

XNOR-Net and Binary Connect Hardware Design

Lucas B., Bibrak C., Adam D., Prateek S.

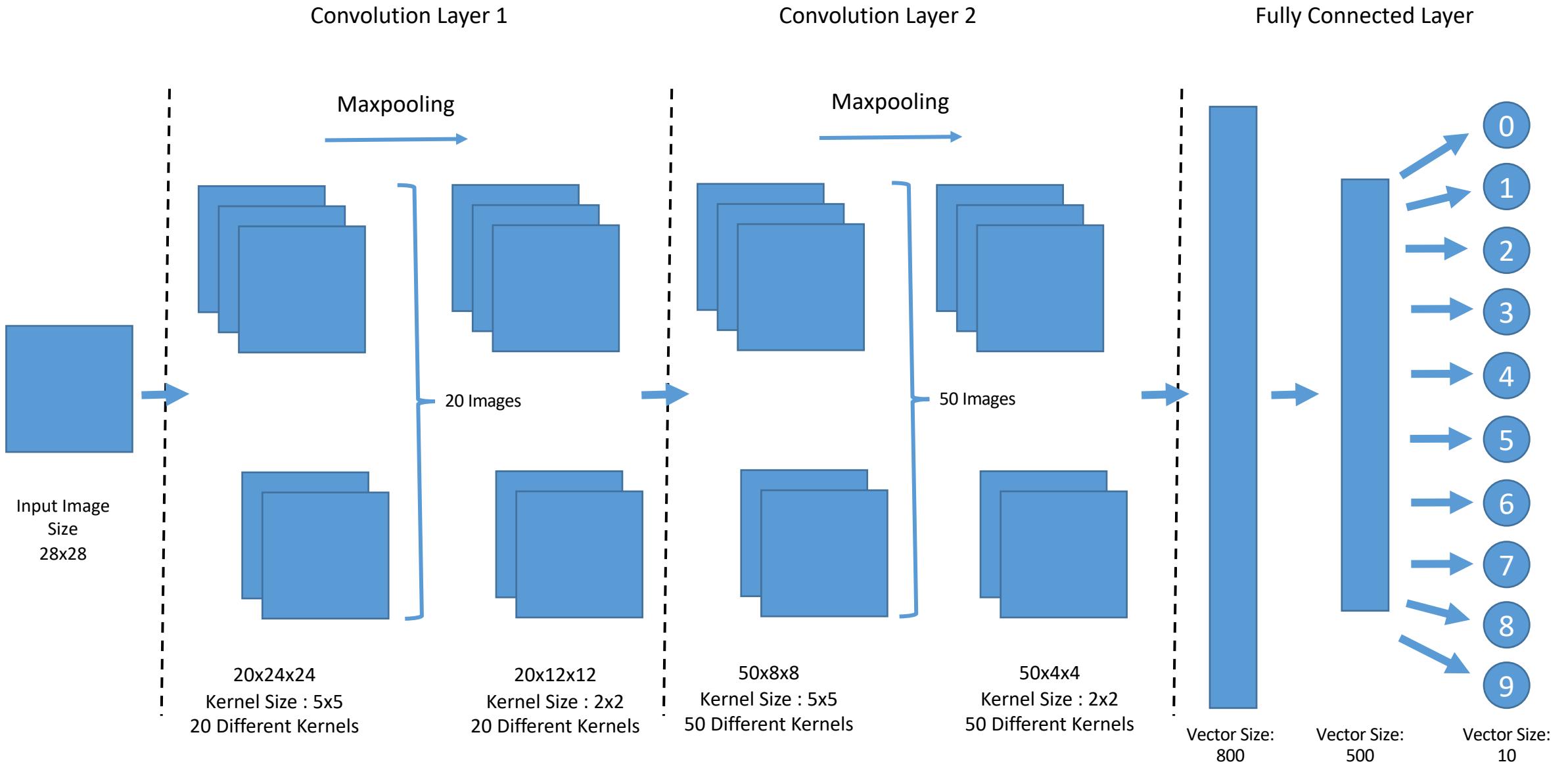
E501 Design Project Presentation
4/16/2018

Introduction

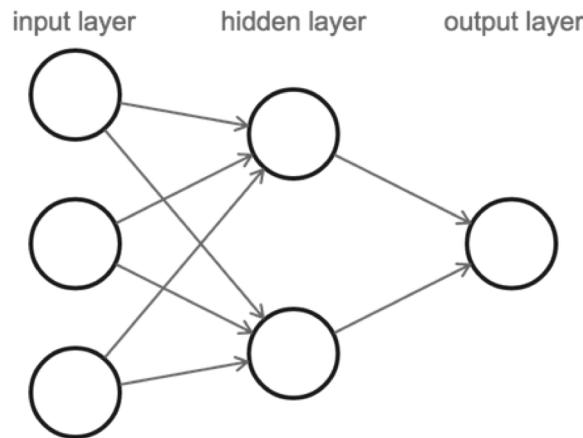
- Hardware designs were created to implement the basic hardware building blocks of the following papers:
 - “XNOR-Net: ImageNet Classification Using Binary Convolutional Neural Networks”
 - “BinaryConnect: Training Deep Neural Networks with binary weights during propagations”

Feed-forward network using a trained neural network model

- For the current scope of the project, we consider MNIST dataset which consists of 10 classes of images i.e. 0 to 9.
- Once the training of the neural network is done using any machine, our hardware design focuses on achieving speedup in the classification phase.
- The weights learned by the neural network will be inserted in the module using a BRAM to the feed forward network.
- These weights along with the images to be classified would be used as an input to the feed forward network on an FPGA.
- The scope of our project is only defined for the implementation of the fully connected layer after the convolutions have been done.



Fully Connected Layer



Input Layer becomes a vector, the weights which are multiplied and accumulated become a matrix. This whole process of Matrix and Vector multiplication along with accumulation is simply a Matrix-Vector Multiplication.

Input Vector (I) : 800×1

Weight Matrix (W_1) : 500×800

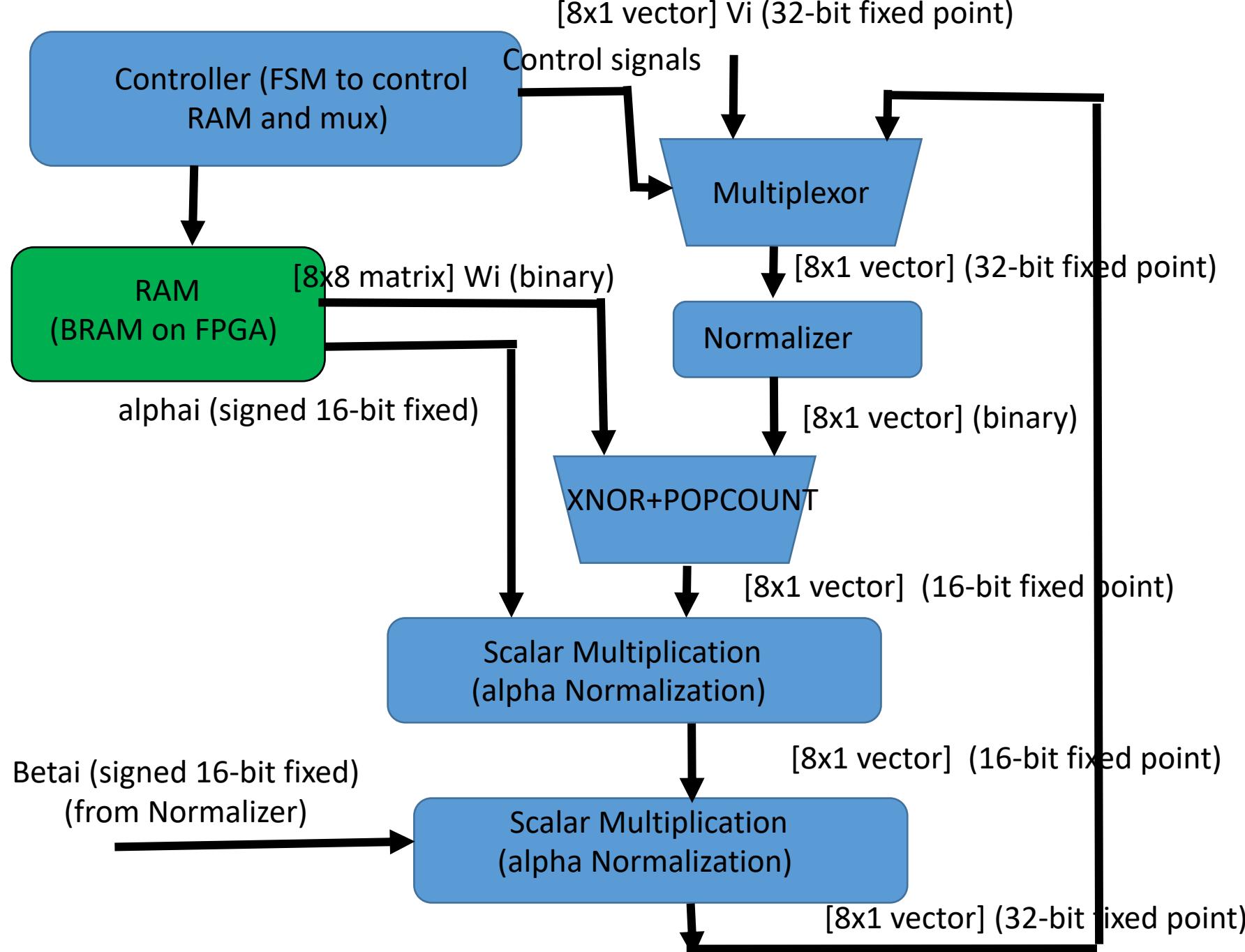
Hidden Layer (H) : $H = W_1 \times I = [500 \times 800] [800 \times 1] = [500 \times 1]$

Weight Matrix (W_2) : 10×500

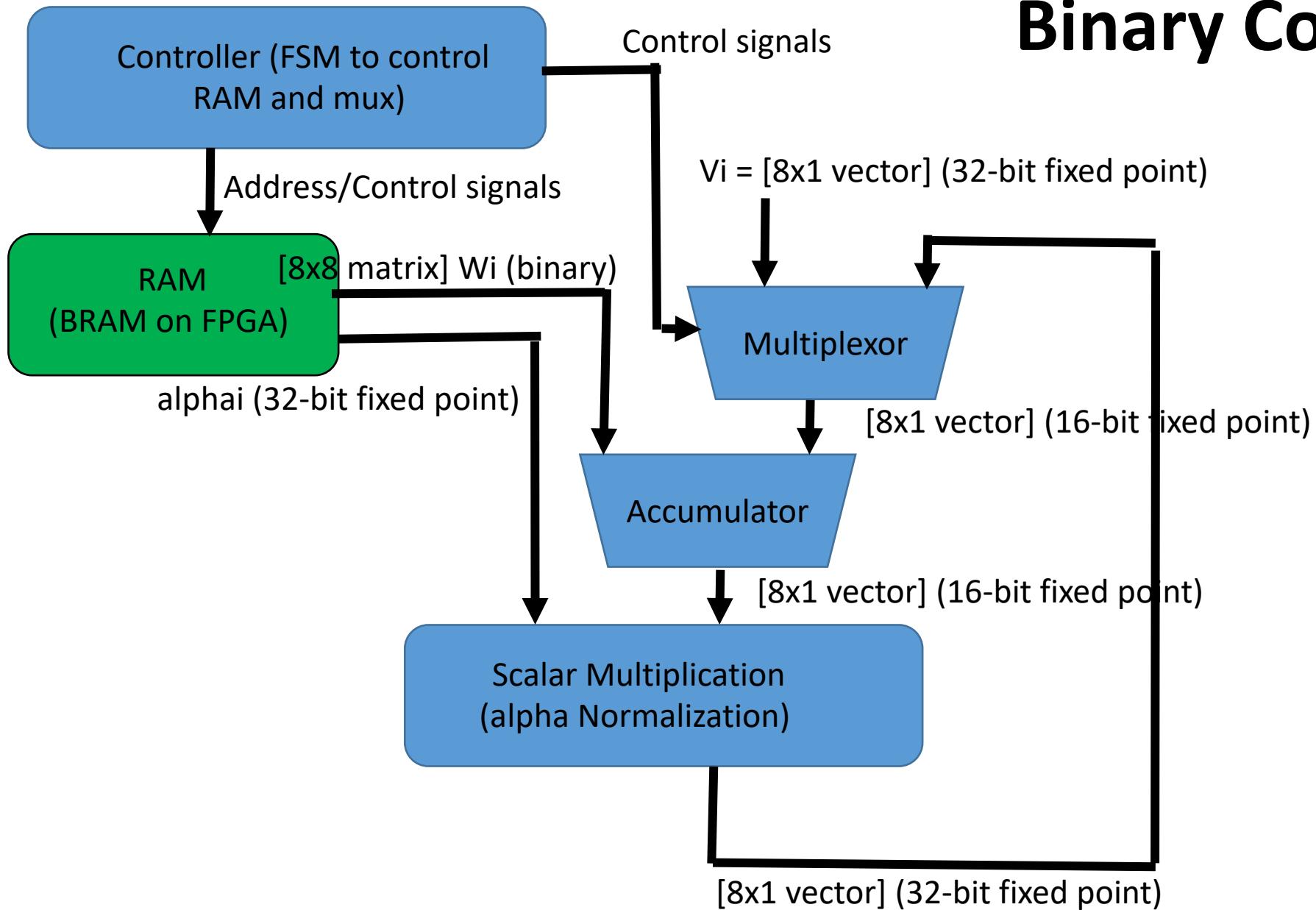
Output Layer (O) : $O = W_2 \times H = [10 \times 500] [500 \times 1] = [10 \times 1]$

This output layer which is a vector of 10 values classifies all the images in the possible 10 numbers according to MNIST Data Set we use.

XNOR-Net



Binary Connect



XNOR-Net Core Components

XNOR + Popcount Design Concept

Performs multiplication and accumulation(MAC) operation optimizations.

Alpha and Beta Multipliers

Performs multiplication of vector to scalar

Binary Connect Core Components

Accumulator

Performs multiplication and accumulation(MAC) operation optimizations.

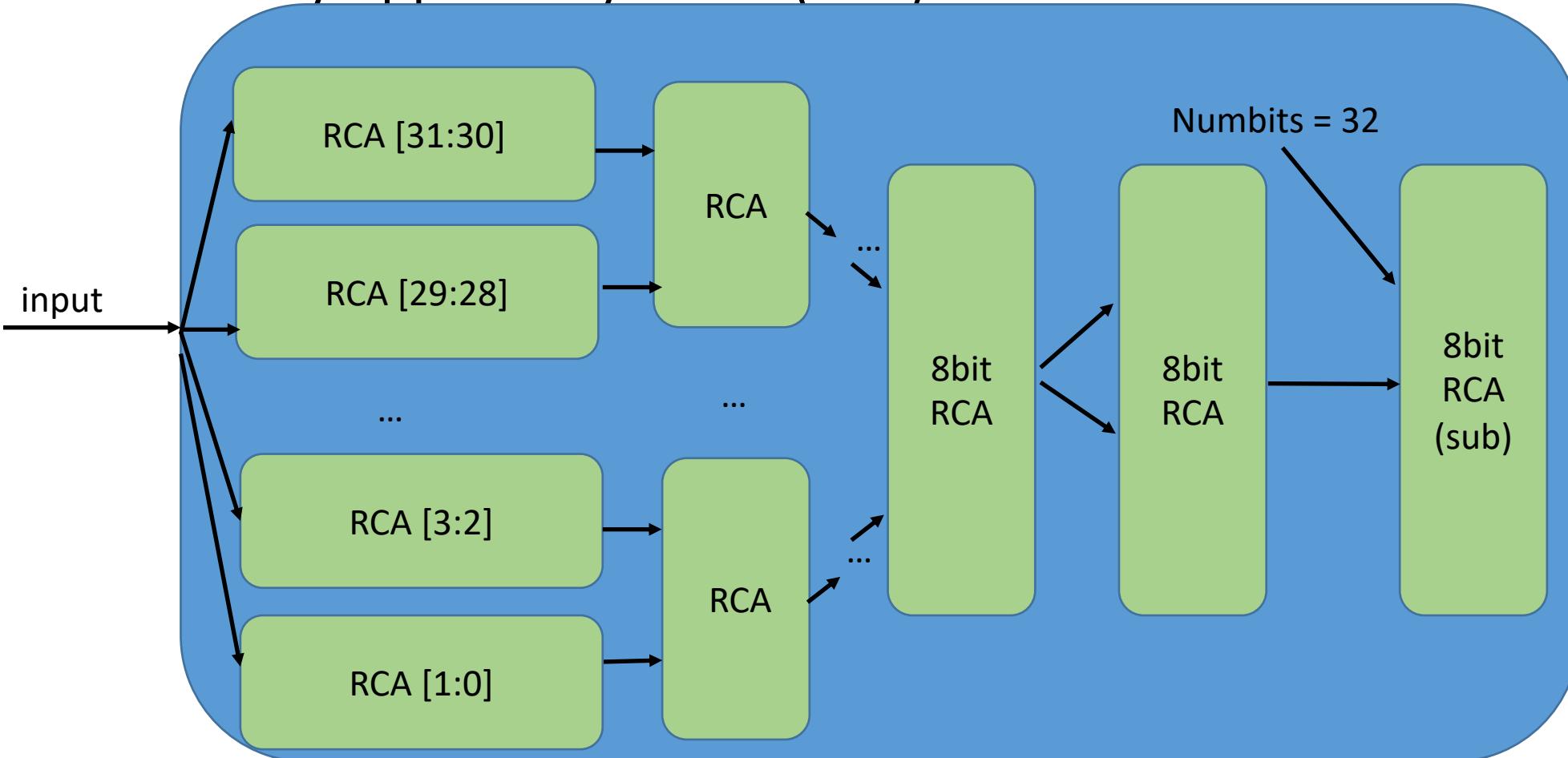
Alpha Multiplier

Performs multiplication of vector to scalar

Modules

Popcount Design Concept 1 (Slow)

- Purely ripple carry adder (RCA) combination



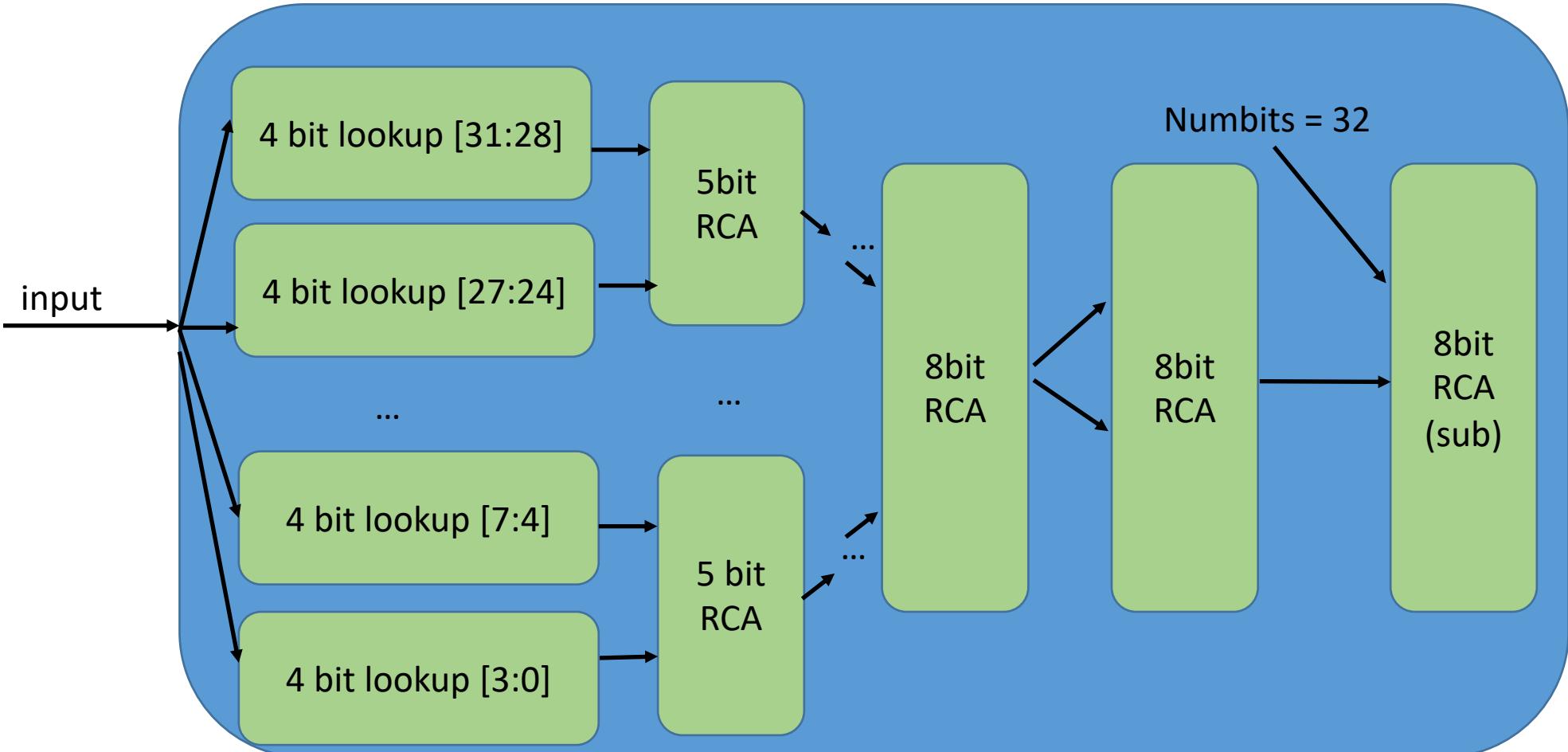
- Also, If you simply count the number of ones in a string, you need to apply:
 - $\text{POPCOUNT} = (\text{NUMBER OF ONES}) * 2 - (\text{NUMBER OF BITS})$

4 bit lookup table

Input [3:0]	Popcount
0000	0
0001	1
0010	1
0011	2
0100	1
0101	2
0110	2
0111	3
1000	1
1001	2
1010	2
1011	3
1100	2
1101	3
1110	3
1111	4

Popcount Design Concept 2 (Better)

- 8 of 4-bit Lookups + ripple carry adder (RCA) combination
- Look up table saves adders and reduces propagation time



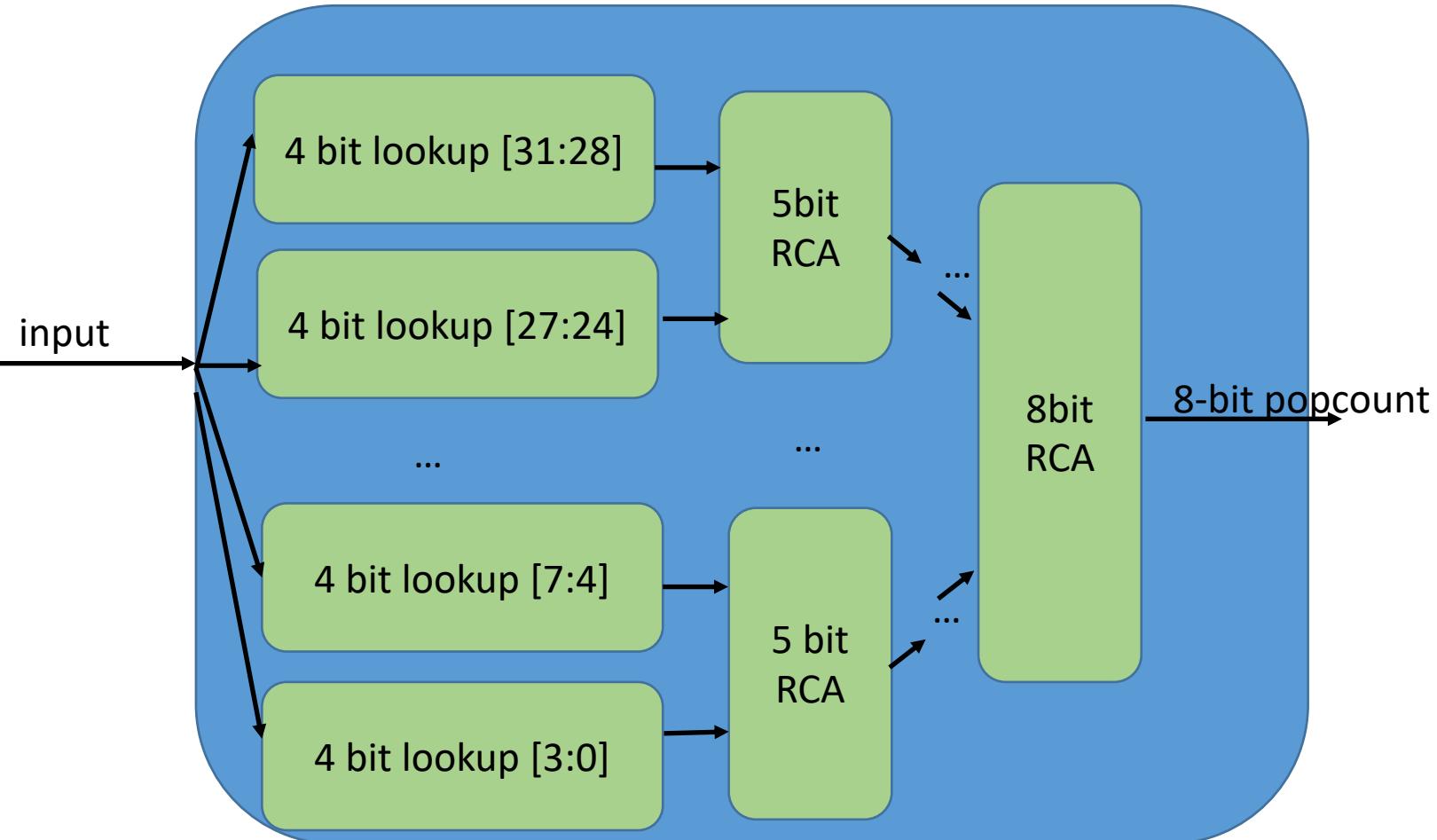
- Still, you simply count the number of ones in a string, you need to apply:
 - $\text{POPCOUNT} = (\text{NUMBER OF ONES}) * 2 - (\text{NUMBER OF BITS})$

4 bit lookup

Input [3:0]	Popcount
0000	-4
0001	-2
0010	-2
0011	0
0100	-2
0101	0
0110	0
0111	2
1000	-2
1001	0
1010	0
1011	2
1100	0
1101	2
1110	2
1111	4

Popcount Design Concept 3 (Best)

- 8 of 4-bit Lookups + ripple carry adder (RCA) combination



- Building the negative popcount values into LUT:
 - Saves 2 RCA components

8-bit Popcount Design

```

entity Popcount_4 is
  Port ( input_vector : in STD_LOGIC_VECTOR(3 downto 0);
         count : out STD_LOGIC_VECTOR(7 downto 0)
        );
end Popcount_4;

architecture Behavioral of Popcount_4 is

begin
  process (input_vector)
  begin
    -- 2s comp notation for 8 bits
    -- -4 = "FC"
    -- -2 = "FE"
    -- 0 = "00"
    -- +2 = "02"
    -- +4 = "04"

    case input_vector is
      when x"0" => count <= x"FC";
      when x"1" | x"2" | x"4" | x"8" => count <= x"FE";
      when x"3" | x"5" | x"6" | x"9" | x"A" | x"C" => count <= x"00";
      when x"7" | x"B" | x"D" | x"E" => count <= x"02";
      when x"E" => count <= x"04";
      when others => count <= x"00";
    end case;
  end process;
end Behavioral;

```

```

entity Popcount_8 is
  Port ( input_vector : in STD_LOGIC_VECTOR(31 downto 0);
         count : out STD_LOGIC_VECTOR(7 downto 0)
        );
end Popcount_8;

architecture Behavioral of Popcount_8 is

component Popcount_4
  port ( input_vector : in STD_LOGIC_VECTOR(3 downto 0);
         count : out STD_LOGIC_VECTOR(7 downto 0)
        );
end component;

component Ripple_Carry_Adder
  Port ( SUB: in STD_LOGIC; -- '0' = add, '1' = sub
         A : in STD_LOGIC_VECTOR (7 downto 0);
         B : in STD_LOGIC_VECTOR (7 downto 0);
         S : out STD_LOGIC_VECTOR (7 downto 0);
         COUT : out STD_LOGIC;
         V : out STD_LOGIC
        );
end component;

signal count1 : std_logic_vector(7 downto 0) := x"00";
signal count0 : std_logic_vector(7 downto 0) := x"00";

begin
  pop1: Popcount_4 PORT MAP (input_vector => input_vector(23 downto 20), count => count1);
  pop0: Popcount_4 PORT MAP (input_vector => input_vector(19 downto 16), count => count0);
  add_0_0: Ripple_Carry_Adder PORT MAP (SUB => '0', A => count1, B => count0, S => count);

end Behavioral;

```

8-bit Style	# RCA	# LUT4	Speed	Area
Just RCA	9	0	slow	large
RCA+LUT	3	2	fast	med
RCA+LUT opt	1	2	fast	small

XNOR and Popcount Results

Name	Value	999,996 ps	999,997 ps	999,998 ps	999,999 ps	1.
► weights[7:0]	[01,33,45]		[01,33,45,67,89,ab,cd,ef]			
► nn_in[7:0]	01100011		01100011			
► xnor_0[7:0]	01110011		01110011			
► xnor_1[7:0]	01010001		01010001			
► xnor_2[7:0]	00110111		00110111			
► xnor_3[7:0]	00010101		00010101			
► xnor_4[7:0]	11111011		11111011			
► xnor_5[7:0]	11011001		11011001			
► xnor_6[7:0]	10101111		10101111			
► xnor_7[7:0]	10011101		10011101			
▼ nn_out[7:0]	[00020000	[00020000,00040000,00020000,00060000,00fe0000,00020000,00fe0000,...				
► [7]	00020000		00020000			
► [6]	00040000		00040000			
► [5]	00020000		00020000			
► [4]	00060000		00060000			
► [3]	00fe0000		00fe0000			
► [2]	00020000		00020000			
► [1]	00fe0000		00fe0000			
► [0]	00020000		00020000			

Scalar Multipliers

Multiplies a fixed point scalar value to an integer

```
8
9 entity scalar_multiplier_ard is
10    Port ( scalar : in  sfixed (7 downto -8);
11          nn_in : in vector_fixed_point_type;
12          nn_out : out sfixed_fixed_point_type
13          );
14 end scalar_multiplier_ard;
15
16 architecture Behavioral of scalar_multiplier_ard is
17
18 begin
19
20    nn_out(0) <= scalar * to_sfixed(to_integer(unsigned(nn_in(0)(31 downto 16))),7,-8);
21    nn_out(1) <= scalar * to_sfixed(to_integer(unsigned(nn_in(1)(31 downto 16))),7,-8);
22    nn_out(2) <= scalar * to_sfixed(to_integer(unsigned(nn_in(2)(31 downto 16))),7,-8);
23    nn_out(3) <= scalar * to_sfixed(to_integer(unsigned(nn_in(3)(31 downto 16))),7,-8);
24    nn_out(4) <= scalar * to_sfixed(to_integer(unsigned(nn_in(4)(31 downto 16))),7,-8);
25    nn_out(5) <= scalar * to_sfixed(to_integer(unsigned(nn_in(5)(31 downto 16))),7,-8);
26    nn_out(6) <= scalar * to_sfixed(to_integer(unsigned(nn_in(6)(31 downto 16))),7,-8);
27    nn_out(7) <= scalar * to_sfixed(to_integer(unsigned(nn_in(7)(31 downto 16))),7,-8);
28
29 end Behavioral;
30
```

Scalar Multipliers

Multiplies a fixed point scalar value to a fixed point

```
o
9  entity scalar_multiplier_last is
10     Port ( scalar : in  sfixed (7 downto -8);
11             nn_in : in sfixed_fixed_point_type;
12             nn_out : out sfixed_fixed_point_type
13         );
14 end scalar_multiplier_last;
15
16 architecture Behavioral of scalar_multiplier_last is
17
18 begin
19
20     nn_out(0) <= scalar * nn_in(0)(7 downto -8);
21     nn_out(1) <= scalar * nn_in(1)(7 downto -8);
22     nn_out(2) <= scalar * nn_in(2)(7 downto -8);
23     nn_out(3) <= scalar * nn_in(3)(7 downto -8);
24     nn_out(4) <= scalar * nn_in(4)(7 downto -8);
25     nn_out(5) <= scalar * nn_in(5)(7 downto -8);
26     nn_out(6) <= scalar * nn_in(6)(7 downto -8);
27     nn_out(7) <= scalar * nn_in(7)(7 downto -8);
28
29 end Behavioral;
30
31
```

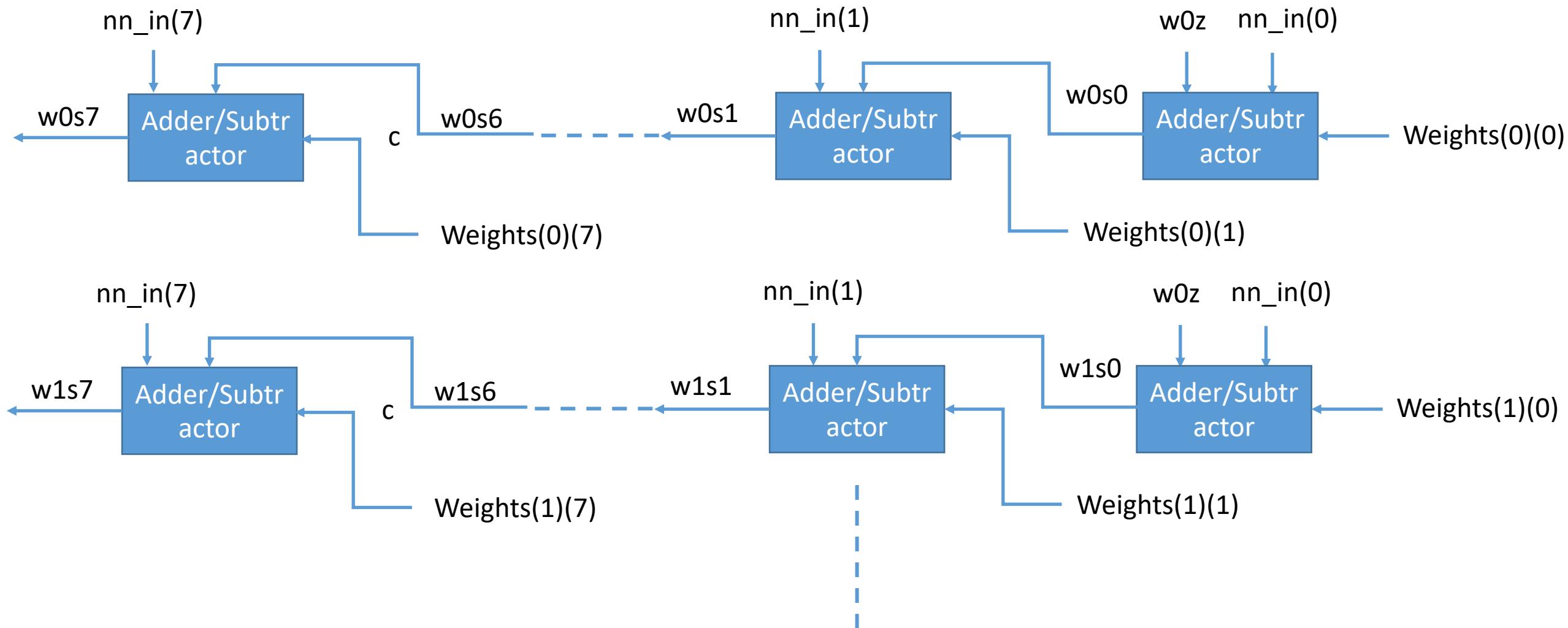
Accumulator

- This module is the core of Binary Connect.
- Multiplying the fixed point numbers to the binary weight matrix essentially is not multiplication. It is addition or subtraction depending on the corresponding weight in the weight matrix. For example, a three dimensional binary weight matrix and an input vector when multiplied behaves as follows:

$$\begin{bmatrix} 1 & 0 & 1 \\ 1 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \times \begin{bmatrix} 1.0 \\ 2.5 \\ 3.5 \end{bmatrix} = \begin{bmatrix} 1 - 2.5 + 3.5 \\ 1 + 2.5 - 3.5 \\ -1 - 2.5 - 3.5 \end{bmatrix}$$

- Accumulator thus, removes the multiplication and replaces it with addition and subtraction operations.
- As the accumulator is cheaper than the multiplier, its usage in the neural network design grants speedup.

Block Diagram for accumulator



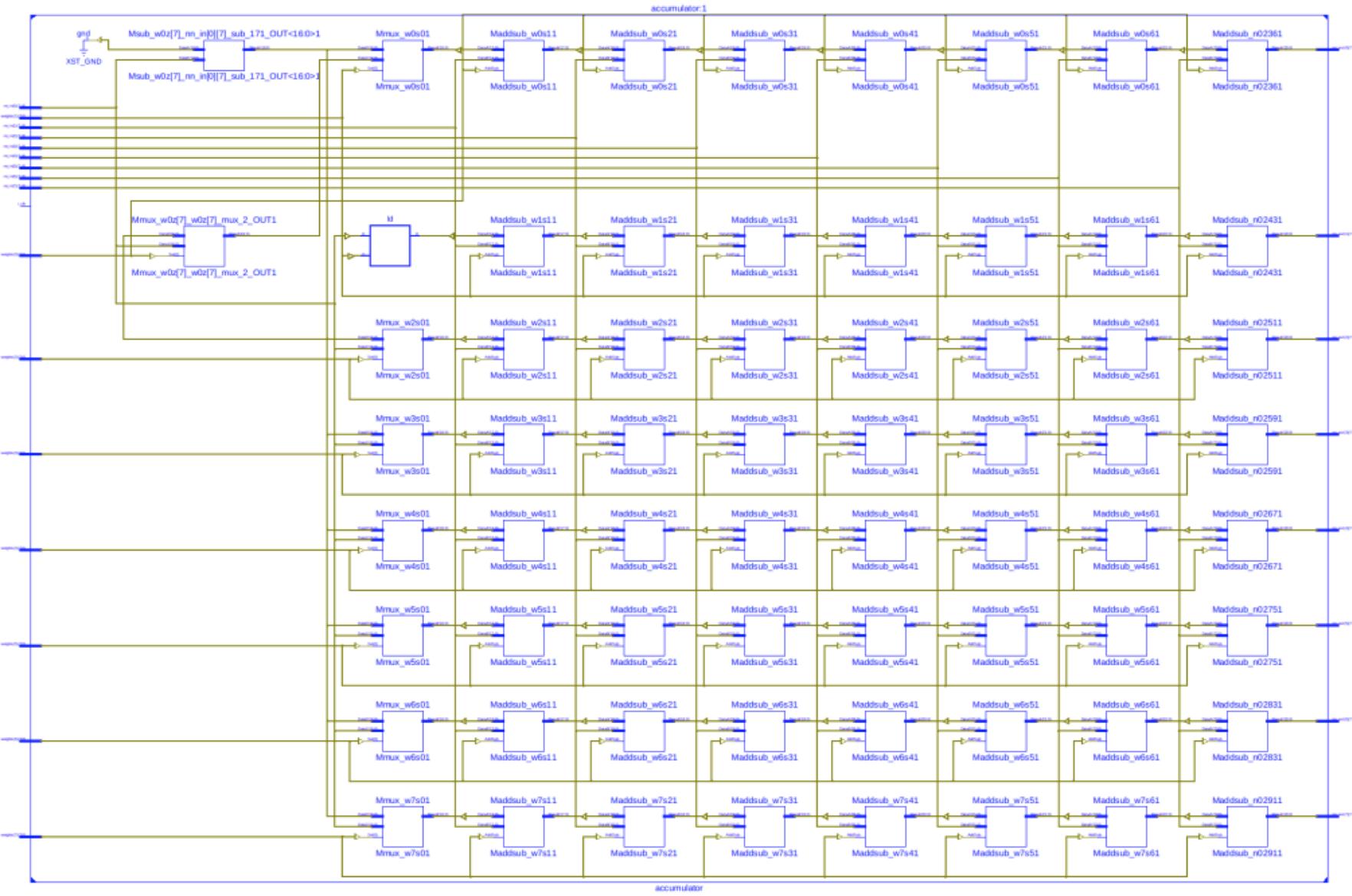
Computation mechanism

The following matrix multiplication (but for 8x1 vector) is computed in the following manner:

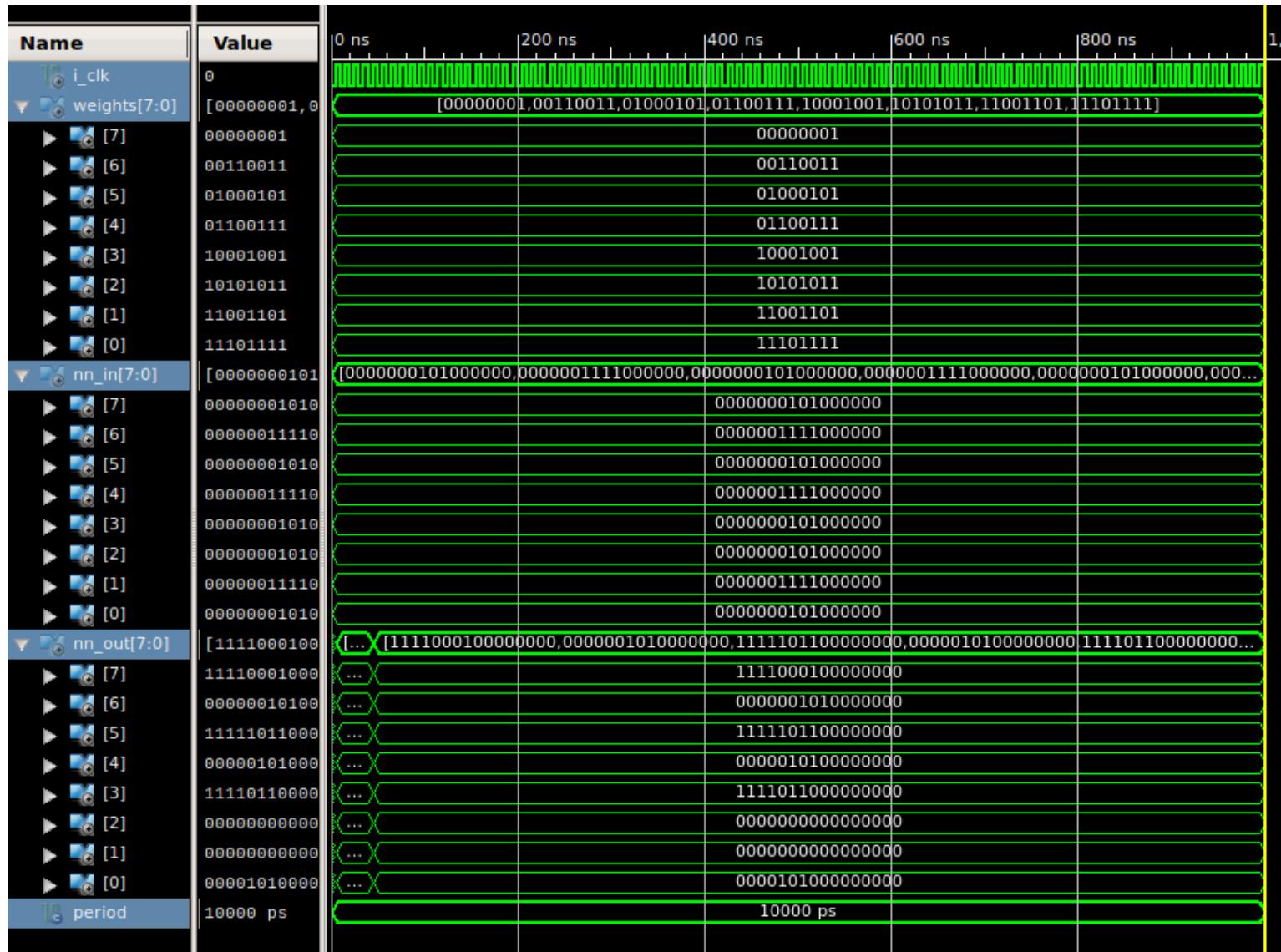
$$\begin{bmatrix} 1 & 0 & 1 \\ 1 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \times \begin{bmatrix} 1.0 \\ 2.5 \\ 3.5 \end{bmatrix} = \begin{bmatrix} 1 - 2.5 + 3.5 \\ 1 + 2.5 - 3.5 \\ -1 - 2.5 - 3.5 \end{bmatrix}$$

```
if (weights(0)(0) = '1') then w0s0 <= w0z + nn_in(0); else w0s0 <= w0z - nn_in(0); end if;
if (weights(0)(1) = '1') then w0s1 <= w0s0 + nn_in(1); else w0s1 <= w0s0 - nn_in(1); end if;
if (weights(0)(2) = '1') then w0s2 <= w0s1 + nn_in(2); else w0s2 <= w0s1 - nn_in(2); end if;
if (weights(0)(3) = '1') then w0s3 <= w0s2 + nn_in(3); else w0s3 <= w0s2 - nn_in(3); end if;
if (weights(0)(4) = '1') then w0s4 <= w0s3 + nn_in(4); else w0s4 <= w0s3 - nn_in(4); end if;
if (weights(0)(5) = '1') then w0s5 <= w0s4 + nn_in(5); else w0s5 <= w0s4 - nn_in(5); end if;
if (weights(0)(6) = '1') then w0s6 <= w0s5 + nn_in(6); else w0s6 <= w0s5 - nn_in(6); end if;
if (weights(0)(7) = '1') then w0s7 <= w0s6 + nn_in(7); else w0s7 <= w0s6 - nn_in(7); end if;
nn_out(0)(7 downto -8) <= w0s7(7 downto -8);
```

Schematic



Waveform



Scalar multiplier

- Implementation of a 32-bit fixed-point multiplication
 - 16-bit inputs: 1-bit signal, 7-bit integer, 8-bit fraction
- To perform fixed-point multiplication
 - Consider inputs as two-complements and do 'normal' multiplication, then determine where the "point" is.

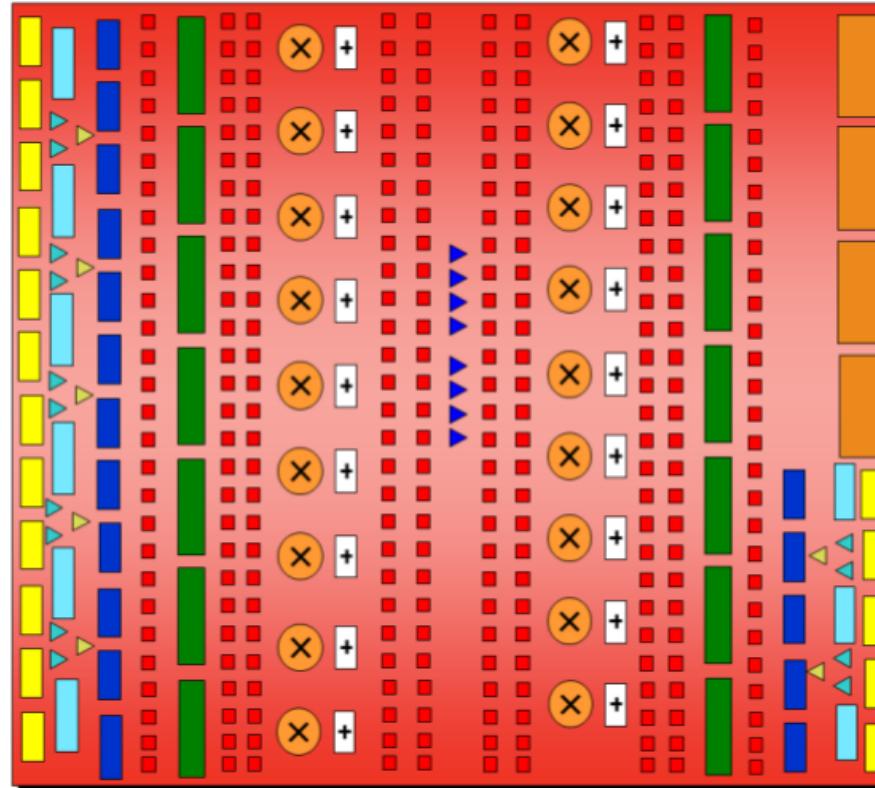
$$\begin{array}{r} 1 \\ 2 \\ \hline 3 \\ 4 \\ 5 \\ 6 \\ 7 \\ \hline 8 \end{array} \times \begin{array}{r} 1 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 1 \\ -7 \\ -14 \\ \hline 98 \end{array}$$

+

$$\begin{array}{r} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 & 1 \\ \hline 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 & 1 & 0 & 98 \end{array}$$

Xilinx 7 Series FPGA

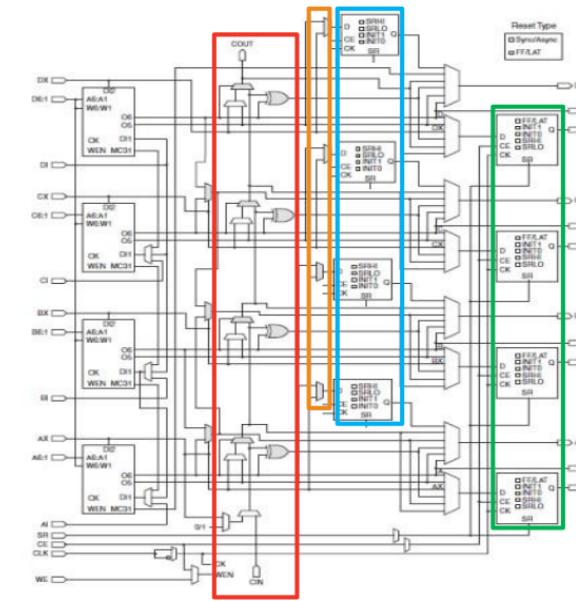
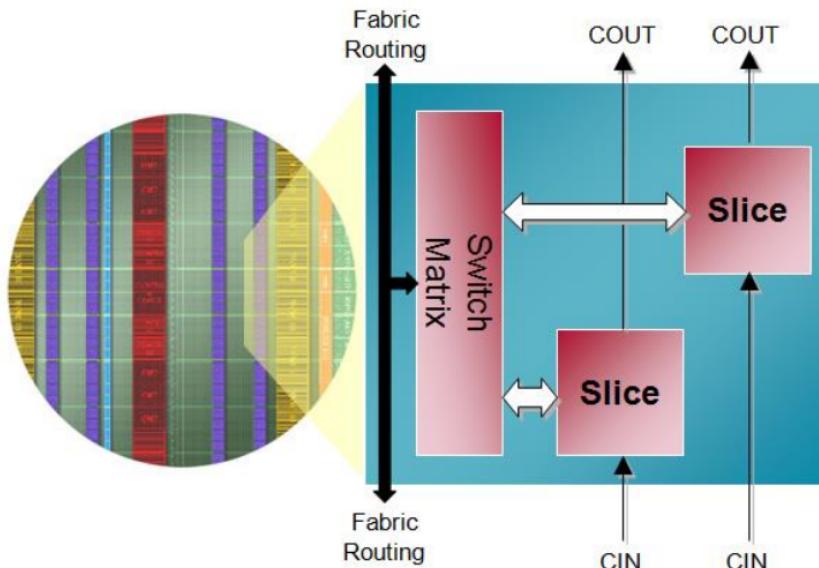
- CLB
- BRAM
- I/O
- CMT
- FIFO Logic
- ▶ BUFG
- × + DSP
- △ ▲ △ ▽ BUFIQ & BUFR
- MGT



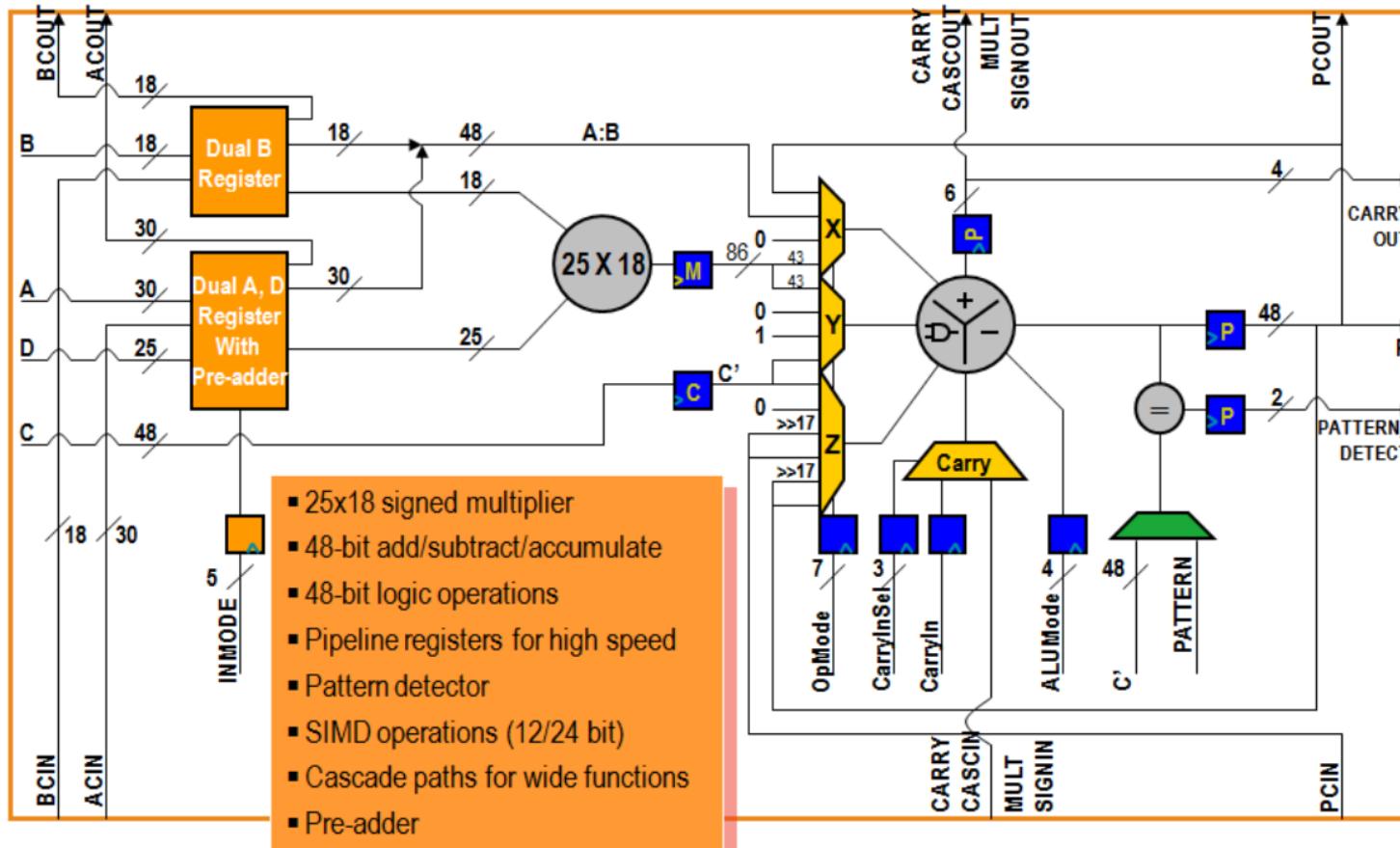
Artix-7 Architecture Overview

Xilinx 7 Series FPGA

- Configurable Logic Block
 - Combinatorial functions
 - Flip-Flops
- Slices
 - Four six-input LUTs
 - Multiplexers
 - Carry Chains
 - Four Flip-Flops/Latches (plus another 4)



Xilinx 7 Series FPGA



Remarks: 3-stage pipeline

Scalar Multiplier

- Besides VHDL version, we designed two more implementations to study resource allocations.
 - Verilog – small latency/area intensive
 - Verilog - DSP

Verilog: Area intensive vs DSP

- Area:

```
always @(scalar,in) begin
    if (~scalar[i])
        partial[i] <= 16'h0;
    else begin
        partial[i] <= tmp_in>>(TMP_WIDTH-i);
    end
end

assign tmp_out =
    (partial[15]<<15)
+ (partial[14]<<14)
+ (partial[13]<<13)
+ (partial[12]<<12)
+ (partial[11]<<11)
+ (partial[10]<<10)
+ (partial[9]<<9)
+ (partial[8]<<8)
+ (partial[7]<<7)
+ (partial[6]<<6)
+ (partial[5]<<5)
+ (partial[4]<<4)
+ (partial[3]<<3)
+ (partial[2]<<2)
+ (partial[1]<<1)
+ partial[0];
```

Note: i indexes bits in 16-bit value

- DSP

```
always @(*s_scalar, s_in_vector[i]) begin
    s_out_vector[i] <= s_in_vector[i] * s_scalar;
end
```

Note: i index 16-bit value in vector

Verilog: Area intensive vs DSP

- Area:

Slice Logic Utilization:

Number of Slice LUTs: 3288 out of 63400 5%
Number used as Logic: 3288 out of 63400 5%

Slice Logic Distribution:

Number of LUT Flip Flop pairs used: 3288
Number with an unused Flip Flop: 3288 out of 3288 100%
Number with an unused LUT: 0 out of 3288 0%
Number of fully used LUT-FF pairs: 0 out of 3288 0%

Multiplexers : 968
1-bit 2-to-1 multiplexer : 968
DSP
Number of DSP48E1s: 0 out of 240 0%

- DSP

Slice Logic Distribution:

Number of LUT Flip Flop pairs used: 0
Number with an unused Flip Flop: 0 out of 0
Number with an unused LUT: 0 out of 0

Multiplexers = : 0
1-bit 2-to-1 multiplexer : 0

DSP
Number of DSP48E1s: 24 out of 240 10%

Results

XNOR-NET Results

Name	Value	15 ns	20 ns	25 ns	30 ns	35 ns
► weights[7:0]	[01,33,45,67,89,ab,cd,ef]	[01,33,45,67,89,ab,cd,ef]	[22,55,66,88,11,12,34,71]	[55,33,45,64,89,ab,cd,ef]		
► nn_in[7:0]	63	63	a9		78	
► alpha[8:-7]	0280	0280	fe80		0580	
► beta[8:-7]	0200	0200	0100		0800	
► nn_after_p	[00020000,00040000,000...	[00000000,00fc0000,00fc0000,...	[00000000,00000000,00f...			
► nn_after_a	[00050000,00050000,000...	[00000000,ff400180,ff400180,ff...	[00000000,00000000,02b...			
▼ nn_out[7:0]	[000a0000,00140000,000...	[0000a0000,00400100,00400100...	[00000000,00000000,fdf...			
► [7]	000a0000	000a0000		00000000		
► [6]	00140000	00140000	00400100		00000000	
► [5]	000a0000	000a0000	00400100		fddfd000	
► [4]	001e0000	001e0000	ffffa000		00580000	
► [3]	007ffa00	007ffa00	00000000		fddfd000	
► [2]	000a0000	000a0000	00400100		fddfd000	
► [1]	007ffa00	007ffa00	00400100		fddfd000	
► [0]	000a0000	000a0000	00000000		fddfd000	
clk	0					
period	10000 ps		10000 ps			

Python Verification of XNOR-Net

```
import numpy as np

# Hyperparameters
batch_size = 8

# FC layer weight matrix
weights = np.asarray([0x01,0x33,0x45,0x67,0x89,0xAB,0xCD,0xEF])

# FC input after binarization
nn_in = 0x63
nn_out = 0x0

# Scalars after binarization
alpha = 2.5
beta = 2.0

# XNOR operation
xnor_results = np.zeros(batch_size)
for idx in range(batch_size):
    xnor_results[idx] = (nn_in ^ weights[idx])
print("xnor_results=", xnor_results)

# POPCOUNT operation
pop_results = np.zeros(batch_size)
for idx in range(batch_size):
    #print(str(bin(int(xnor_results[idx]))))
    num_ones = 8 - int(str(bin(int(xnor_results[idx]))).count("1"))
    pop_results[idx] = 2*num_ones-8
print("pop_results = ", pop_results)

# Scalar multiplication (alpha)
nn_results = np.zeros(batch_size)
for idx in range(batch_size):
    nn_results[idx] = pop_results[idx] * alpha
print("nn_results after alpha multiplication= ", nn_results)

# Scalar multiplication (beta)
for idx in range(batch_size):
    nn_results[idx] = nn_results[idx] * beta
print("nn_results after beta multiplication= ", nn_results)
```

Name	Value	15 ns
► weights[7:0]	[01,33,45,67,89,ab,cd,ef]	[01,33,45,67,89,ab,cd,ef]
► nn_in[7:0]	63	63
► alpha[8:-7]	0280	0280
► beta[8:-7]	0200	0200
► nn_after_p	[00020000,00040000,000...	[00020000,00040000,000...
► nn_after_a	[00050000,000a0000,000...	[00050000,000a0000,000...
▼ nn_out[7:0]	[000a0000,00140000,000...	[000a0000,00140000,000...
► [7]	000a0000	000a0000
► [6]	00140000	00140000
► [5]	000a0000	000a0000
► [4]	001e0000	001e0000
► [3]	007ffa00	007ffa00
► [2]	000a0000	000a0000
► [1]	007ffa00	007ffa00
► [0]	000a0000	000a0000
clk	0	
period	10000 ps	

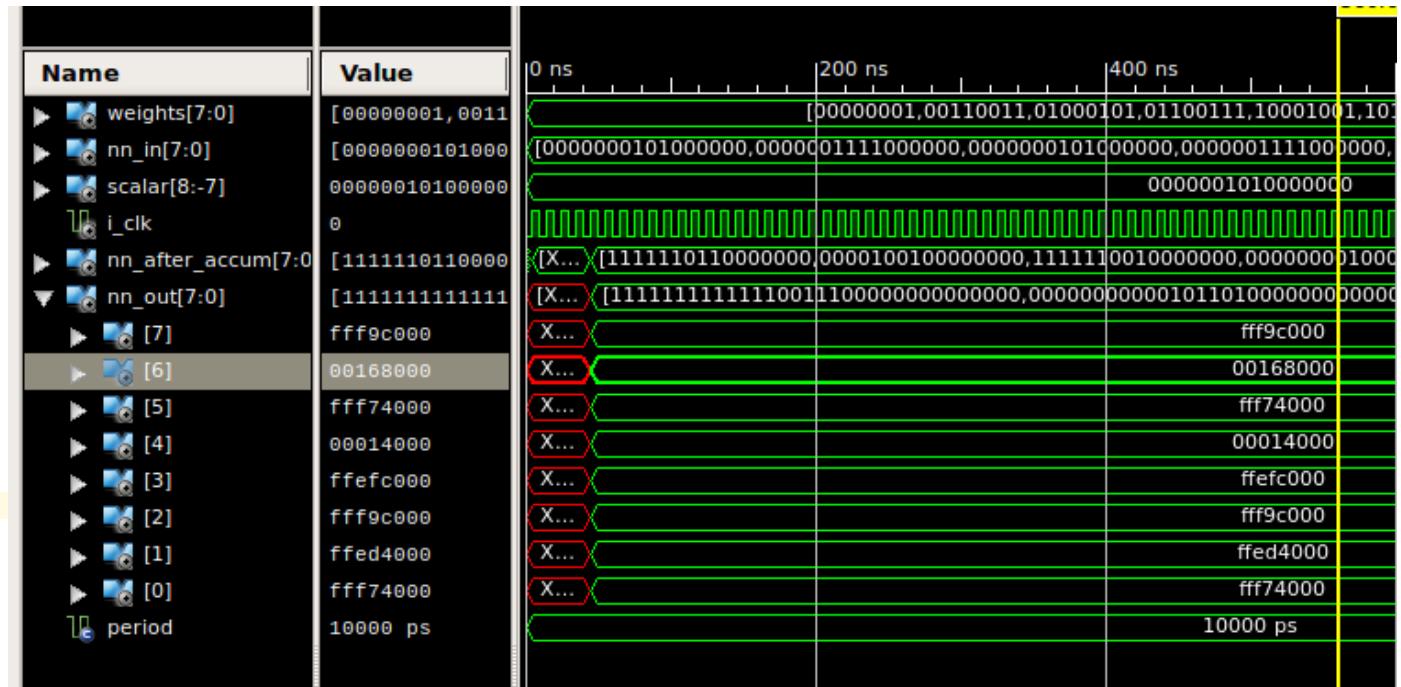
```
>>> reload(xnor_net)
('xnor_results=', array([ 98.,   80.,   38.,    4.,  234.,  200.,  174.,  140.]))
('pop_results = ', array([ 2.,   4.,   2.,   6.,  -2.,   2.,  -2.,   2.]))
('nn_results after alpha multiplication= ', array([ 5.,  10.,   5.,  15.,  -5.,   5.,  -5.,   5.]))
('nn_results after beta multiplication= ', array([ 10.,  20.,  10.,  30., -10.,  10., -10.,  10.]))
<module 'xnor_net' from 'xnor_net.pyc'>
```

Binary Connect Results

Name	Value	0 ns	200 ns	400 ns	600 ns	800 ns	1000 ns
► weights[7:0]	[00000001, 00110011, 01000101, 011001]	[00000001, 00110011, 01000101, 01100111, 10001001, 10101011, 11001101, 11101111]					
► nn_in[7:0]	[0000000101000000, 0000001111000000]	[0000000101000000, 0000001111000000, 0000000101000000, 0000001111000000, 0000000111100000, 0000001111110...					
► scalar[8:-7]	0000001010000000			0000001010000000			
► i_clk	0						
► nn_after_accum[7:0]	[1111110110000000, 0000100100000000]	[...][1111110110000000, 0000100100000000, 1111110110000000, 0000000010100000, 0000000010100000, 0000000010100000, 1111...					
▼ nn_out[7:0]	[111111111111001110000000000000, 0000000000000000]	[...][111111111111001110000000000000, 0000000000000000, 0000000000000000101101000000000000000000, 1111...					
► [7]	111111111111001110000000000000	...	111111111111001110000000000000				
► [6]	000000000010110100000000000000	...		000000000000101101000000000000			
► [5]	111111111110111010000000000000	...		111111111110111010000000000000			
► [4]	000000000000001010000000000000	...		000000000000001010000000000000			
► [3]	111111111110111110000000000000	...		111111111110111110000000000000			
► [2]	11111111111011110011100000000000	...		11111111111011110011100000000000			
► [1]	11111111111011101010000000000000	...		11111111111011101010000000000000			
► [0]	11111111111011101000000000000000	...		11111111111011101000000000000000			
► period	10000 ps			10000 ps			

Python Verification of Binary Connect

```
1 import numpy as np
2 from decimal import *
3
4 # Hyperparameters
5 batch_size = 8
6
7 # FC layer weight matrix
8 weights = np.asarray([0x01, 0x33, 0x45, 0x67, 0x89, 0xAB, 0xCD, 0xEF])
9
10 # FC input after binarization
11 nn_in = [1.25, 3.75, 1.25, 3.75, -3.25, -4.25, 0.75, 0.75]
12 nn_out = list(np.zeros(8))
13
14 # Scalars after binarization
15 alpha = 2.5
16
17 # accumulation
18 weights = [format(x, '#010b')[2:] for x in weights]
19 for x in range(0, len(weights)):
20     for y in range(0, len(weights[x])):
21         if weights[x][y] == '1':
22             nn_out[x] += nn_in[y]
23         else:
24             nn_out[x] -= nn_in[y]
25
26 print("Matrix vector multiplication after accumulator = ", nn_out)
27
28 nn_results = np.zeros(batch_size)
29 for idx in range(batch_size):
30     nn_results[idx] = nn_out[idx] * alpha
31 print("nn_results after alpha multiplication = ", nn_results)
32
```



```
Matrix vector multiplication after accumulator = [-2.5, 9.0, -3.5, 0.5, -6.5, -2.5, -7.5, -3.5]
nn_results after alpha multiplication = [ -6.25  22.5   -8.75    1.25 -16.25  -6.25 -18.75  -8.75]
```