

1 Hardware Implementation of a Neural Network node

1.1 Milestone 1: Multiplier and Accumulator (MAC) Module Design

A MAC module has been designed in this step that computes following function : $y = \sum_{i=1}^{16} w_i \cdot x_i$

The schematic of this module is shown in figure 1:

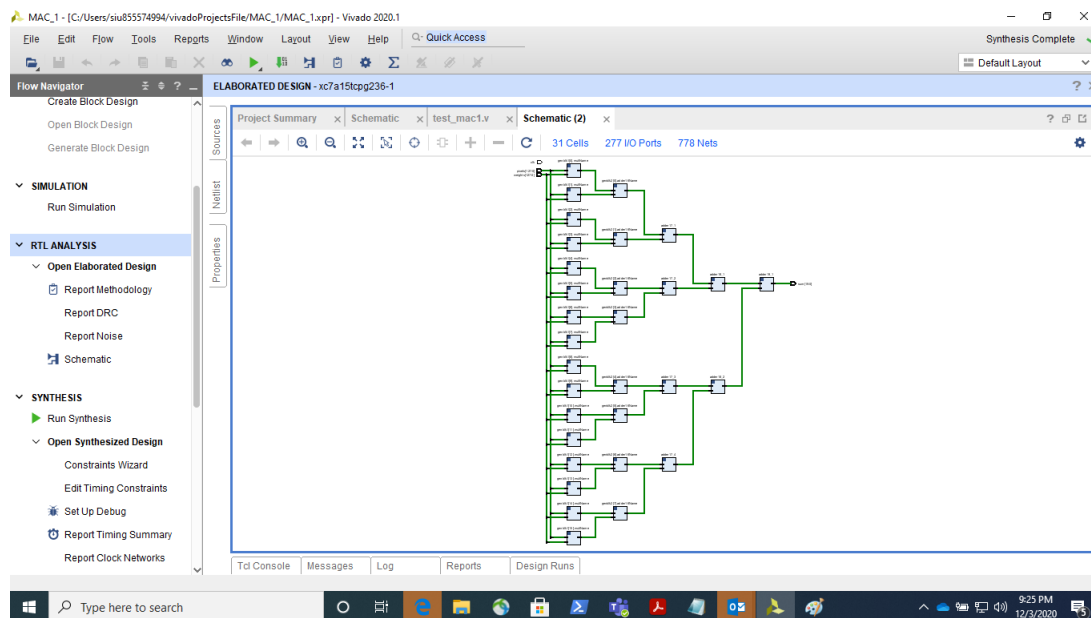


Figure 1: The schematic for MAC module

This MAC consist of 16 multiplier 8 bit multiplier which outputs is 16 bit. The module also consist of 8 16 bit adder, 4 17 bit adder, 2 18 bit adder and 1 19 bit adder. The Module's output is 20 bit and it takes $16 \times 8 = 128$ bit weight and 128 bit pixel value. The 128 bit inputs are distributed into 16, 8 bit chunks and multiplied with 16 multiplier.

The behavioral and after synthesize, obtained waveform are given in figure 2 and 3.

There are some timing glitch observed in the output when I did post synthesis timing simulation. This is due to datapath variation for different elements of the module.

One important observation I found is if I implement pipeline (include pipeline registers in the data-path to make circuit faster in terms of supported clock frequency), the output of the module is wrong for the first instance during the synthesis phase of the module. This is likely related to initialization timing used in vivado environment. If I get rid of the pipeline registers, the output is as expected. I was implementing 3 stage of pipelines.

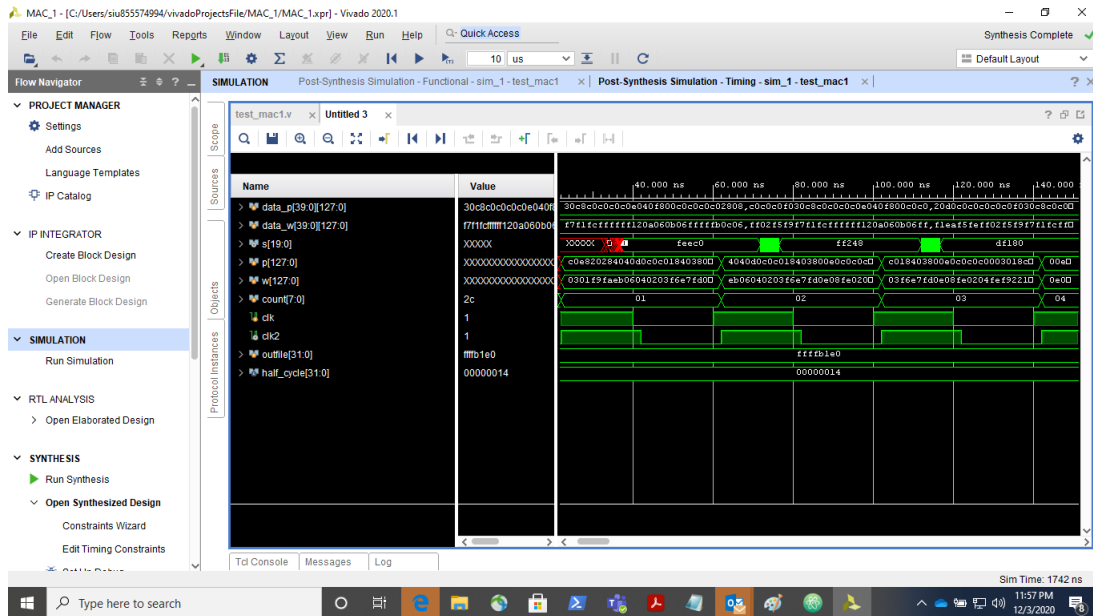


Figure 2: Behavioral waveform obtained in MAC module

Top module for this part in mac1.v and testbench is mac2_tb.v

```

1
2 module multiplier(
3     input signed [7:0] a,
4     input signed [7:0] b,
5     output signed [15:0] p
6 );
7     assign p=a*b;
8 endmodule

```

```

1
2 module adder #(
3     parameter WIDTH=16,
4     parameter A_WIDTH=WIDTH,
5     parameter B_WIDTH=WIDTH)
6 (
7     input signed [A_WIDTH-1:0] a,
8     input signed [B_WIDTH-1:0] b,
9     output signed [WIDTH:0] s
10 );
11
12     assign s=a+b;
13
14 endmodule

```

```

1
2 module mac1(
3     input [127:0] pixels ,
4     input [127:0] weights ,
5     output [19:0] sum);
6
7     // reg [127:0] pixels;

```

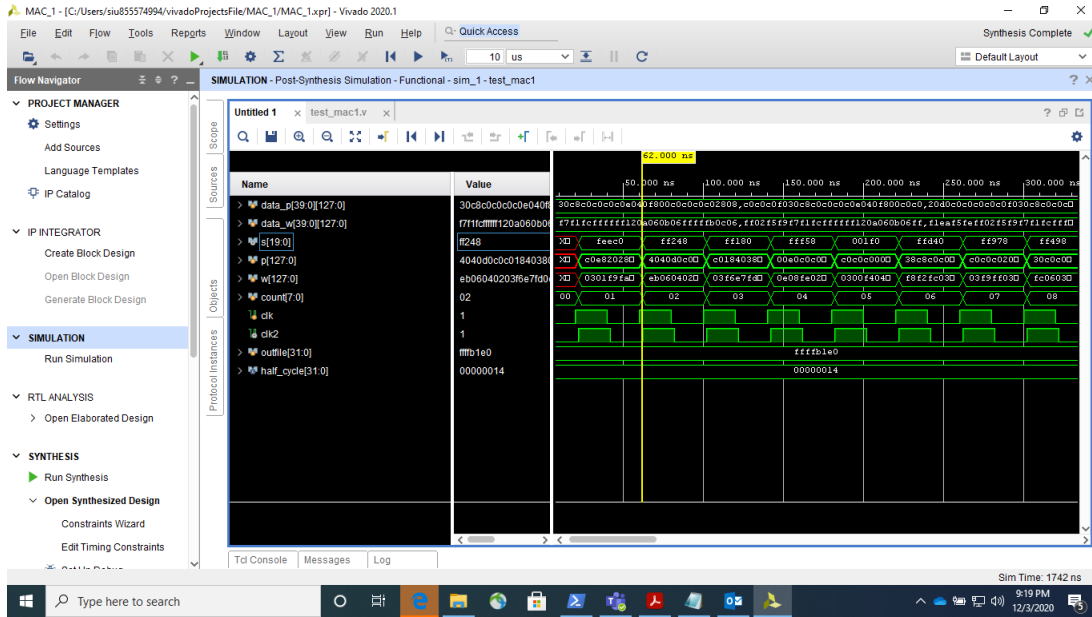


Figure 3: waveform obtained after synthesis for MAC module

```

8 // reg [127:0] weights;
9 // reg [19:0] sumOUT;
10 // wire [19:0] sum;
11
12 wire[255:0] p;
13 wire[135:0] s1;
14 wire[71:0] s2;
15 wire[37:0] s3;
16
17 //Generate the multipliers
18 genvar i;
19 generate
20     for(i=0;i<=15;i=i+1)
21         multiplier multName(pixels[(127-8*i):(127-8*i-7)],
22                             weights[(127-8*i):(127-8*i-7)],
23                             p[(255-16*i):(255-16*i-15)]);
24 endgenerate
25
26 //Generate adders
27 genvar k;
28 generate
29     for(k=0;k<=7;k=k+1)
30         adder #(16) adderName(p[(255-16*2*k):(255-16*2*k-15)],
31                             p[(255-16*(2*k+1)): (255-16*(2*k+1)-15)],
32                             s1[(135-17*k):(135-17*k-16)]);
33 endgenerate
34
35 adder #(17) adder_17_1(s1[135:119], s1[118:102], s2[71:54]);
36 adder #(17) adder_17_2(s1[101:85], s1[84:68], s2[53:36]);
37 adder #(17) adder_17_3(s1[67:51], s1[50:34], s2[35:18]);
38 adder #(17) adder_17_4(s1[33:17], s1[16:0], s2[17:0]);
39

```

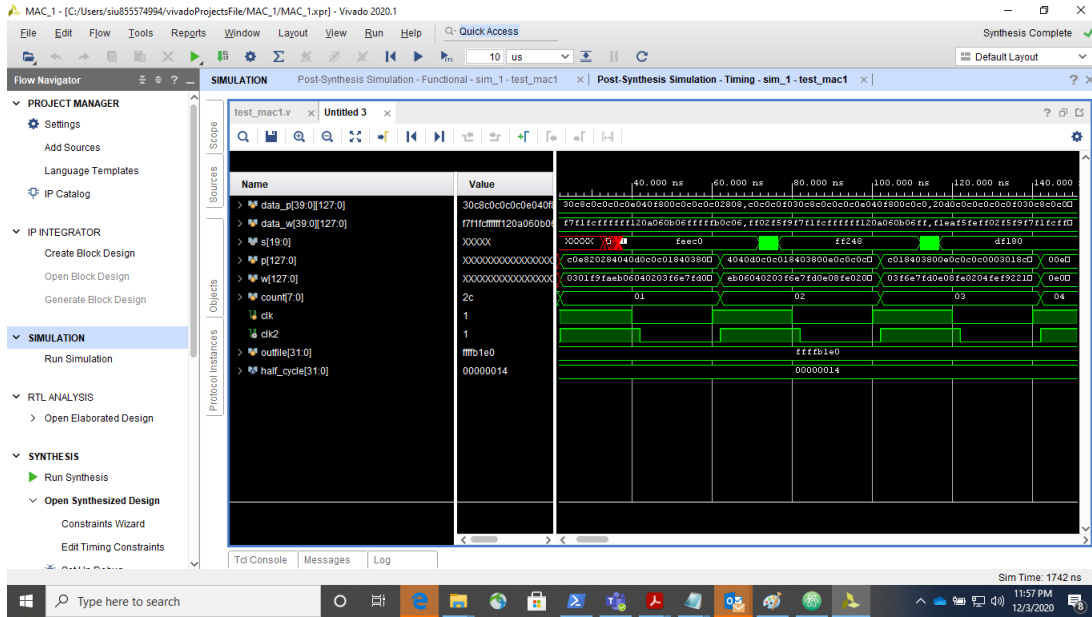


Figure 4: Timing waveform obtained after synthesis for MAC module

```

40   adder #(18) adder_18_1(s2[71:54], s2[53:36], s3[37:19]);
41   adder #(18) adder_18_2(s2[35:18], s2[17:0], s3[18:0]);
42
43   adder #(19) adder_19_1(s3[37:19], s3[18:0], sum);
44
45
46   // //Pipeline
47   // always @(posedge clk) begin
48   //     pixels <= pixelsIN;
49   //     weights <= weightsIN;
50   //     sumOut <= sum;
51   // end
52 endmodule

```

```

1
2 // `timescale 1ns/1ps
3
4 module mac2_tb();
5
6     parameter half_cycle = 20;
7     reg[127:0] data_p[39:0];
8     reg[127:0] data_w[39:0];
9     wire [19:0] s;
10
11     reg[127:0] p,w;
12     reg[7:0] count;
13     reg clk;
14     wire clk2;
15
16     integer outfile;
17
18     assign #2 clk2=clk;
19

```

```

20 mac2 uut( clk ,p,w,s );
21
22 initial begin
23     $readmemh("digits_hex.txt", data_p);
24     $readmemh("weights_hex.txt", data_w);
25     outfile=$fopen("simout.txt","w");
26     clk=0;
27     count=0;
28 end
29
30 always #half_cycle clk=!clk;
31
32 //Wriete to a file
33 always @(posedge clk)
34     if (count>0)
35         $fdisplay(outfile,"%h",s);
36
37 always @(posedge clk2)
38 begin
39     p=data_p[count];
40     w=data_w[count];
41     count=count+1;
42     if (count==41) begin
43         $fclose(outfile);
44         $finish;
45     end
46 end
47
48 endmodule

```

synthesize results expected result

00000.	
00000	feec0
00000	ff248
00000	ff180
ff248	fff58
ff180	001f0
fff58	ffd40
001f0	ff978
ffd40	.
ff978	.
.	.
.	.
.	.

Figure 5: With pipeline registers, the MAC output glitch

As shown above, I am missing feec0 in synthesized result instead I am getting a 00000 on 4th clock cycle . The next results are consistent and correct. This is not happening to my behavioral simulation though. This is avoided by using no pipeline registers.

1.2 Milestone 2: Accumulator (ACC) Module Design

Computing the output of a hidden node in the neural network requires adding 64 multiplication results. However, the MAC module I have designed earlier only compute 16 multiplication and add them together. Thus, I plan to use the same MAC module four times and accumulate the results as shown in Figure 6. Since MAC output is 20 bits and four MAC results are to be accumulated, the ACC will take four clock cycles to produce a result and the result has 22 bits.

The reason for accumulation is to decrease hardware cost by re-utilizing same module in expense of a slower circuit.

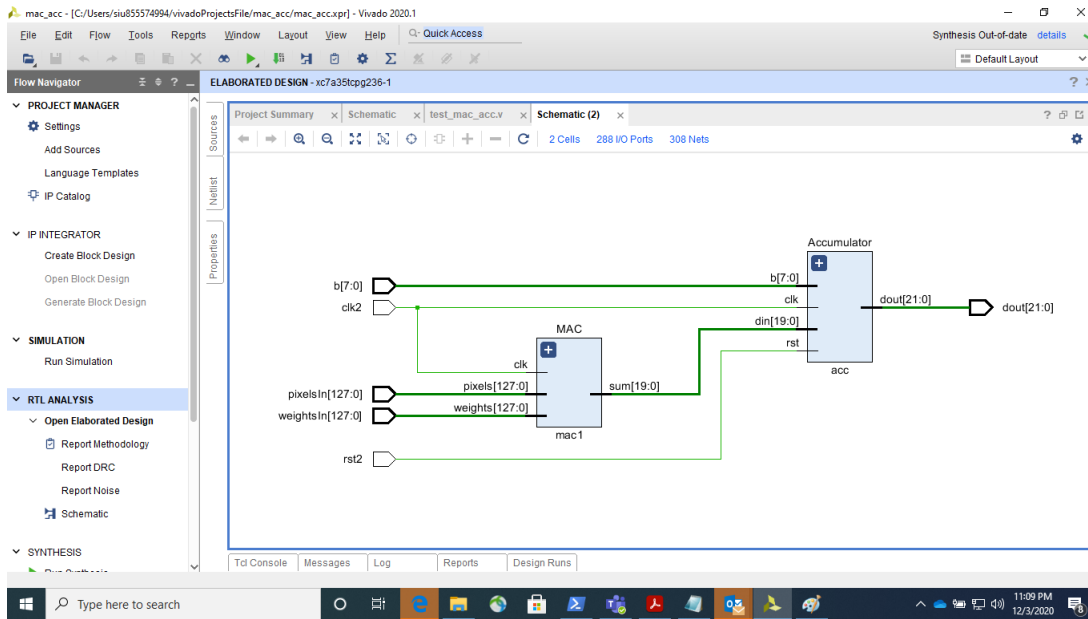


Figure 6: MAC and ACC module together

the bias was set to 11 and the given weight and pixel values are read through the testbench file. A accumulator control circuit was written to implement a state machine. The state machine controls the accumulation and cycles around a 4 clock cycle period. As the accumulator takes 4 MAC instance and accumulate them, for 1st clock cycle the ACC will add MAC output and bias b . For the remaining three clock cycle the adder accumulates the 3 instance of MAC module output and latch the results into the output register.

The schematic of the accumulator is shown in figure 7.

Obtained behavioral and post synthesis simulation results are shown in below two figures :

Some interesting finding was the vivado synthesizer requires some initialization time so the reset signal should be close to 120ns.

the top module and components code and testbench is given below:

```

1
2 `timescale 1ns / 1ps
3 module mac_acc(
4   input clk2 ,

```

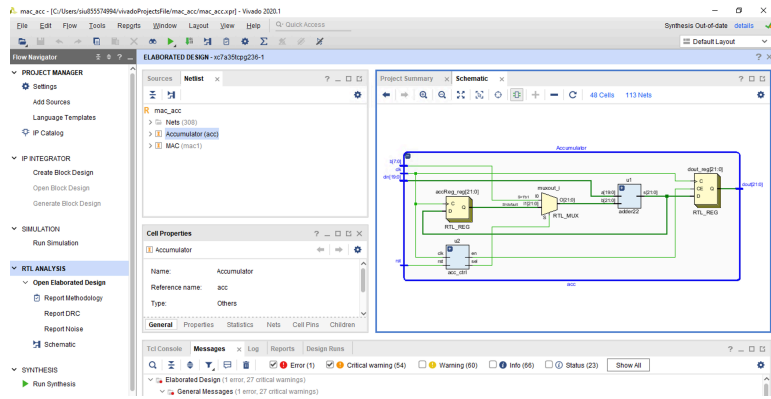


Figure 7: Accumulator Schematic

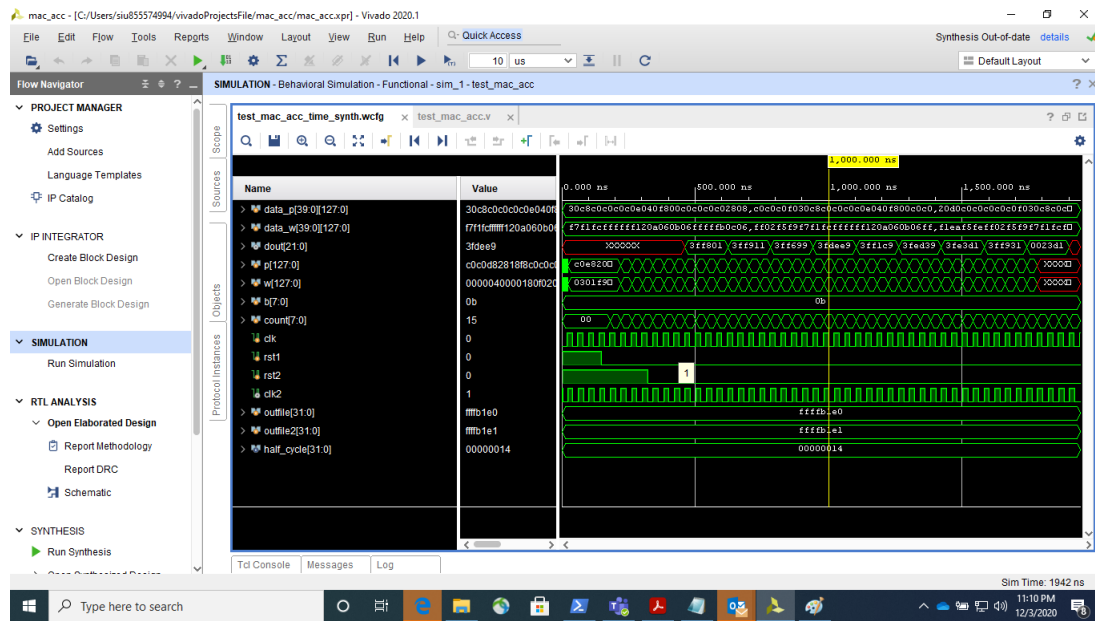


Figure 8: Behavioral wave-forms obtained from MAC+ACC module

```

5  input rst2 ,
6  input [127:0] pixelsIn ,
7  input [127:0] weightsIn ,
8  input [7:0] b ,
9  output [21:0] dout );
10
11  wire [19:0] macOut;
12  wire [21:0] dout;
13
14  mac1 MAC(clk2 , pixelsIn , weightsIn , macOut);
15  acc Accumulator(macOut,b , clk2 , rst2 , dout );
16
17  endmodule

```

```

1  `timescale 1ns / 1ps
2  module acc(
3      input [19:0] din ,

```

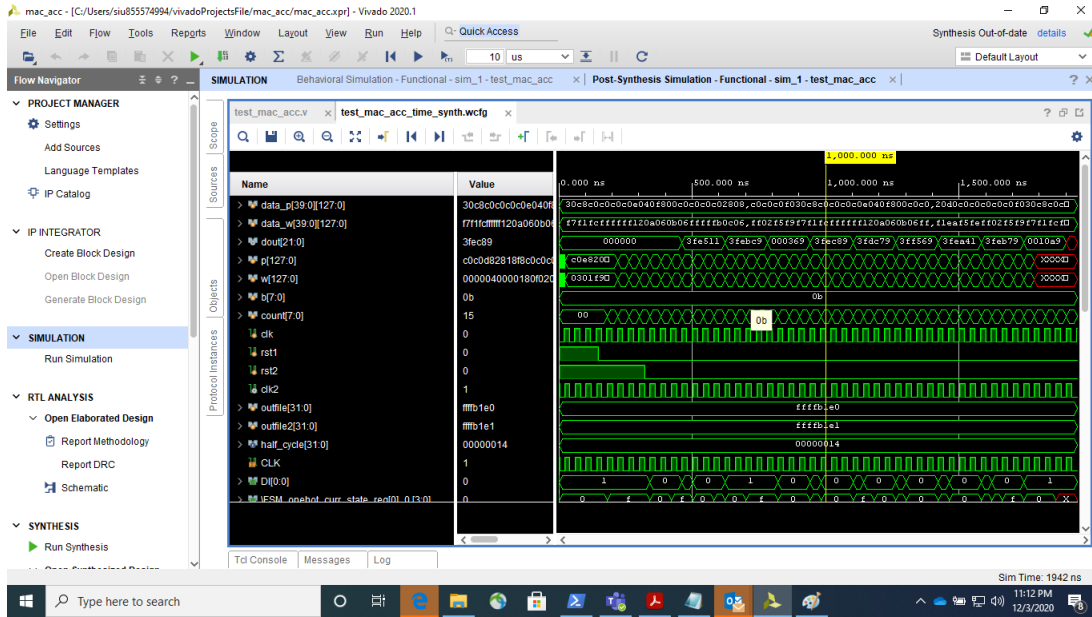


Figure 9: Post Synthesis wave-forms obtained from MAC+ACC module

```

4   input [7:0] b,
5   input clk ,
6   input rst ,
7   output [21:0] dout);
8
9   reg [21:0] accReg , muxout;
10  reg [21:0] dout;
11
12  wire sel ,en;
13  wire [21:0] b_ext ,sum;
14
15  adder22 u1(din , muxout,sum);
16  acc_ctrl u2(clk ,rst ,sel ,en);
17
18  assign b_ext = {{14{b[7]}} ,b}; //This is sign extension
19
20  always @(posedge clk) begin
21      if(en) begin
22          dout <=sum;
23      end
24  end
25
26
27  always @(posedge clk)
28      accReg <=sum;
29
30  always @(*) begin
31      if(sel)
32          muxout = b_ext;
33      else
34          muxout = accReg;
35  end

```



```

36     end
37 endmodule

1
2 `timescale 1ns / 1ps
3 module adder22(
4     input signed[19:0] a,
5     input signed[21:0] b,
6     output signed[21:0] s
7 );
8
9     assign s=a+b;
10
11 endmodule

1
2 `timescale 1ns / 1ps
3 module acc_ctrl(
4     input clk,
5     input rst,
6     output sel,
7     output en
8 );
9
10     reg[1:0] next_state, curr_state;
11     parameter s1 = 2'b00;
12     parameter s2 = 2'b01;
13     parameter s3 = 2'b10;
14     parameter s4 = 2'b11;
15
16     always @(posedge clk) begin
17         if(rst)
18             curr_state = s1;
19         else
20             curr_state = next_state;
21     end
22
23
24     always @(curr_state) begin
25         case(curr_state)
26             s1: next_state = s2;
27             s2: next_state = s3;
28             s3: next_state = s4;
29             s4: next_state = s1;
30         endcase
31     end
32
33     assign sel = (curr_state == s1) ? 1'b1 : 1'b0;
34     assign en = (curr_state == s4) ? 1'b1 : 1'b0;
35
36 endmodule

1
2 `timescale 1ns / 1ps
3 module test_mac_acc();
4     parameter half_cycle = 20;

```

```

5      reg [127:0] data_p[39:0];
6      reg [127:0] data_w[39:0];
7      wire [21:0] dout;
8      reg [127:0] p, w;
9      reg [7:0] b;
10     reg [7:0] count;
11     reg clk, rst1, rst2;
12     wire clk2;
13     integer outfile;
14     integer outfile2;
15
16     assign #2 clk2=clk; // delayed clk
17
18     mac_acc uut (clk2, rst2, p, w, b, dout);
19
20
21     initial begin
22         $readmemh("C:\\Users\\siu855574994\\OneDrive - Southern
Illinois University\\Fall 2020\\ECE-528\\projects\\Codes\\digits_hex.txt",data_p);
23         $readmemh("C:\\Users\\siu855574994\\OneDrive - Southern
Illinois University\\Fall 2020\\ECE-528\\projects\\Codes\\weights_hex.
txt",data_w);
24         outfile = $fopen("C:\\Users\\siu855574994\\OneDrive - Southern
Illinois University\\Fall 2020\\ECE-528\\projects\\Codes\\simout_mac_acc
.txt","w");
25         outfile2 = $fopen("C:\\Users\\siu855574994\\OneDrive - Southern
Illinois University\\Fall 2020\\ECE-528\\projects\\Codes\\simout_mac.
txt","w");
26         clk=0;
27         count=0;
28         rst1=1;
29         rst2=1;
30         b=11;
31         #150 rst1=0;
32     end
33
34     //Generate clock
35     always #half_cycle clk=!clk;
36     // write acc output to file
37     always @(posedge clk)
38         if ((count>7)&(count[1:0]==2'b00))
39             $fdisplay(outfile, "%h", dout);
40
41     // write mac output to file
42
43     always @(posedge clk) begin
44         if(count>3)begin
45             $fdisplay(outfile2, "%h", uut.dout);
46             //rst2=1;
47         end
48     end
49     always @(posedge clk2) begin
50         p=data_p[count];
51         w=data_w[count];

```

```

52         if (!rst1) begin
53             count= count+1;
54             if (count==4)
55                 #half_cycle rst2=0;
56
57             if (count==45) begin
58                 $fclose(outfile);
59                 $fclose(outfile2);
60                 $finish;
61             end
62         end
63     end
64 endmodule

```

1.3 Milestone 3: Integrating with sigmoid IP block

So far with MAC+ACC module I have computed following function.

$$y = \sum_{i=1}^N w_i \cdot x_i + b_i$$

But for a neural network node, a sigmoid non-linearity is needed to be incorporated to fulfill following expression.

$$f(y) = \frac{1}{1+e^{-y}} - 0.5$$

as the sigmoid module was given, I wrote a IP-wrapper module which takes the 22 bit output from the MAC+ACC module and pass it through the Sigmoid IP block to generate 8 bit final output.

The schematic for this module is shown in figure 10.

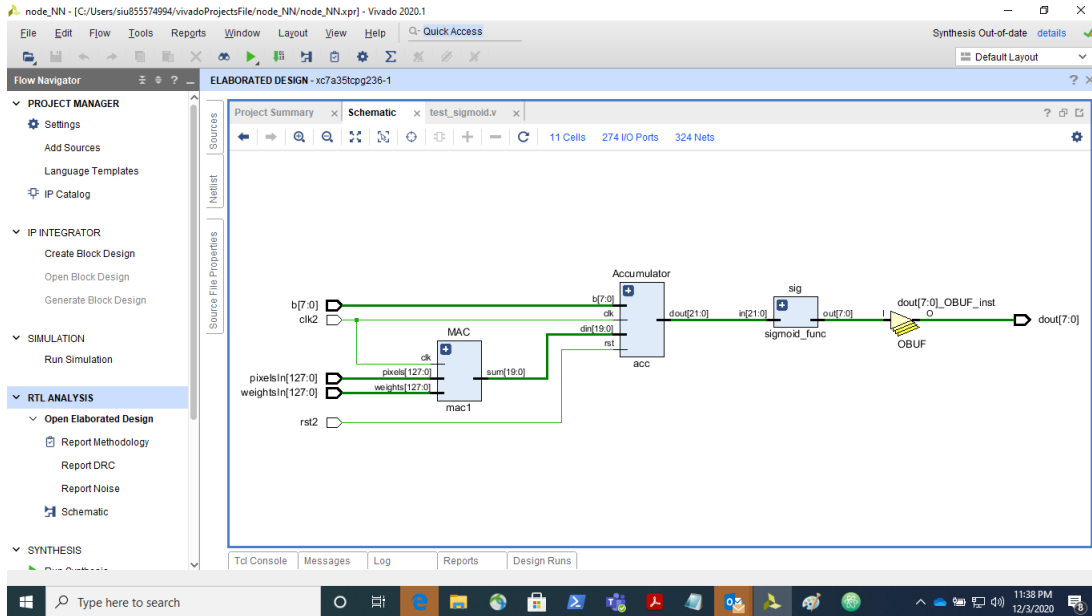


Figure 10: Schematic of MAC+ACC+Sigmoid block or a single neuron

The schematic of wrapper and sigmoid function is given in below figure :

The output obtained from this combination module : MAC+ACC+Sigmoid is shown in below:

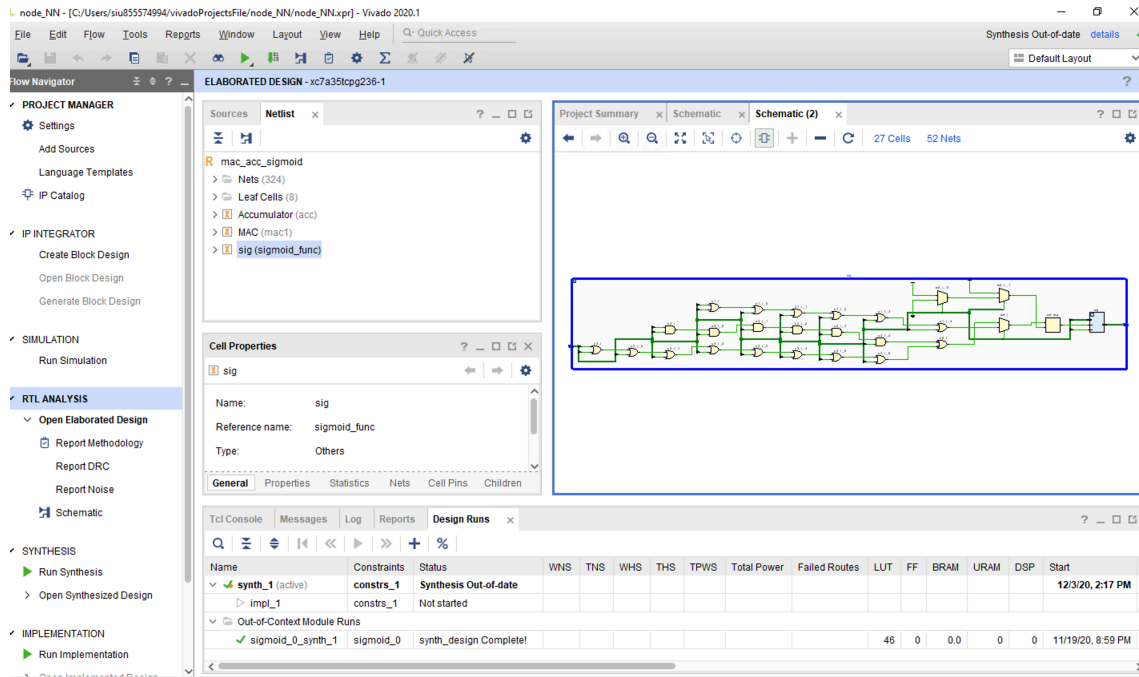


Figure 11: Schematic of IP wrapper and Sigmoid block

Top module, the IP-wrapper and sigmoid (sigmoid_func) and testbench codes are given below:

```

1
2 `timescale 1ns / 1ps
3 module mac_acc_sigmoid(
4     input clk2 ,
5     input rst2 ,
6     input [127:0] pixelsIn ,
7     input [127:0] weightsIn ,
8     input [7:0] b,
9     output [7:0] dout);
10
11     wire [19:0] macOut;
12     wire [21:0] dout_acc;
13     wire [7:0] dout;
14
15     mac1 MAC(clk2 , pixelsIn , weightsIn , macOut);
16     acc Accumulator(macOut,b, clk2 , rst2 , dout_acc);
17     sigmoid_func sig(dout_acc , dout);
18
19 endmodule

```

```

1
2 module sigmoid_func(
3     input [21:0] in ,
4     output [7:0] out
5 );
6
7 wire [16:0] y;
8 reg ovf;

```

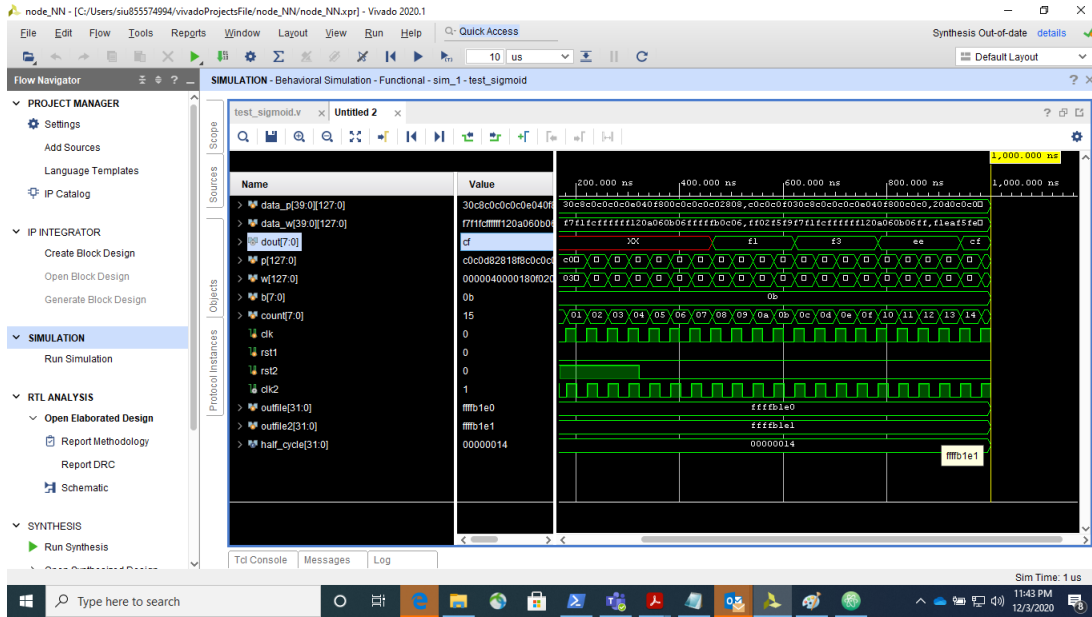


Figure 12: Final waveform obtained from a single node

```

9 wire sign;
10 wire out1,out2,out3,out4,out5,out6;
11 wire [8:0] data;
12
13 assign y=in[21:5];
14 assign sign=y[16];
15 assign data=y[8:0];
16 or u1(out1,y[15],y[14],y[13],y[12],y[11],y[10],y[9]);
17 and u2(out2,y[15],y[14],y[13],y[12],y[11],y[10],y[9]);
18 or u3(out3,y[8],y[7],y[6],y[5],y[4],y[3],y[2],y[1],y[0]);
19 not u4(out4,out2);
20 not u5(out5,out3);
21 or u6(out6,out4,out5);
22
23 always @(*) begin
24     if (sign ==1'b0) begin
25         ovf<=out1;
26     end
27     if (sign ==1'b1) begin
28         ovf<=out6;
29     end
30 end
31
32 sigmoid_0 sig(data,sign,ovf,out);
33
34 endmodule

```

```

1
2 `timescale 1ns / 1ps
3 module test_sigmoid( );
4
5 parameter half_cycle = 20;
6 reg [127:0] data_p[39:0];

```

```

7 reg [127:0] data_w [39:0];
8 wire [7:0] dout;
9 reg [127:0] p, w;
10 reg [7:0] b;
11 reg [7:0] count;
12 reg clk, rst1, rst2;
13 wire clk2;
14 integer outfile;
15 integer outfile2;
16
17 assign #2 clk2=clk; // delayed clk
18
19 mac_acc_sigmoid uut (clk2, rst2, p, w, b, dout);
20
21
22 initial begin
23 $readmemh("C:\\Users\\siu855574994\\OneDrive - Southern Illinois
    University\\Fall 2020\\ECE-528\\projects\\Codes\\digits_hex.txt", data_p)
    ;
24 $readmemh("C:\\Users\\siu855574994\\OneDrive - Southern Illinois
    University\\Fall 2020\\ECE-528\\projects\\Codes\\weights_hex.txt", data_w
    );
25 outfile = $fopen("C:\\Users\\siu855574994\\OneDrive - Southern Illinois
    University\\Fall 2020\\ECE-528\\projects\\Codes\\simout_sigmoid_node_syn
    .txt", "w");
26 outfile2 = $fopen("C:\\Users\\siu855574994\\OneDrive - Southern Illinois
    University\\Fall 2020\\ECE-528\\projects\\Codes\\simout_mac_sigmoid.txt
    ", "w");
27 clk=0;
28 count=0;
29 rst1=1;
30 rst2=1;
31 b=11;
32 #150 rst1=0;
33 end
34
35 //Generate clock
36 always #half_cycle clk=!clk;
37 // write acc output to file
38 always @(posedge clk)
39     if ((count>7)&(count[1:0]==2'b00))
40         $fdisplay(outfile, "%h", dout);
41
42 // write mac output to file
43
44 always @(posedge clk) begin
45     if(count>3)begin
46         $fdisplay(outfile2, "%h", uut.dout);
47         //rst2=1;
48     end
49 end
50 always @(posedge clk2) begin
51     p=data_p[count];
52     w=data_w[count];
53     if (!rst1) begin

```

```

54     count= count+1;
55     if (count==4)
56         #half_cycle rst2=0;
57
58     if (count==45) begin
59         $fclose(outfile);
60         $fclose(outfile2);
61         $finish;
62     end
63 end
64 end
65 endmodule

```

1.4 Design with considerations

Beyond engineering one critical consideration was implemented in this design and that is economic consideration. In the accumulation unit, I utilize same one MAC module to do 4 different multiplication. The price paid was instead of multiplying each weight and pixel value pairs con-currently, 4 clock cycle is required to compute same computation sequentially. So the reason for accumulation unit is to decrease hardware cost by re-utilizing same module in expense of a slower circuit.