

VIETNAM NATIONAL UNIVERSITY HO CHI MINH
CITY
HO CHI MINH CITY UNIVERSITY OF TECHNOLOGY



REPORT

Logic Design with HDL

Topic:
MATRIX MULTIPLICATION

Instructor: Mr. Nguyễn Thành Lộc

Group No: 6

Class: CC02

Completion Day May 29, 2025

Group members: Lê Lương Trung Hiếu (2452330)
Trần Nguyễn Nguyên Khoa (2452563)
Hoàng Trọng Huy Minh (2452743)
Lê Hoàng Phúc (2452990)
Cao Thiên (2453192)



Contents

1	Chapter 1: Abstract and Introduction	2
2	Chapter 2: Backgrounds and applications	3
3	Chapter 3: Design	5
3.1	Ideal Design – Matrix Multiplication System	5
3.2	Block diagram	6
3.3	Flow chart	10
3.4	Algorithms	12
4	Chapter 4: Implement	13
4.1	Module: Baugh-Wooley multiplier	13
4.2	Booth multiplier	16
4.3	Comparison	26
5	Chapter 5: Results	29
5.1	Experiment setups	29
5.2	Waveform	29
6	Conclusion	40
7	Chapter 7: References	41

1 Chapter 1: Abstract and Introduction

Introduction:

- Efficient hardware implementation of matrix multiplication is a critical task in many computational domains, including signal processing, computer vision, and neural networks. To support signed arithmetic in these applications, high-performance signed multipliers are essential. This project explores and compares two popular signed multiplier architectures—Baugh-Wooley and Booth multipliers—through the implementation of a 4×4 matrix multiplier in Verilog HDL.
- Signed multiplication is particularly important in matrix operations where negative values are common. Two of the most prominent signed multiplication algorithms for hardware implementation are the **Baugh-Wooley multiplier** and the **Booth multiplier**. The **Baugh-Wooley algorithm** is optimized for **two's complement arithmetic** and offers a regular partial product structure, which simplifies logic design and timing closure.
- In contrast, **the Booth-based architecture** employs a sequential Booth multiplication algorithm, which processes the multiplier in 4-bit chunks. For 8-bit operands, this approach requires three clock cycles per multiplication. The control FSM includes a COMPUTE state to account for this latency, followed by an OUTPUT state where the sum of the four products is assigned. While Booth's recoding technique can reduce the number of partial products and may lead to area savings, this implementation trades off throughput by serializing each multiplication step and avoiding pipelining.
- This project presents a comparative HDL implementation of matrix multiplication using both Baugh-Wooley and Booth multipliers, written in Verilog. A 4×4 matrix multiplication module is developed, where each matrix element is an 8-bit signed number. The design incorporates four parallel multipliers to calculate each matrix element's dot product, increasing throughput. The core multiplication logic is modular, allowing for the integration of either Baugh-Wooley or Booth implementations.
- **And in this report**, using simulation and waveform analysis, we would analyze and test not only the correctness but also effectiveness of the theory and the code created based on it.

2 Chapter 2: Backgrounds and applications

Background

- Matrix multiplication is a foundational operation in digital signal processing, computer vision, robotics, cryptography, and machine learning. In these domains, data is frequently represented in matrix form, and processing often involves extensive matrix computations. Hardware acceleration of such computations is critical to meet real-time performance requirements, especially when operating on large datasets or performing complex tasks.
- Signed multipliers are a key component in these operations, especially when dealing with two's complement representation—commonly used in digital systems to handle negative values. Among the various multiplier architectures, **Baugh-Wooley** and **Booth** multipliers are notable for their suitability in signed arithmetic. Each architecture offers different trade-offs in terms of speed, area, power consumption, and design complexity.
- The **Baugh-Wooley multiplier** leverages the regularity of two's complement arithmetic to simplify the generation of partial products and improve timing performance. Its highly parallel structure makes it well-suited for high-throughput applications.
- The **Booth multiplier**, on the other hand, reduces the number of partial products through recoding, potentially leading to reduced hardware complexity. However, its sequential nature and dependency on control logic (e.g., FSM) can limit throughput if not properly pipelined.
- In Verilog-based digital system design, understanding the strengths and weaknesses of these architectures is essential for choosing the most appropriate multiplier for a given application.

Applications

Efficient signed matrix multiplication has wide-ranging applications, such as:

- **Digital Signal Processing (DSP):** Filters, transforms (e.g., DFT/FFT), and correlation algorithms rely heavily on fast matrix operations.
- **Machine Learning and Neural Networks:** Matrix multiplication forms the core of operations in dense layers, convolutional layers, and attention mechanisms. Signed multipliers are crucial when using quantized signed data.
- **Image and Video Processing:** Operations such as convolution, edge detection, and color space transformations use signed matrix multiplications extensively.
- **Control Systems and Robotics:** State estimation, sensor fusion, and motion planning often involve signed matrix arithmetic.



- **Embedded Systems and ASIC/FPGA Accelerators:** For real-time, low-power applications, implementing matrix multiplication in hardware using efficient signed multipliers is vital.

By exploring and comparing the Baugh-Wooley and Booth multipliers in the context of a 4×4 matrix multiplier, this project addresses both the theoretical and practical aspects of high-performance arithmetic circuit design. The results can inform design decisions in a wide range of applications that depend on signed, high-throughput, and resource-efficient matrix multiplications.

3 Chapter 3: Design

3.1 Ideal Design – Matrix Multiplication System

The system is designed to compute the product of two 4×4 signed integer matrices, A and B, using Verilog HDL. Each matrix signal has a size of 8 bits, in which MSB will be the sign bit and the remaining 7 bits will be the magnitude). The result matrix C will contain 17-bit signed numbers, is generated using a custom-designed multiply-accumulate process controlled by an FSM.

The system operates in **four main stages**:

3.1.1 Input Interface and Data Storage

- The input signal `din` port. receives 32 signed 8-bit integers sequentially, where:
 - The first 16 values are stored into matrix A.
 - The next 16 values are stored into matrix B.
- The signal `wrt_en` is used to trigger data storage at each clock cycle. Two register arrays (`A[0:15]` and `B[0:15]`) are used to store the matrices.

3.1.2 Control Unit (FSM Controller)

- The process is managed by a finite state machine (FSM) with four states:
 - `IDLE`: Waiting for data input.
 - `LOAD_A`: Loads matrix A.
 - `LOAD_B`: Loads matrix B.
 - `OUTPUT`: Performs matrix multiplication and outputs results.
- State transitions are triggered by `wrt_en` and internal counters (`load_count`, `output_count`).

3.1.3 Multiply-Accumulate Core

- the system uses four parallel instances of a custom 8×8 signed multiplier module, `baugh_wooley_multiplier`, based on the Baugh–Wooley algorithm. These modules compute the partial products:
- For each result element `C[i][j]`, the system calculates:

$$C[i][j] = \sum_{k=0}^3 A[i][k] \times B[k][j]$$

- The partial results are then summed to form the final output.

3.1.4 Output Handling

- Once the computation begins, the system outputs one result per clock cycle via the 17-bit `dout` port. After 16 results, the `done` signal is raised to indicate that the full 4×4 matrix multiplication is complete.

3.1.5 Design Highlights

- **Custom multiplier:** The Baugh–Wooley multiplier is implemented manually to support signed multiplication efficiently in hardware.
- **Parallelism:** Four multipliers run concurrently to increase throughput.
- **Simplicity:** The design does not yet apply pipelining, prioritizing clarity and functional correctness.
- **Expandability:** Can be upgraded to pipelined architecture or reused in larger matrix multipliers.

3.2 Block diagram

The matrix multiplier system consists of the following functional blocks:

- **Input Register:** Receives 8-bit signed data via the `din` port and stores it in matrix A or B based on `wrt_en` control.
- **Matrix A Storage:** Holds 16 elements of matrix A, each 8-bit signed.
- **Matrix B Storage:** Holds the next 16 elements for matrix B.
- **FSM Controller:** A finite state machine that manages the stages: loading A, loading B, computing, and outputting.
- **MAC Core (Multiply and Accumulate):** Contains four parallel custom Baugh–Wooley multipliers to compute $A[i][k] \times B[k][j]$ and accumulate the result.
- **Output Register:** Stores and outputs each 17-bit result `C[i][j]` sequentially via the `dout` port.

Input/Output signals:

- `clk, rst` : Clock and synchronous reset signals.
- `din [7:0]` : 8-bit signed input data.
- `wrt_en` : Write enable signal for loading data.
- `dout [16:0]` : 17-bit signed result output.
- `done` : Indicates completion of the computation.

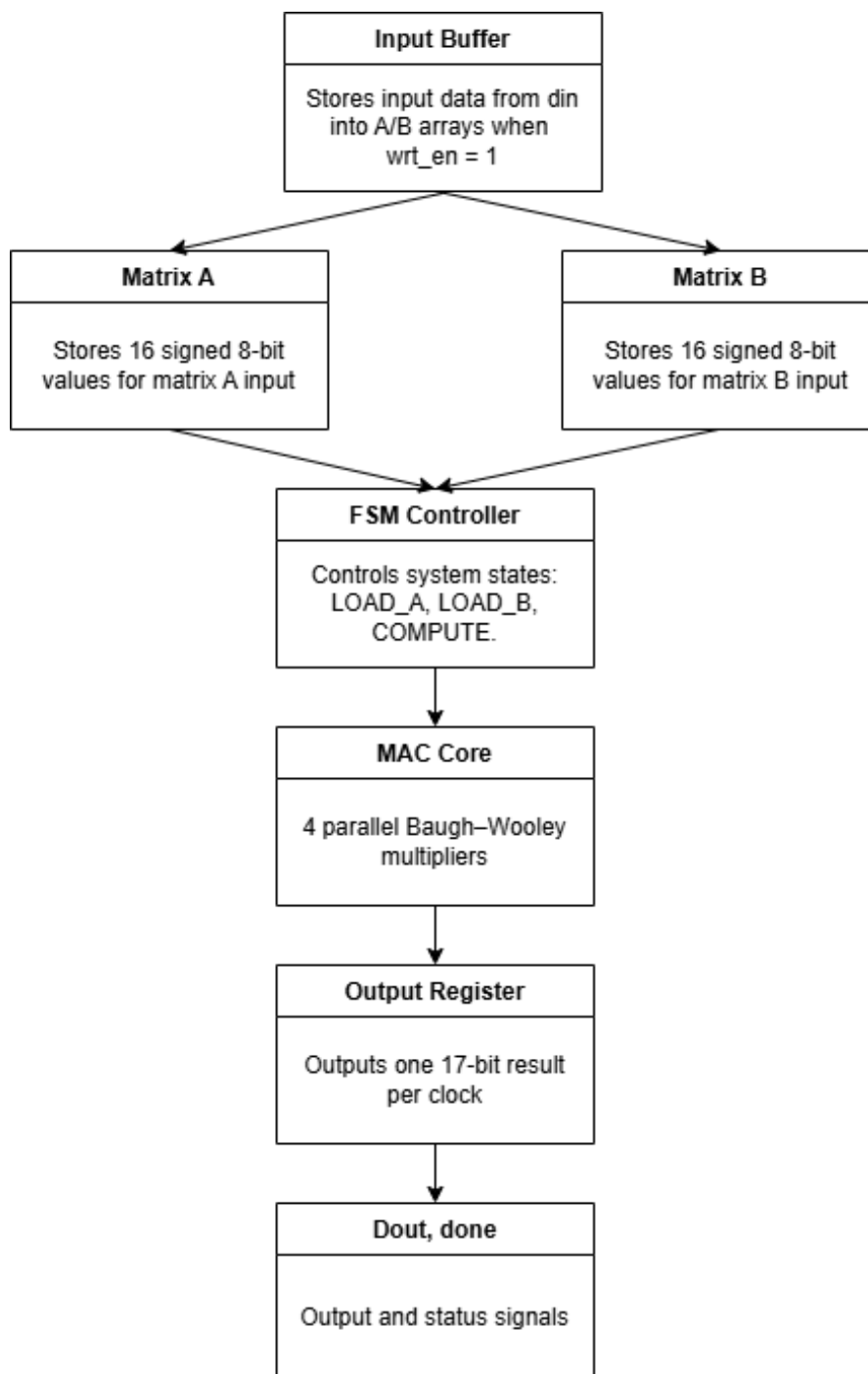


Figure 1: The Matrix Multiplier System Blocks

The Booth multiplier system consists of the following functional blocks:

- **Input Registers:** Receive two signed 8-bit inputs — **M** (multiplicand) and **R** (multiplier).
- **Register Setup:** The multiplicand is sign-extended to form **A**; the multiplier remains in **R**.
- **Booth Encoding Logic:** Analyzes the multiplier bits using the Booth algorithm to generate encoding patterns that reduce the number of partial products.
- **Control Signal Generator:** Produces control signals such as **PnM**, **M_Sel**, and **En** to guide partial product selection and operations.
- **Partial Product Generator:** Generates multiple shifted versions of the multiplicand: **Mx1**, **Mx2**, **Mx4**, and **Mx8**.
- **Booth Operations:** Includes B-op and C-op blocks that select and align partial products based on the Booth encoding and control logic.
- **Adder Tree Stages:** Multi-level adder structure that computes intermediate sums using carry-in and high-bit propagation: $T = H_i + B + C_{i_B}$, followed by $S = T + C + C_{i_C}$.
- **Guard Logic and Shift Bit:** Ensures correct bit-width alignment and overflow handling during shifting.
- **Shift and Accumulate:** Performs iterative shifting and accumulation with rounding/guard logic.
- **Final Result:** The final 16-bit signed product **P** is produced after all accumulation steps.

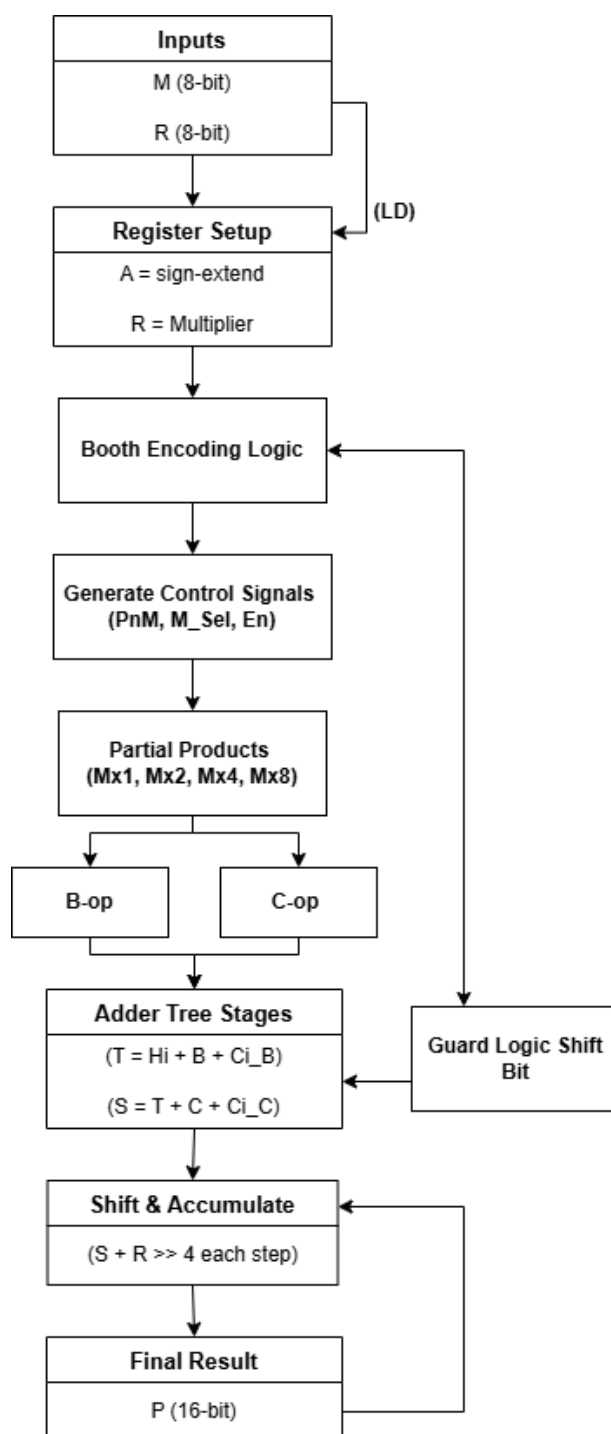


Figure 2: Booth Multiplier Functional Blocks

3.3 Flow chart

To coordinate the operation of the matrix multiplier system, a finite state machine (FSM) is implemented. The FSM consists of four main states:

- **IDLE:** Waits for the `wrt_en` signal to begin data loading.
- **LOAD_A:** Receives and stores the first 16 input values into Matrix A.
- **LOAD_B:** Continues receiving the next 16 values for Matrix B.
- **COMPUTE:** Performs the matrix multiplication using four parallel multipliers and outputs one result per clock cycle.

The FSM transitions between these states based on internal counters and input control signals. Once all 16 results are computed, the `done` signal is asserted, and the FSM returns to the IDLE state.

Figure 3 shows the overall control flow of the FSM.

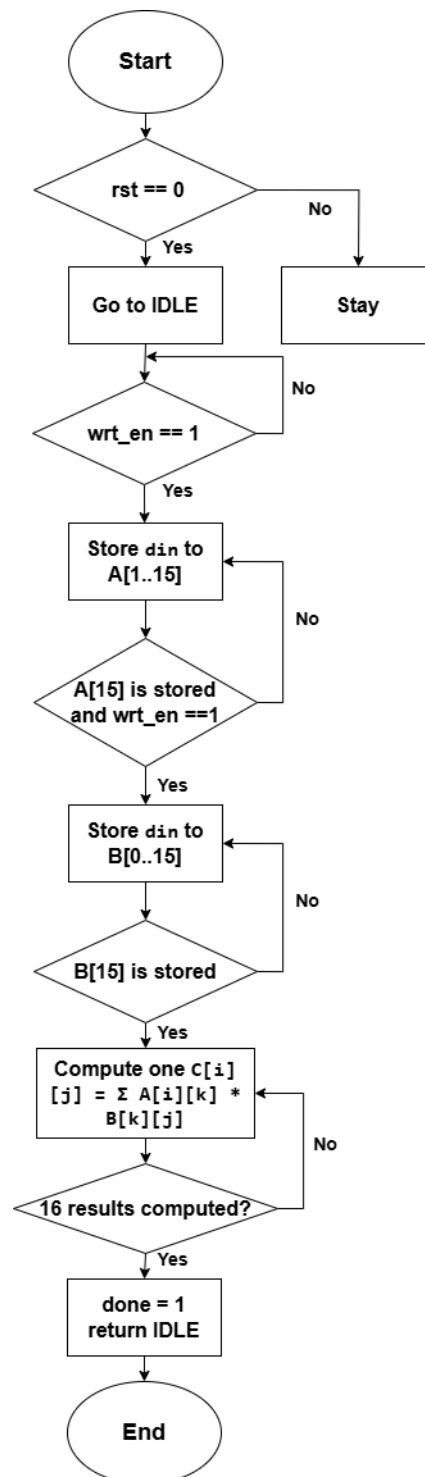


Figure 3: FSM flowchart for matrix multiplier control logic

3.4 Algorithms

Algorithm 1: FSM Controller for 4x4 Matrix Multiplication

Input: 32 signed 8-bit values via `din`, controlled by `wrt_en`

Output: 16 signed 17-bit values via `dout`

```
1 Initialize: state  $\leftarrow$  IDLE, load_count  $\leftarrow$  0, output_count  $\leftarrow$  0
2 while true do
3   if rst == 1 then
4     | Reset all states and counters
5   end
6   else
7     if state == IDLE and wrt_en then
8       | Store din into A[0], load_count  $\leftarrow$  1, state  $\leftarrow$  LOAD_A
9     if state == LOAD_A and wrt_en then
10      | Store din into A[load_count], increment load_count
11      if load_count == 16 then
12        | state  $\leftarrow$  LOAD_B
13      end
14    if state == LOAD_B and wrt_en then
15      | Store din into B[load_count - 16], increment load_count
16      if load_count == 32 then
17        | output_count  $\leftarrow$  0, state  $\leftarrow$  COMPUTE
18      end
19    if state == COMPUTE then
20      | Compute  $C[i][j] = \sum A[i][k] \times B[k][j]$ 
21      | Output to dout, increment output_count
22      if output_count == 16 then
23        | Set done  $\leftarrow$  1, state  $\leftarrow$  IDLE
24      end
25    end
26 end
```

4 Chapter 4: Implement

4.1 Module: Baugh-Wooley multiplier

```
1 'timescale 1ns / 1ps
2
3 module Matrix_multi(
4     input wire clk, rst, wrt_en,
5     input wire signed [7:0] din,
6     output reg signed [16:0] dout
7 );
8     reg signed [7:0] A [0:15];
9     reg signed [7:0] B [0:15];
10
11     reg [2:0] state;
12     reg [5:0] load_count;
13     reg [3:0] output_count;
14
15     localparam IDLE = 3'd0,
16                 LOAD_A = 3'd1,
17                 LOAD_B = 3'd2,
18                 OUTPUT = 3'd3;
19
20     wire [1:0] row = output_count[3:2];
21     wire [1:0] col = output_count[1:0];
22
23     wire signed [15:0] p0, p1, p2, p3;
24     wire signed [16:0] sum;
25
26     baugh_wooley_multiplier bw0(.a(A[{row,2'b00}]), .b(B[{2'b00,col}]), .p(p0));
27     baugh_wooley_multiplier bw1(.a(A[{row,2'b01}]), .b(B[{2'b01,col}]), .p(p1));
28     baugh_wooley_multiplier bw2(.a(A[{row,2'b10}]), .b(B[{2'b10,col}]), .p(p2));
29     baugh_wooley_multiplier bw3(.a(A[{row,2'b11}]), .b(B[{2'b11,col}]), .p(p3));
30
31     assign sum = p0 + p1 + p2 + p3;
32
33     always @(posedge clk) begin
34         if (rst) begin
35             state <= IDLE;
36             load_count <= 0;
37             output_count <= 0;
38             dout <= 0;
39         end else begin
40             case (state)
41                 IDLE:
42                     if (wrt_en) begin
43                         A[0] <= din;
44                         load_count <= 1;
45                         state <= LOAD_A;
46                     end
47                 LOAD_A:
48                     if (wrt_en) begin
49                         A[load_count] <= din;
```

```

50         load_count <= load_count + 1;
51         if (load_count == 15)
52             state <= LOAD_B;
53         end
54     LOAD_B:
55         if (wrt_en) begin
56             B[load_count - 16] <= din;
57             load_count <= load_count + 1;
58             if (load_count == 31) begin
59                 state <= OUTPUT;
60                 output_count <= 0;
61             end
62         end
63     OUTPUT:
64         begin
65             dout <= sum;
66             output_count <= output_count + 1;
67             if (output_count == 15)
68                 state <= IDLE;
69         end
70     endcase
71 end
72 end
73 endmodule

```

Listing 1: Matrix_multi Module with Baugh-Wooley Multipliers

```

1 module baugh_wooley_multiplier(
2     input wire signed [7:0] a,
3     input wire signed [7:0] b,
4     output reg signed [15:0] p
5 );
6     integer i, j;
7     reg signed [15:0] acc;
8     reg bit;
9
10    always @* begin
11        acc = 16'sd0;
12
13        for (i = 0; i <= 6; i = i + 1) begin
14            for (j = 0; j <= 6; j = j + 1) begin
15                bit = a[i] & b[j];
16                acc = acc + (bit << (i + j));
17            end
18        end
19
20        for (j = 0; j <= 6; j = j + 1) begin
21            bit = a[7] & b[j];
22            acc = acc - (bit << (7 + j));
23        end
24
25        for (i = 0; i <= 6; i = i + 1) begin
26            bit = a[i] & b[7];
27            acc = acc - (bit << (i + 7));

```

```
28     end
29
30     bit = a[7] & b[7];
31     acc = acc + (bit << 14);
32
33     p = acc;
34 end
35 endmodule
```

Listing 2: Baugh-Wooley Multiplier Module

```
1  `timescale 1ns / 1ps
2  module tb_Matrix_multi_scaled_time;
3
4      // Inputs
5      reg clk, rst, wrt_en;
6      reg signed [7:0] din;
7
8      // Outputs
9      wire signed [16:0] dout;
10
11     // Internal signals for monitoring
12     wire [2:0] state_monitor;
13     wire [3:0] output_count_monitor;
14
15     // Instantiate the Unit Under Test (UUT)
16     Matrix_multi uut (
17         .clk(clk),
18         .rst(rst),
19         .wrt_en(wrt_en),
20         .din(din),
21         .dout(dout)
22     );
23
24     // Expose internal signals for monitoring
25     assign state_monitor = uut.state;
26     assign output_count_monitor = uut.output_count;
27
28     // Clock generation: 2ns period (for faster simulation)
29     always #1 clk = ~clk;
30
31     // Matrix and result storage
32     reg signed [7:0] matrix_A [0:15];
33     reg signed [7:0] matrix_B [0:15];
34     integer i;
35
36     initial begin
37         // Initialize signals
38         clk = 0;
39         rst = 1;
40         wrt_en = 0;
41         din = 0;
42
43         // Brief reset (2ns)
```



```
44     #2 rst = 0;
45
46     // Initialize matrix_A
47     matrix_A[0] = 8'd3; matrix_A[1] = 8'd7; matrix_A[2] = 8'd12; matrix_A[3] = 8'd2;
48     matrix_A[4] = 8'd1; matrix_A[5] = 8'd4; matrix_A[6] = 8'd6; matrix_A[7] = 8'd8;
49     matrix_A[8] = 8'd9; matrix_A[9] = 8'd5; matrix_A[10] = 8'd10; matrix_A[11] =
50         8'd11;
51     matrix_A[12] = 8'd7; matrix_A[13] = 8'd2; matrix_A[14] = 8'd4; matrix_A[15] =
52         8'd8;
53
54     // Initialize matrix_B
55     matrix_B[0] = 8'd5; matrix_B[1] = 8'd1; matrix_B[2] = 8'd3; matrix_B[3] = 8'd8;
56     matrix_B[4] = 8'd2; matrix_B[5] = 8'd7; matrix_B[6] = 8'd4; matrix_B[7] = 8'd6;
57     matrix_B[8] = 8'd9; matrix_B[9] = 8'd5; matrix_B[10] = 8'd2; matrix_B[11] = 8'd1;
58     matrix_B[12] = 8'd3; matrix_B[13] = 8'd6; matrix_B[14] = 8'd1; matrix_B[15] =
59         8'd7;
60
61     // --- Write matrix A ---
62     #2 wrt_en = 1;
63     for (i = 0; i < 16; i = i + 1) begin
64         din = matrix_A[i];
65         #2;
66     end
67     wrt_en = 0;
68
69     // --- Write matrix B ---
70     wrt_en = 1;
71     for (i = 0; i < 16; i = i + 1) begin
72         din = matrix_B[i];
73         #2;
74     end
75     wrt_en = 0;
76
77     // Wait for OUTPUT state
78     wait(state_monitor == 3);
79
80     // Display outputs
81     for (i = 0; i < 16; i = i + 1) begin
82         @(posedge clk);
83         $display("Output_%0d:_%0d", i, dout);
84     end
85
86     #10 $finish;
87 end
88 endmodule
```

Listing 3: Testbench for Matrix_multi (Scaled Time)

4.2 Booth multiplier

```
1 module Matrix_multi(
```

```
2   input wire clk, rst, wrt_en,
3   input wire signed [7:0] din,
4   output reg signed [16:0] dout
5 );
6   reg signed [7:0] A [0:15];
7   reg signed [7:0] B [0:15];
8
9   reg [2:0] state;
10  reg [5:0] load_count;
11  reg [3:0] output_count;
12  reg [1:0] compute_count; // Counter for 3-cycle multiplication latency
13
14  localparam IDLE = 3'd0,
15             LOAD_A = 3'd1,
16             LOAD_B = 3'd2,
17             COMPUTE = 3'd3,
18             OUTPUT = 3'd4;
19
20  wire [1:0] row = output_count[3:2];
21  wire [1:0] col = output_count[1:0];
22
23  wire signed [15:0] p0, p1, p2, p3;
24  wire valid0, valid1, valid2, valid3;
25  reg ld; // Load signal for Booth multipliers
26  wire signed [16:0] sum;
27
28  // Instantiate four Booth multipliers
29  Booth_Multiplier_4xA #(
30      .N(8)
31  ) bw0 (
32      .Rst(rst),
33      .Clk(clk),
34      .Ld(ld),
35      .M(A[{row,2'b00}]),
36      .R(B[{2'b00,col}]),
37      .Valid(valid0),
38      .P(p0)
39  );
40  Booth_Multiplier_4xA #(
41      .N(8)
42  ) bw1 (
43      .Rst(rst),
44      .Clk(clk),
45      .Ld(ld),
46      .M(A[{row,2'b01}]),
47      .R(B[{2'b01,col}]),
48      .Valid(valid1),
49      .P(p1)
50  );
51  Booth_Multiplier_4xA #(
52      .N(8)
53  ) bw2 (
54      .Rst(rst),
```

```
55     .Clk(clk),
56     .Ld(ld),
57     .M(A[{row,2'b10}]),
58     .R(B[{2'b10,col}]),
59     .Valid(valid2),
60     .P(p2)
61 );
62 Booth_Multiplier_4xA #(
63     .N(8)
64 ) bw3 (
65     .Rst(rst),
66     .Clk(clk),
67     .Ld(ld),
68     .M(A[{row,2'b11}]),
69     .R(B[{2'b11,col}]),
70     .Valid(valid3),
71     .P(p3)
72 );
73
74 assign sum = p0 + p1 + p2 + p3;
75
76 always @(posedge clk) begin
77     if (rst) begin
78         state <= IDLE;
79         load_count <= 0;
80         output_count <= 0;
81         compute_count <= 0;
82         ld <= 0;
83         dout <= 0;
84     end else begin
85         case (state)
86             IDLE: begin
87                 if (wrt_en) begin
88                     A[0] <= din;
89                     load_count <= 1;
90                     state <= LOAD_A;
91                 end
92             end
93             LOAD_A: begin
94                 if (wrt_en) begin
95                     A[load_count] <= din;
96                     load_count <= load_count + 1;
97                     if (load_count == 15) state <= LOAD_B;
98                 end
99             end
100             LOAD_B: begin
101                 if (wrt_en) begin
102                     B[load_count-16] <= din;
103                     load_count <= load_count + 1;
104                     if (load_count == 31) begin
105                         state <= COMPUTE;
106                         output_count <= 0;
107                         compute_count <= 0;
```

```

108         ld <= 1; // Start first multiplication
109     end
110 end
111 end
112 COMPUTE: begin
113     ld <= 0; // Pulse ld for one cycle
114     compute_count <= compute_count + 1;
115     if (compute_count == 2) begin // 3 cycles total (0,1,2)
116         state <= OUTPUT;
117         compute_count <= 0;
118     end
119 end
120 OUTPUT: begin
121     dout <= sum;
122     output_count <= output_count + 1;
123     ld <= 1; // Start next multiplication
124     if (output_count == 15) begin
125         state <= IDLE;
126         ld <= 0;
127     end else begin
128         state <= COMPUTE; // Return to COMPUTE for next element
129     end
130 end
131 endcase
132 end
133 end
134 endmodule

```

Listing 4: Matrix_multi Module Implementation for Booth module

```

1 module Booth_Multiplier_4xA #(
2     parameter N = 8
3 )
4     input Rst,
5     input Clk,
6     input Ld,
7     input [(N - 1):0] M,
8     input [(N - 1):0] R,
9     output reg Valid,
10    output reg [((2*N) - 1):0] P
11 );
12
13    localparam pNumCycles = ((N + 1)/4);
14
15    reg [4:0] Cntr;
16    reg [4:0] Booth;
17    reg Guard;
18    reg [(N + 3):0] A;
19    wire [(N + 3):0] Mx8, Mx4, Mx2, Mx1;
20    reg PnM_B, M_Sel_B, En_B;
21    reg PnM_C, M_Sel_C, En_C;
22    wire [(N + 3):0] Hi;
23    reg [(N + 3):0] B, C;
24    reg Ci_B, Ci_C;

```

```
25 wire [(N + 3):0] T, S;
26 reg [((2*N) + 3):0] Prod;
27
28 always @(posedge Clk)
29 begin
30     if(Rst)
31         Cntr <= #1 0;
32     else if(Ld)
33         Cntr <= #1 pNumCycles;
34     else if(!Cntr)
35         Cntr <= #1 (Cntr - 1);
36 end
37
38 always @(posedge Clk)
39 begin
40     if(Rst)
41         A <= #1 0;
42     else if(Ld)
43         A <= #1 {{4{M[(N - 1)]}}, M};
44 end
45
46 assign Mx8 = {A, 3'b0};
47 assign Mx4 = {A, 2'b0};
48 assign Mx2 = {A, 1'b0};
49 assign Mx1 = A;
50
51 always @(*) Booth <= {Prod[3:0], Guard};
52 assign Hi = Prod[((2*N) + 3):N];
53
54 // Control logic for operand B
55 always @(*)
56 begin
57     case(Booth)
58         5'b00000, 5'b00001, 5'b00010, 5'b00011, 5'b00100,
59         5'b11011, 5'b11100, 5'b11101, 5'b11110, 5'b11111 :
60             {PnM_B, M_Sel_B, En_B} <= 3'b000;
61         5'b00101, 5'b00110, 5'b00111, 5'b01000, 5'b01001, 5'b01010 :
62             {PnM_B, M_Sel_B, En_B} <= 3'b001;
63         5'b01101, 5'b01110, 5'b01111 :
64             {PnM_B, M_Sel_B, En_B} <= 3'b011;
65         5'b10000, 5'b10001, 5'b10010 :
66             {PnM_B, M_Sel_B, En_B} <= 3'b111;
67         default :
68             {PnM_B, M_Sel_B, En_B} <= 3'b101;
69     endcase
70 end
71
72 // Control logic for operand C
73 always @(*)
74 begin
75     case(Booth)
76         5'b00000, 5'b01111, 5'b10111, 5'b11000, 5'b11111 :
77             {PnM_C, M_Sel_C, En_C} <= 3'b000;
```

```
78      5'b00001, 5'b00010, 5'b01001, 5'b01010,
79      5'b10001, 5'b10010, 5'b11001, 5'b11010 :
80          {PnM_C, M_Sel_C, En_C} <= 3'b001;
81      5'b00011, 5'b00100, 5'b01011, 5'b01100 :
82          {PnM_C, M_Sel_C, En_C} <= 3'b011;
83      5'b00101, 5'b00110, 5'b01101, 5'b01110,
84      5'b10101, 5'b10110, 5'b11101, 5'b11110 :
85          {PnM_C, M_Sel_C, En_C} <= 3'b101;
86      default :
87          {PnM_C, M_Sel_C, En_C} <= 3'b111;
88      endcase
89  end
90
91  // Operand B mux logic
92  always @(*)
93  begin
94      case({PnM_B, M_Sel_B, En_B})
95          3'b001: {Ci_B, B} <= {1'b0, Mx4};
96          3'b011: {Ci_B, B} <= {1'b0, Mx8};
97          3'b101: {Ci_B, B} <= {1'b1, ~Mx4};
98          3'b111: {Ci_B, B} <= {1'b1, ~Mx8};
99          default: {Ci_B, B} <= 0;
100      endcase
101  end
102
103  // Operand C mux logic
104  always @(*)
105  begin
106      case({PnM_C, M_Sel_C, En_C})
107          3'b001: {Ci_C, C} <= {1'b0, Mx1};
108          3'b011: {Ci_C, C} <= {1'b0, Mx2};
109          3'b101: {Ci_C, C} <= {1'b1, ~Mx1};
110          3'b111: {Ci_C, C} <= {1'b1, ~Mx2};
111          default: {Ci_C, C} <= 0;
112      endcase
113  end
114
115  assign T = Hi + B + Ci_B;
116  assign S = T + C + Ci_C;
117
118  always @(posedge Clk)
119  begin
120      if(Rst)
121          Prod <= #1 0;
122      else if(Ld)
123          Prod <= #1 R;
124      else if(!Cntr)
125          Prod <= #1 {{4{S[(N + 3)]}}, S, Prod[(N - 1):4]};
126  end
127
128  always @(posedge Clk)
129  begin
130      if(Rst)
```

```
131     Guard <= #1 0;
132     else if(Ld)
133         Guard <= #1 0;
134     else if(!Cntr)
135         Guard <= #1 Prod[3];
136 end
137
138 always @(posedge Clk)
139 begin
140     if(Rst)
141         P <= #1 0;
142     else if(Cntr == 1)
143         P <= #1 {S, Prod[(N - 1):4]};
144 end
145
146 always @(posedge Clk)
147 begin
148     if(Rst)
149         Valid <= #1 0;
150     else
151         Valid <= #1 (Cntr == 1);
152 end
153 endmodule
```

Listing 5: Booth Multiplier 4xA

```
1  `timescale 1ns / 1ps
2
3  module tb_Matrix_multi_scaled_time;
4
5      // Inputs
6      reg clk, rst, wrt_en;
7      reg signed [7:0] din;
8
9      // Outputs
10     wire signed [16:0] dout;
11
12     // Internal signals for monitoring
13     wire [2:0] state_monitor;
14     wire [3:0] output_count_monitor;
15
16     // Instantiate the Unit Under Test (UUT)
17     Matrix_multi uut (
18         .clk(clk),
19         .rst(rst),
20         .wrt_en(wrt_en),
21         .din(din),
22         .dout(dout)
23     );
24
25     // Expose internal signals for monitoring
26     assign state_monitor = uut.state;
27     assign output_count_monitor = uut.output_count;
28
```

```
29 // Clock generation: 2ns period (for faster simulation)
30 always #1 clk = ~clk;
31
32 // Matrix and result storage
33 reg signed [7:0] matrix_A [0:15];
34 reg signed [7:0] matrix_B [0:15];
35 integer i;
36
37 initial begin
38     // Initialize signals
39     clk = 0;
40     rst = 1;
41     wrt_en = 0;
42     din = 0;
43
44     // Brief reset (2ns)
45     #2 rst = 0;
46
47     // Initialize matrix_A with 1- and 2-digit decimal numbers
48     matrix_A[0] = 8'd3; matrix_A[1] = 8'd7; matrix_A[2] = 8'd12; matrix_A[3] = 8'd2;
49     matrix_A[4] = 8'd1; matrix_A[5] = 8'd4; matrix_A[6] = 8'd6; matrix_A[7] = 8'd8;
50     matrix_A[8] = 8'd9; matrix_A[9] = 8'd5; matrix_A[10] = 8'd10; matrix_A[11] =
51         8'd11;
52     matrix_A[12] = 8'd7; matrix_A[13] = 8'd2; matrix_A[14] = 8'd4; matrix_A[15] =
53         8'd8;
54
55     // Initialize matrix_B with 1- and 2-digit decimal numbers
56     matrix_B[0] = 8'd5; matrix_B[1] = 8'd1; matrix_B[2] = 8'd3; matrix_B[3] = 8'd8;
57     matrix_B[4] = 8'd2; matrix_B[5] = 8'd7; matrix_B[6] = 8'd4; matrix_B[7] = 8'd6;
58     matrix_B[8] = 8'd9; matrix_B[9] = 8'd5; matrix_B[10] = 8'd2; matrix_B[11] = 8'd1;
59     matrix_B[12] = 8'd3; matrix_B[13] = 8'd6; matrix_B[14] = 8'd1; matrix_B[15] =
60         8'd7;
61
62     // --- Write matrix A ---
63     #2 wrt_en = 1;
64     for (i = 0; i < 16; i = i + 1) begin
65         din = matrix_A[i];
66         #2;
67     end
68     wrt_en = 0;
69
70     // --- Write matrix B ---
71     wrt_en = 1;
72     for (i = 0; i < 16; i = i + 1) begin
73         din = matrix_B[i];
74         #2;
75     end
76     wrt_en = 0;
77
78     // Wait for OUTPUT state
79     wait(state_monitor == 3);
80
81     // Print each output
```



```
79     for (i = 0; i < 16; i = i + 1) begin
80         @(posedge clk);
81         $display("Output_%0d:_%0d", i, dout);
82     end
83
84     #200 $finish;
85 end
86
87 endmodule
```

Listing 6: Testbench for Booth module

Code overview

This Verilog module performs matrix multiplication of two 4×4 matrices, A and B , using Booth multipliers. The result is a matrix $C = A \times B$, with each element $C[i][j]$ computed one at a time.

Inputs and Outputs

- **Inputs:**

- clk: Clock signal.
- rst: Reset signal.
- wrt_en: Write enable for loading data.
- din (signed 8-bit): Input data stream for matrices A and B .

- **Output:**

- dout (signed 17-bit): Resulting element $C[i][j]$.

Internal Registers and Signals

- **Matrices:**

- A[0:15], B[0:15]: Flattened row-major storage of 4×4 matrices.

- **FSM State Register:**

- state (3-bit): Tracks module states: IDLE, LOAD_A, LOAD_B, COMPUTE, OUTPUT.

- **Counters:**

- load_count (6-bit): Counts loaded elements.
- output_count (4-bit): Indexes output matrix elements.
- compute_count (2-bit): Handles 3-cycle multiplier latency.

- **Indexing:**

- `row = output_count[3:2], col = output_count[1:0]`

- **Multiplier Outputs:**

- `p0-p3` (signed 16-bit): Partial products.
 - `valid0-valid3`: Validity flags for multiplier outputs.

- **Sum:**

$$\text{sum} = p0 + p1 + p2 + p3 \quad (17\text{-bit})$$

Finite State Machine (FSM)

1. **IDLE**: Waits for `wrt_en`, transitions to `LOAD_A`.
2. **LOAD_A**: Loads 16 elements of *A*.
3. **LOAD_B**: Loads 16 elements of *B*.
4. **COMPUTE**: Pulses `ld`, waits 3 cycles.
5. **OUTPUT**: Assigns `sum` to `dout`, increments `output_count`.

Booth Multiplier Module: `Booth_Multiplier_4xA`

A sequential Booth multiplier module for two 8-bit signed inputs producing a 16-bit signed product over 3 cycles.

Parameters and I/O

- **Parameter:** $N = 8$

- **Inputs:**

- `Rst`, `Clk`, `Ld`
 - `M`, `R` (8-bit signed): Multiplicand and multiplier

- **Outputs:**

- `Valid`: High when product is ready.
 - `P`: 16-bit signed product.

Operation

- On Ld:
 - Cntr set to 3
 - A loaded with M
 - Prod loaded with R
 - Guard set to 0
- Each cycle:
 - Booth = { Prod[3:0], Guard }
 - Decoded to determine operation using shifted A
 - Operand B and C selected or disabled
- Adder Tree:
 - $T = H_i + B + Ci_B$
 - $S = T + C + Ci_C$
- Update:
 - Prod shifts right by 4, S inserted at top
 - Guard updated from Prod[3]
- Completion:
 - When Cntr == 1, P set to product, Valid is asserted

4.3 Comparison

Booth Multiplier

- Reference from Booth_Multiplier_4xA[1], a sequential multiplier that processes the multiplier in 4-bit chunks.
- For an 8-bit multiplier ($N = 8$), requires $\lceil (N + 1)/4 \rceil = 3$ cycles.
- FSM includes a COMPUTE state for 3-cycle latency, followed by OUTPUT to assign the sum.
- Four instances compute $A[i][k] \times B[k][j]$ for $k = 0$ to 3, and their outputs are summed as:

$$\text{sum} = p_0 + p_1 + p_2 + p_3$$

- No pipelining; each output element waits for previous multiplication to complete.

Baugh-Wooley Multiplier

- Uses the Baugh-Wooley technique for matrix multiplication.
- Product computed in one clock cycle after inputs stabilize.
- FSM lacks a **COMPUTE** state; the **OUTPUT** state directly assigns the result.
- Four multiplier instances compute partial products, with immediate summation.

Latency and Throughput

Booth Multiplier

- **Per Output Latency:** 4 clock cycles per element:
 - 1 cycle in **OUTPUT** to assign 1d.
 - 3 cycles in **COMPUTE** to generate valid output.
- **Total Time for 16 Outputs:**
 $2 (\text{reset}) + 16 (\text{LOAD_A}) + 16 (\text{LOAD_B}) + 3 (\text{initial COMPUTE}) + 16 \times 4 (\text{per output}) = 101 \text{ cycles}$
- **Total Time (2 ns clock):** $101 \times 2 \text{ ns} = 202 \text{ ns}$
- **Throughput:** 1 output every 4 cycles.

Baugh-Wooley Multiplier

- **Per Output Latency:** 1 clock cycle per element (fully combinational multiplication).
- **Total Time for 16 Outputs:**
 $2 (\text{reset}) + 16 (\text{LOAD_A}) + 16 (\text{LOAD_B}) + 16 (\text{OUTPUT}) = 50 \text{ cycles}$
- **Total Time (2 ns clock):** $50 \times 2 \text{ ns} = 100 \text{ ns}$
- **Throughput:** 1 output every cycle.

Resource Usage

Booth Multiplier

- Each `Booth_Multiplier_4xA` is sequential and uses registers and a pipelined adder tree.
- Estimated LUTs per 8-bit multiplier: $\sim 50\text{--}100$
- Total for 4 multipliers: $\sim 200\text{--}400$ LUTs

- Flip-flops per multiplier: ~ 20 – 50 , total ~ 80 – 200
- Lower combinational delay may support faster clock (e.g., < 2 ns), but limited by sequential latency.

Baugh-Wooley Multiplier

- Each `baugh_wooley_multiplier` is combinational, with full partial product arrays and correction terms.
- Estimated LUTs per 8-bit multiplier: ~ 200 – 300
- Total for 4 multipliers: ~ 800 – 1200 LUTs
- Minimal flip-flops used (~ 0 – 20 total)
- Higher combinational delay may necessitate slower clock (e.g., 5 – 10 ns depending on synthesis).

5 Chapter 5: Results

5.1 Experiment setups

The matrix multiplier design was verified using a custom testbench written in Verilog and simulated using Vivado Simulator. The experimental environment and conditions are summarized as follows:

- **Simulation Tool:** Vivado Simulator
- **Testbench:** `tb_Matrix_multi`
- **Clock Period:** 2 ns (500 MHz)
- **Reset Duration:** 0–4 ns (2 clock cycles)
- **Input Signal:** `din` (8-bit signed), loaded sequentially
- **Write Control:** `wrt_en` signal enabled during `LOAD_A` and `LOAD_B`
- **Data Count:** 16 values for Matrix A and 16 values for Matrix B
- **Output Signal:** `dout` (17-bit signed), 1 element per cycle during `COMPUTE`
- **FSM States:** `IDLE` → `LOAD_A` → `LOAD_B` → `COMPUTE` → `IDLE`

The goal of the simulation is to verify the functional correctness of data loading, multiplication using four parallel cores, and the sequential output of 16 computed values. Waveform snapshots were captured for each FSM stage and are presented in the following section.

5.2 Waveform

5.2.1 Waveform for Baugh-Wooley method

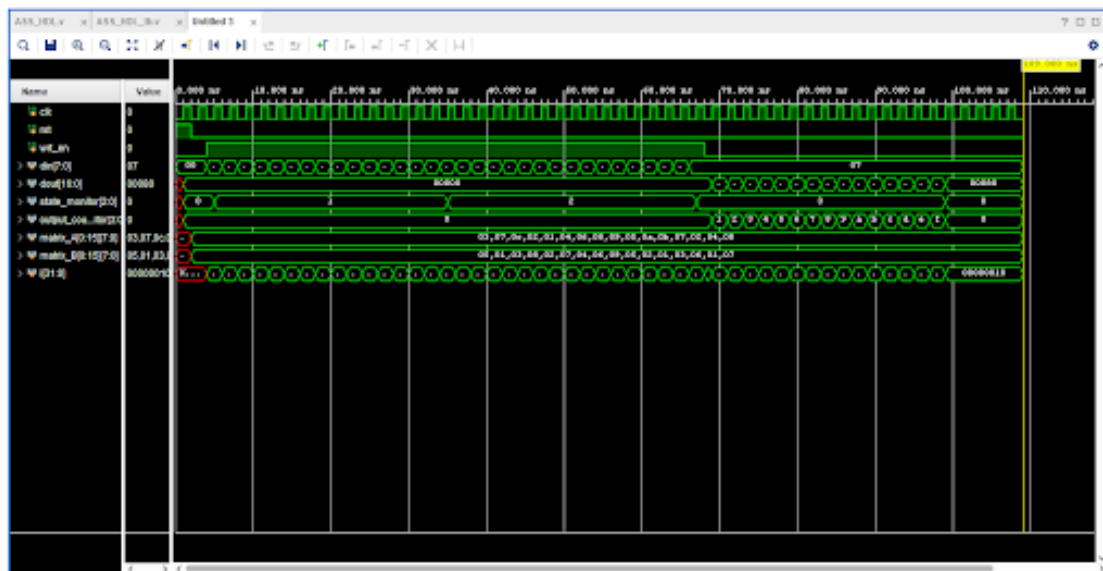


Figure 4: Waveform

5.2.2 Explanation:

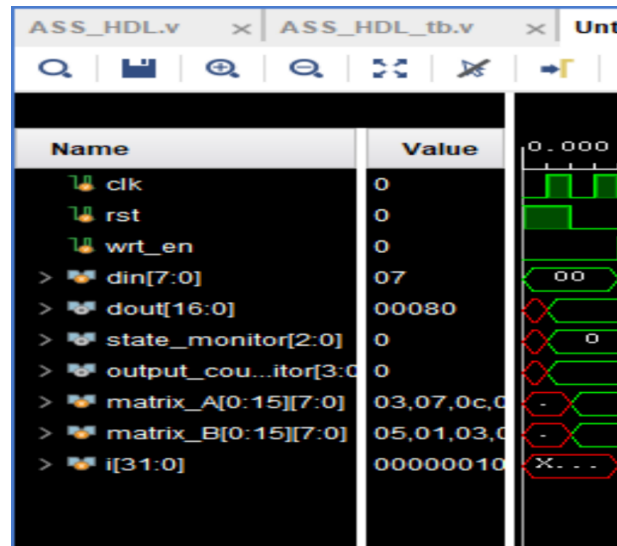


Figure 5: STATE IDLE

- The system remains in the IDLE state and waits for the `wrt_en` signal to be asserted.
- Once `wrt_en` is high, the first value from `din` is stored in `A[0]`.
- The internal counter `load_count` is initialized to 1 to track the next input index.
- The `done` signal is also reset to 0 to indicate that the system is ready for a new computation cycle.
- After this initialization step, which takes exactly one clock cycle (2 ns), the FSM transitions to the `LOAD_A` state.

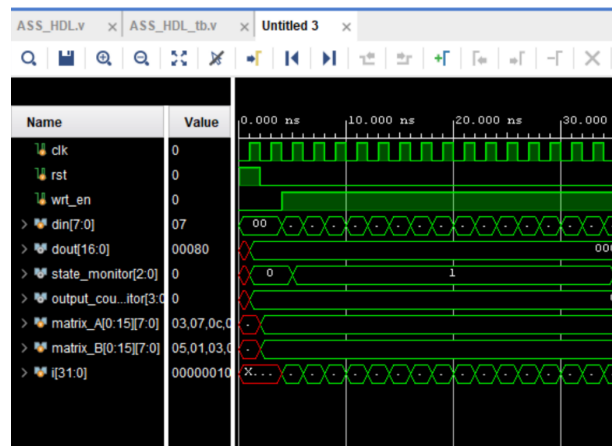


Figure 6: STATE LOAD_A

- The system waits for the `wrt_en` signal to be asserted.
- Once active, input values from the `din` port are transferred sequentially into matrix A, starting from A[1] and continuing to A[15].
- (Note: A[0] was already loaded during the IDLE state.)
- The entire LOAD_A stage spans from **5 ns to 35 ns**, which corresponds to 15 clock cycles at 2 ns per cycle.
- After successfully loading all 16 elements into matrix A, the FSM transitions to the LOAD_B state.

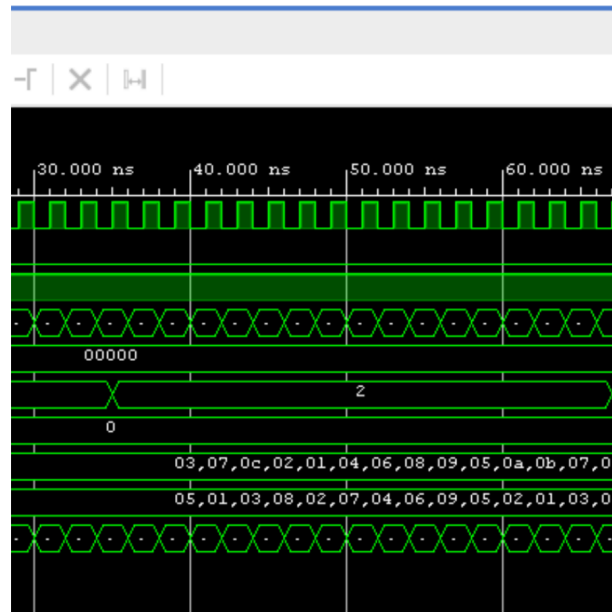


Figure 7: STATE LOAD_B

- The system waits for the `wrt_en` signal to be enabled.
- Once active, the values from the input port `din` are sequentially transferred into matrix B, starting from B[0] up to B[15].
- The process is controlled by an internal counter that ensures exactly 16 values are stored.
- The entire LOAD_B stage occurs from **35 ns to 67 ns**, which corresponds to 16 clock cycles with a 2 ns period per cycle.
- After completing the loading of matrix B, the FSM automatically transitions to the COMPUTE state.

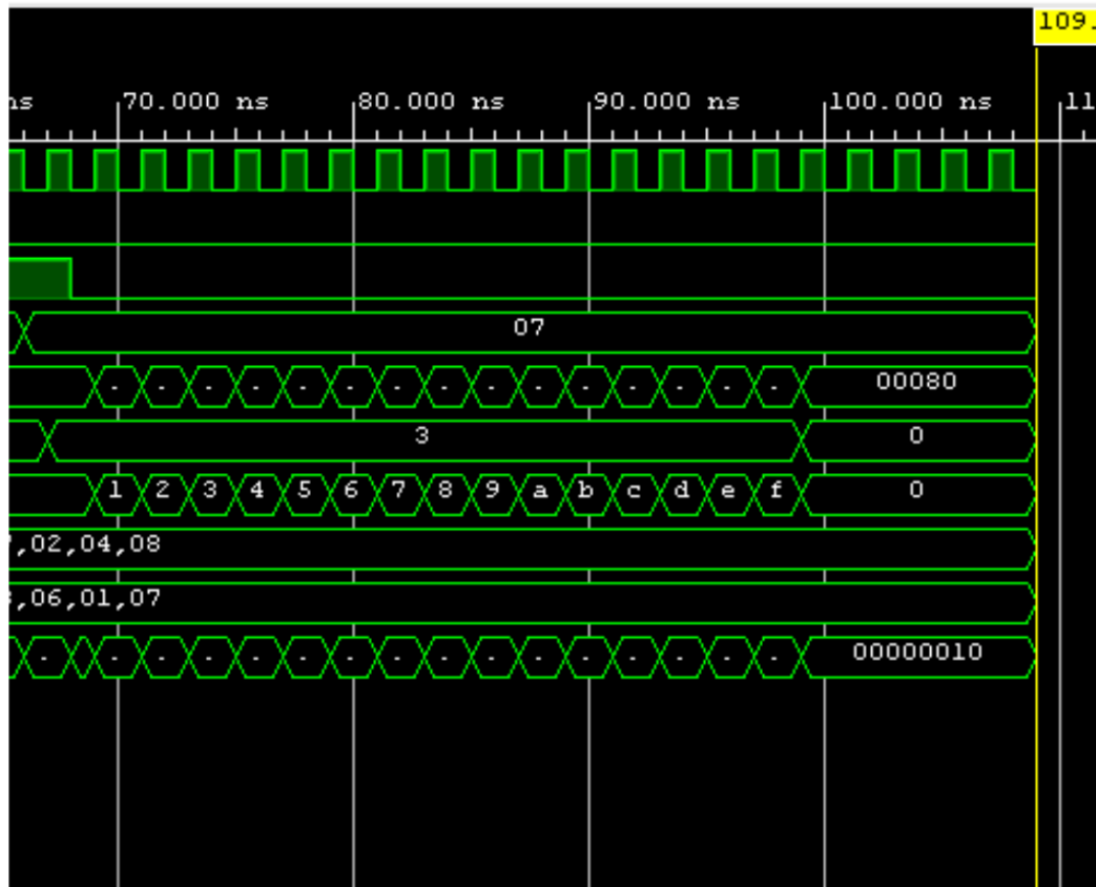


Figure 8: STATE COMPUTE & OUTPUT

- During the COMPUTE state, the system calculates one result element at position `C[i][j]` in the output matrix.
- Each element is computed using four partial products, obtained from the `baugh_wooley_multiplier` module instances.
- These four multipliers operate in parallel to compute:
$$\text{sum} = p_0 + p_1 + p_2 + p_3$$
 where each p_k is a signed 16-bit product.
- The result of this computation is assigned to the output signal `dout`.
- The full COMPUTE phase occurs from **67 ns to 99 ns**, producing 16 output values sequentially, one per clock cycle.
- After completing all outputs, the FSM transitions back to the IDLE state and remains there for 10 ns before simulation ends.

5.2.3 Waveform for Booth multiplier method

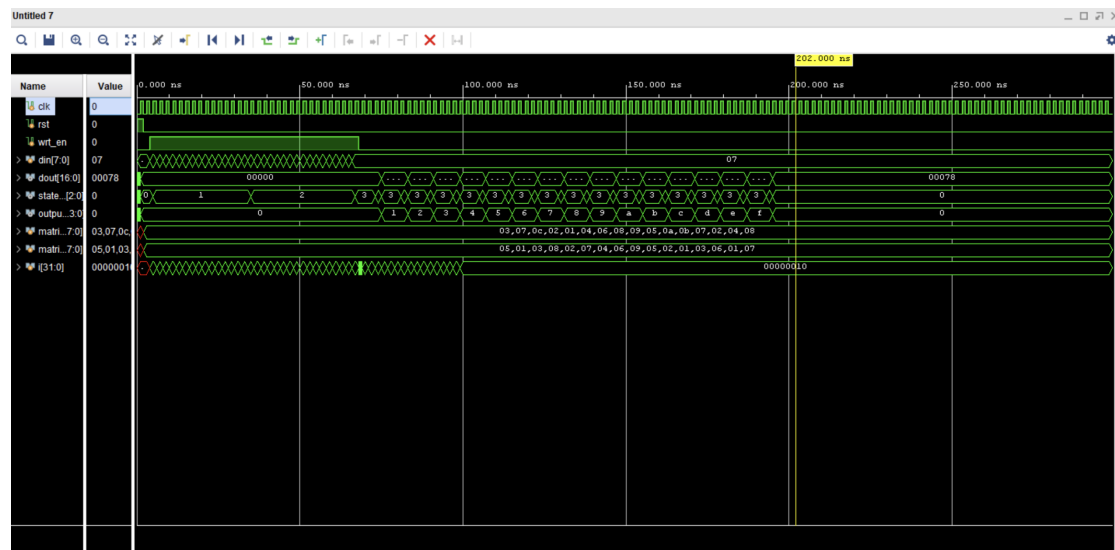


Figure 9: Waveform

5.2.4 Explanation:

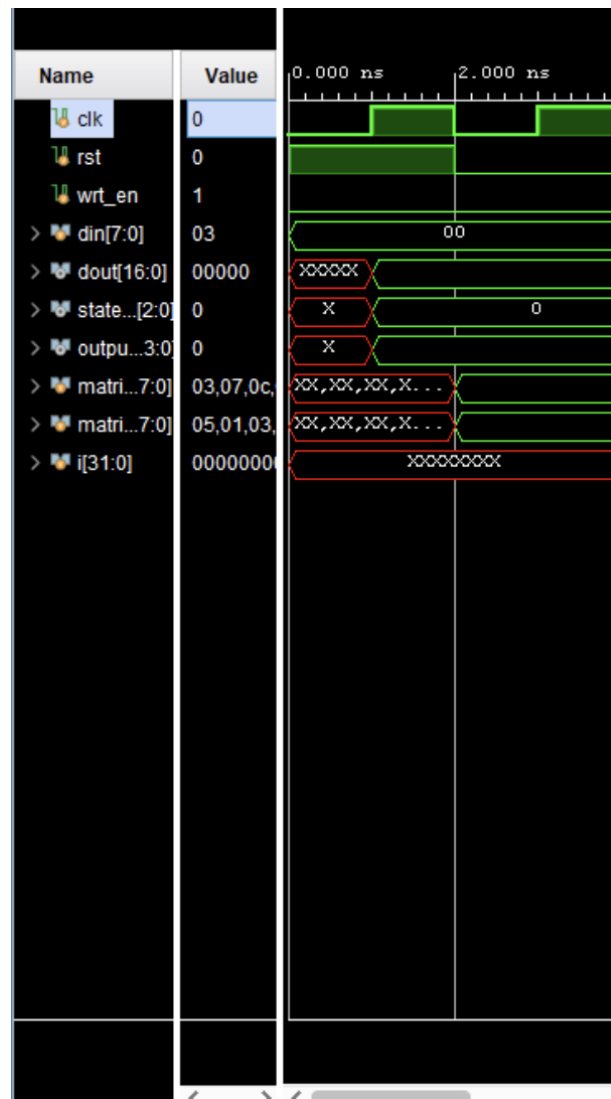


Figure 10: STATE IDLE

- 2 cycles of Reset (reset and wait for wrt_en) from 0-4 ns.

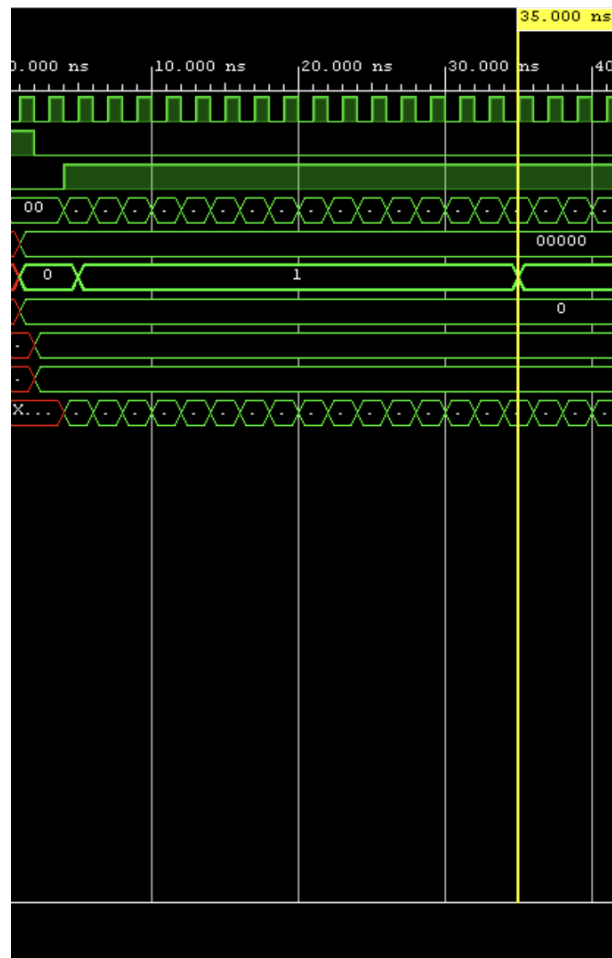


Figure 11: LOAD A

- wrt_en turns on, wait for state LOAD A to 5ns.
- LOAD A for 16 clock cycles = 32ns
- All processes are from 4ns -> 35ns.

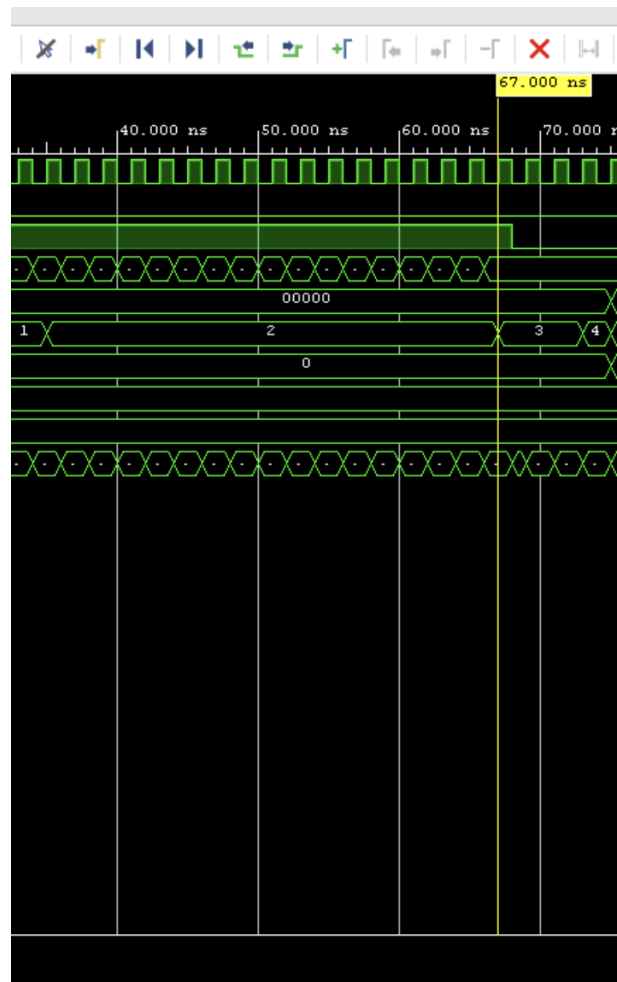


Figure 12: LOAD B

- LOAD B need 16 clock cycles to take in 16 elements from the matrix. needs 32 ns. The process takes place from 35 ns -> 67 ns.

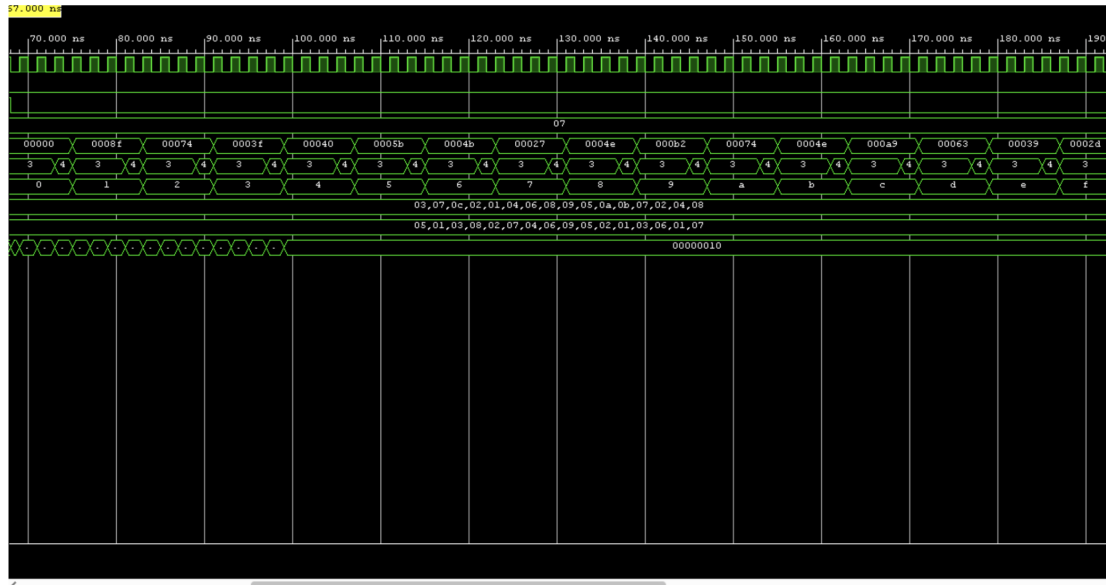


Figure 13: COMPUTE & OUTPUT

- COMPUTE state is to handle the 3-cycle latency of the Booth multipliers.
- The ld signal is pulsed for one cycle to start the four multiplications, and a counter (compute_count) waits for 3 cycles before transitioning to OUTPUT.
- The sum of the four partial products is assigned to dout, and ld is pulsed again for the next element's multiplications, maintaining one output per cycle after the initial latency.
- All processes need 3 (initial COMPUTE)+164 (OUTPUT + COMPUTE per element) = 67clock cycles \Rightarrow 134ns.
- The process takes place from 67 ns -> 201 ns with the IDLE state at the end of the process.

6 Conclusion

- The Booth design uses fewer LUTs and adds flip-flops, making it more resource-efficient for FPGAs with limited combinational resources but requiring more sequential elements. Moreover, the Booth design may consume less power in dynamic scenarios, while the Baugh-Wooley design might be more power-efficient in static conditions.
- In terms of timing constraints, Booth design can operate at a higher clock frequency, but the overall execution time, as observed in the waveform, is longer due to latency. The Baugh-Wooley design is better for applications where latency is critical and clock frequency can be adjusted.
- In Booth Multiplier, the sequential nature makes it scalable to larger bit widths with manageable resource growth, as the number of cycles increases linearly with $N/4$, while the Baugh-Wooley technique, with its combinational nature, leads to an exponential growth with bit width. This makes the Baugh-Wooley technique less suitable for large matrices due to the delay and area constraints.



7 Chapter 7: References

References

- [1] *Book.*
https://github.com/MorrisMA/Booth_Multipliers/blob/master/Src/Booth_Multiplier_4xA.ucf
- [2] *Book: VLSI Design 2011 by Gerald E. Sobelman.*
<http://dce.hust.edu.vn/wpcontent/uploads/2017/05/MultiplierBaughWooley.pdf>