

## Performance analysis of Floating point adder using Sequential Processing on Reconfigurable hardware

Meenu Talwar\*, Karan Gumber\*\*, Sharmelee Thangjam

\*(M. Tech (CSE), Department of Computer Science, CGC landran, PTU)

\*\* (M. E (ECE), Department of Electronics and communication, Panjab University, Chandigarh)

\*\*\* (Assistant Professor, Department of Electronics and communication, Panjab University, Chandigarh)

### ABSTRACT

The main objective of implementation of floating point adder using sequential processing on the reconfigurable hardware i.e. on FPGAs is to utilize less chip area, less number of destination paths/ports, less clock period, less combinational delay and faster speed. Less chip area means less number of slices is used in reconfigurable hardware i.e. on FPGAs. Less combinational delay means less latency i.e. less time is required to appear an output after the input response is applied and if there is less latency then there will be the faster speed. If there is less number of components are used on FPGAs then less number of paths are used to connect them. Floating point adder implementation on FPGA utilizes 349 slices with a combinational delay of 69.987 nsec consuming 187244 Kbytes of memory with 1 Global clock.

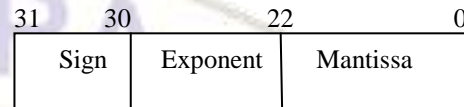
*Keywords* – Floating point adder, Xilinx, FPGAs.

### I. INTRODUCTION

Many scientific problems require floating point arithmetic with high level of accuracy in their calculations. Therefore VHDL programming for IEEE single precision floating point adder in have been explored. For implementation of floating point adder using Sequential processing on FPGAs module various parameters i.e. clock period, latency, area (number of slices used), total number of paths/ destination ports, combinational delay, modeling formats etc will be outline in the synthesis report. VHDL code for floating point adder is written in Xilinx 8.1i and its synthesis report is shown in Design process of Xilinx which will outline various parameters like number of slices used, number of slice flip flop used, number of 4 input LUTs, number of bonded IOBs, number of global CLKs. Floating point addition is most widely used operation in DSP/Math processors, Robots, Air traffic controller, Digital computers because of its raising application the main emphasis is on the implementation of floating point adder effectively such that it uses less chip area with more clock speed[1][2].

### II. FLOATING POINT FORMAT

Floating point number is composed of three fields and can be of 16, 18, 32 and 64 bit. Figure shows the IEEE standard for floating point numbers consists of 32 bits [3] [4]. 1 bit sign of signifies whether the number is positive or negative. '0' indicate positive number whether '1' indicates negative positive number. 8 bit exponent provides the exponent range from E



(min) = -126 to E (max) = 127. 23 bit mantissa signifies the fractional part of a number the mantissa must not be confused with the significand. The leading '1' in the significand is made implicit [2].

#### 2.1 Conversion of Decimal to Floating numbers

Conversion of Decimal to Floating point 32 bit format is explained with example. Let us take an example of a decimal number that how could it will be converted into floating format. Enter a decimal number suppose 129.85 before converting into floating format this number is converted into binary value which is 1000001.110111. After conversion move the radix point to the left such that there will be only one bit which is left of the radix point and this bit must be 1 this bit is known as hidden bit and also made above number of 24 bit including hidden bit which is always '1' like 1.0000011101110000000000 the number which is after the radix point is called mantissa which is of 23 bits and the whole number is called significand which is of 24 bits. Count the number of times the radix point is shifted say 'x'. But in above case there is 7 times shifting of radix point to the left. This value must be added to 127 to get the exponent value i.e. original exponent value is 127 + 'x'. In above case exponent is 127 + 7 = 134 which is 1000110. Sign bit i.e. MSB is '0' because number is +ve. Now assemble result into 32 bit format which is sign, exponent, mantissa. 010001100000001110111000000000. Now take another example which is totally different from above let us enter a decimal number -0.5 which is converted into binary value which is .000011. After conversion move the radix point to the right in this case such that there will be only one bit which is left of the radix point and this bit must be 1 this bit is known as hidden bit and also made above number of 24 bit including hidden bit which is always '1' 1.100000000000000000000000 the number which is after the radix point is called mantissa which is of 23 bits and the whole number is called significand which is of 24 bits. Count the number of times the radix point is shifted to the right say 'x'. In this case there is 5 times shifting of radix point to the right. This value must be subtracted to 127 to get the exponent value i.e. original exponent value is 127 - 'x'. In above case exponent is 127 - 5 = 122 which is 01111010. Sign bit i.e. MSB is '1' because number is -ve. Now assemble result into 32 bit format which is sign, exponent, mantissa 10111101010000000000000000000000.

#Note: - Hidden bit is not included into 32 bit format this bit is implicit. When performing operation with this format this implicit bit is made explicit [5].

### III. ADDITION ALGORITHM FOR FLOATING POINT NUMBERS

VHDL coding for floating point adder is done in Xilinx 8.1i and simulation waveform is shown in Model Sim. Synthesis report in Xilinx 8.1i will tell us all the information about the project summary, Device utilization summary (estimated values) and project status. The floating point addition is the most complex operation then the floating point multiplication since the alignment of mantissa is required before mantissa addition. I would like to explain floating point addition algorithm in 2 cases with example. Case I is when both the numbers are of same sign i.e. when both the numbers are either +ve or -ve means the MSB of both the numbers are either 1 or 0. Case II when both the numbers are of different sign i.e. when one number is +ve and other number is -ve means the MSB of one number is 1 and other is 0. The flowchart of the algorithm is given below in next page and it is explained in following steps with proper example.

#### A. Case I: - When both numbers are of same sign

Step 1:- Enter two numbers N1 and N2. E1, S1 and E2, S2 represent exponent and significand of N1 and N2.  
 Step 2:- Is E1 or E2 = '0'. If yes set hidden bit of N1 or N2 is zero. If not then check is E2 > E1 if yes swap N1 and N2 now contents of N2 in N1 and N1 in N2 and if E1 > E2 make contents of N1 and N2 same there is no need to swap.  
 Step 3:- Calculate difference in exponents d=E1-E2. If d = '0' then there is no need of shifting the significand and if d is more than '0' say 'y' then shift S2 to the right by an amount 'y' and fill the left most bits by zero. Shifting is done with hidden bit.  
 Step 4:- Amount of shifting i.e. 'y' is added to exponent of N2 value. New exponent value of E2= previous E2 + 'y'. Now result is in normalize form because E1 = E2.  
 Step 5:- Is N1 and N2 have different sign 'no'. In this case N1 and N2 have same sign.  
 Step 6:- Add the significands of 24 bits each including hidden bit S=S1+S2.  
 Step 7:- Is there is carry out in significand addition. If yes then add '1' to the exponent value of either E1 or new E2 and shift the overall result of significand addition to the right by one by making MSB of S is '1' and dropping LSB of significand.  
 Step 8:- If there is no carry out in step 6 then previous exponent is the real exponent.  
 Step 9:- Sign of the result i.e. MSB = MSB of either N1 or N2.  
 Step 10:- Assemble result into 32 bit format excluding 24<sup>th</sup> bit of significand i.e. hidden bit [6][7].

**Example** Step 1: Enter N1 and N2.

N1=2.3=0 1000000 100100100000000000000000  
 N2=7.4=0 1000001 111011000000000000000000  
 E1= 10000000  
 E2= 10000001  
 S1=100100100000000000000000

S2= 111011000000000000000000

Step 2: If E2>E1. Yes then swap N1 & N2.

New N1=0 1000010 111011000000000000000000

New N2=0 1000000 100100100000000000000000

Step 3: Calculate d =E1-E2.

10000001-10000000 = 1

Step 4: Shifting of S2 to the right by one and also add 1 to E2.

N2=0 1000000 100100100000000000000000(original)

N2= 0 1000000 010010010000000000000000 (one time shifted)

Shifting by 1 time means add '1' to exponent.

Step 5: New exponent value E2 = 10000001, new significand value S2 = 010010010000000000000000 here E1 = E2.

Step 6: S=S1+S2.

S1=111011000000000000000000

S2=010010010000000000000000

---

S=100110101000000000000000

Step 7: Here is carry out add '1' to exponent and shift result to the right by one bit and discard the LSB of 'S'.

Original exponent=10000010

Original significand=100110101000000000000000

Step 8: MSB of result is '0'.

Step 9: Assemble into 32 bit format.

0 1000010 001101010000000000000000

#### B. Case II: - When both numbers are of different sign

Step 1, 2, 3 & 4 are same as done in case I.

Step 5:- Is N1 and N2 have different sign 'Yes'.

Step 6:- Take 2's complement of S2 and then add it to S1 i.e. S=S1+2's complement of S2.

Step 7:- Is there is carry out in significand addition. If yes then discard the carry and also shift the result to left until there is '1' in MSB also counts the amount of shifting say 'z'.

Step 8:- Subtract 'z' from exponent value either from E1 or E2. Now the original exponent is E1-'z'. Also append the 'z' amount of zeros at LSB.

Step 9:- If there is no carry out in step 6 then MSB must be '1' and in this case simply replace 'S' by 2's complement.

Step 10:- Sign of the result i.e. MSB = Sign of the larger number either MSB of N1 or it can be MSB of N2.

Step 10:- Assemble result into 32 bit format excluding 24<sup>th</sup> bit of significand i.e. hidden bit [6][7].

#### Example

Step 1: Enter N1 and N2.

N1=128.5=0 1000110 100000010000000000000000

N2=-18.25=1 1000011 100100100000000000000000

E1=1000110

E2=1000011

S1=100000010000000000000000

S2=100100100000000000000000

Step 2:  $E1 > E2$  no need to swap.  
 Step 3: Calculate 'd'= $E1 - E2$ .

10000110-10000011=00000011= $\Rightarrow 3$  in decimal.

Step 4: Shifting of S2 to the right by three and also add 3 to E2.

$N2 = 0\ 10000011\ 100100100000000000000000$  (original)

$N2 = 0\ 10000100\ 010010010000000000000000$  (1 time shifting)

$N2 = 0\ 10000101\ 001001001000000000000000$  (2 time shifting)

$N2 = 0\ 10000110\ 000100100100000000000000$  (3 time shifting)

Shifting by 1 time means add '1' to exponent.

Step 5: New exponent value  $E1 = 10000110$ , new significand value  $S2 = 000100100100000000000000$  here  $E1 = E2$  i.e. result is in normalized form.

Step 6: Take 2's complement of S2 because S2 is -ve i.e.  $S = S1 + 2$ 's complement of S2.

$S1 = 100000001000000000000000$

$S2 = 111011011100000000000000$

---

$S = 101101110010000000000000$

Step 7: Here is carry out add lets discard the carry and shift result to the left by one bit to make MSB '1' and then subtract the amount of shifting from E1 or E2 to form original exponent of result.

Original exponent= $10000110 - 1 = 10000101$

Original significand= $110111001000000000000000$

Step 8:- Sign bit of result i.e. MSB= Sign of 128.5 which is larger number.

Step 9: Assemble into 32 bit format.

0 10000101 110111001000000000000000

**C. Special Conditions**

There are some special conditions while implementing floating point adder which needs to be handle these are explained below

- 1: If  $N1 = N2 = '0'$  then overall result is '0'.
- 2: If  $E1 = E2$  and sign bit of  $E1 \neq E2$  then again overall result is '0'.
- 3: If  $E1 = '0'$  and  $E2 \neq '0'$  then overall result is equal to  $E2$ .
- 4: If  $E2 = '0'$  and  $E1 \neq '0'$  then overall result is equal to  $E1$
- 5: If  $d = E1 - E2 \geq 24$  then overall result is larger of  $E1$  or  $E2$  [3].

**D. Problems associated in addition**

There are two problems which occurs when we are going to add two floating point numbers

- 1: When the exponent of two numbers are different this can be solved by shifting the significand of smaller number to the right by an amount equal to exponent difference and this amount is added to exponent value of smaller number to make exponent of both the numbers are same means in normalized form
- 2: When there is carry out in significand addition if both the number are of different sign then add '1' to the exponent

and shift the result of significand to the right by one discarding LSB and if both the number are of different sign then discard the carry and shift the result to the left until there is '1' at MSB the amount of shifting is subtracted from exponent to form real exponent [7][8].

**IV. SYNTHESIS REPORT**

Parameters	Sequential Processing
Number of Slices	349 (7% utilization)
Number of GCLKs	1(4% utilization)
Combinational Delay	69.987nsec
Total number of paths/Destination ports	33/33
Memory	187244 Kbytes
Flip Flop/Latches	103
Clock Buffers	1
I/O Buffers	99
Global Fan out	500
Input Format	Mixed
Comparator	4(1-24 bit, 3-8 bit)
Xor	133
Cell Usage	749

**V. SIMULATION WAVEFORM(USING MODEL SIM)**

Input  
 $1 = 7.5 = 0100000011111100000000000000000000$

Input  
 $2 = 9.25 = 0100000101001010000000000000000000$

/tpa_procedure/n1	UUUUUUUUUUUU	0100000111100000000000000000
/tpa_procedure/n2	UUUUUUUUUUUU	01000001010010100000000000000000
/tpa_procedure/sum	UUUUUUUUUUUU	01000001110000110000000000000000
/tpa_procedure/sub_e	00000000	10000010
/tpa_procedure/c_tem	0	
/tpa_procedure/shift_c	0	1
/tpa_procedure/num2	0000000000000000	01000001111000000000000000000000
/tpa_procedure/s33	0000000000000000	00001100000000000000000000000000
/tpa_procedure/s2_tem	0000000000000000	10000110000000000000000000000000
/tpa_procedure/diff	00000000	00000000

**REFERENCES**

- [1] Ali malik, Dongdong chen and Soek bum ko, "Design tradeoff analysis of floating point adders in FPGAs," Can. J. elect. Comput. Eng., ©2008IEEE.
- [2] Loucas Louca, Todd A cook and William H. Johnson, "Implementation of IEEE single precision floating point addition and multiplication on FPGAs," ©1996 IEEE.
- [3] Alexandru, Mircea, Lucian and Oana, "Exploiting parallelism in double path adder structure for increase throughput of floating point addition," ©2007 IEEE.
- [4] V. Y. Gorshtein, A. I Grushin, S>R Shevtsov, "Floating point addition method and apparatus," Sun microsystem U.S patent 5808926,1998.
- [5] IEEE std. 1076-2002, "IEEE standard VHDL reference manual," Sponsored by Design Automation standards Committee published by IEEE.
- [6] Metin Mete, Mustafa Gok, "A multiprecision floating point adder," ©2011 IEEE.
- [7] Florent de Dinechin, "Pipelined FPGA adders," ©2010 IEEE.
- [8] Ali malik, Soek bum ko, "Effective implementation of floating point adder using pipelined LOP in FPGAs," ©2010 IEEE.

Desired

Result=16.75=0100000111000011000000000000000000

Simulation Result=0100000111000011000000000000000000

Input 1=-

3.5=11000000011100000000000000000000000000

Input 2=-

120.75=110000101111100011000000000000000000

/tpa_procedure/n1	UUUUUUUUUUUU	11000001111000000000000000000000
/tpa_procedure/n2	UUUUUUUUUUUU	11000101111100011000000000000000000000
/tpa_procedure/sum	UUUUUUUUUUUU	11000101111100011000000000000000000000
/tpa_procedure/sub_e	00000000	10000101
/tpa_procedure/c_tem	0	
/tpa_procedure/shift_c	0	5
/tpa_procedure/num2	0000000000000000	11000001111000000000000000000000000000
/tpa_procedure/s33	0000000000000000	00000000000000000000000000000000000000
/tpa_procedure/s2_tem	0000000000000000	00000000000000000000000000000000000000
/tpa_procedure/diff	00000000	00000000

Desired Result=-

124.25=110000101111110001000000000000000000

Simulation

Result=110000101111110001000000000000000000

**ACKNOWLEDGEMENTS**

I would like to thanks the anonymous users for their insightful comments.