Design*, 2006, pp. 211-220.
- [10] J.M. Muller "Handbook for floating point arithmetic, 1st ed. Basel, 2009.