

Final Project Report – ECE 564

Name: Rushiraj Chaitanyakumar Sheth
Unityid: rsheth@ncsu.edu
StudentID: 200596988

Delay (ns to run provided provided example).
Clock period: **5.80 ns**
cycles": **3354**
Delay (clock period * #cycles):
19453.2 ns

Logic Area:
(um²)
12866.6859
Memory: N/A

1/(delay*area) (ns⁻¹.um⁻²)
4.00 x 10⁻⁹ (ns⁻¹. um⁻¹)

Delay (TA provided example. TA to complete)

1/(delay.area) (TA)

RTL Design: Self-attention matrix – LLM Transformers

Rushiraj Chaitanya Kumar Sheth

Abstract

This project focuses on the hardware implementation of the scaled dot-product attention mechanism $(\frac{QK^T}{\sqrt{d_k}})V$.

The objective is to design and verify Register Transfer Level (RTL) architecture that efficiently computes the attention scores, leveraging Verilog for its implementation. The project involves three main computational stages:

- forming the query, key, and value matrices by multiplying input embeddings with weight matrices,
- computing the score matrix through the multiplication of the query and the transposed key matrices,
- obtaining the scaled dot-product attention output by multiplying the score matrix with the value matrix.

The design is verified using a testbench that handles the communication and synchronization of data via SRAM interfaces, considering timing constraints and ensuring accuracy. The synthesized design is optimized for performance and verified for correctness through Synopsys and ModelSim tools.

Report Structure: The report proceeds as under:

Section 1: Introduction

Section 2: Micro-architecture

- 2.1: Hardware algorithmic approach implemented
 - 2.1.1: Matrix multiply
 - 2.1.2: Matrix transpose
 - 2.1.3: Datapath design

Section 3: Interface specification

- 3.1: List of signals, their width and functions
- 3.2: Timing diagram

Section 4: Finite state machine design

Section 5: Results achieved

Section 6: Conclusion

1. Introduction

The module performs three main operations:

- **Matrix Multiplications:** Compute query (Q), key (K), and value (V) matrices by multiplying input embeddings with corresponding weight matrices.
- **Score Calculation:** Multiply the query matrix (Q) with the transposed key matrix (K^T) to produce the attention scores.
- **Scaled Dot-Product Attention:** Compute the final output by multiplying the score matrix with the value matrix (V).

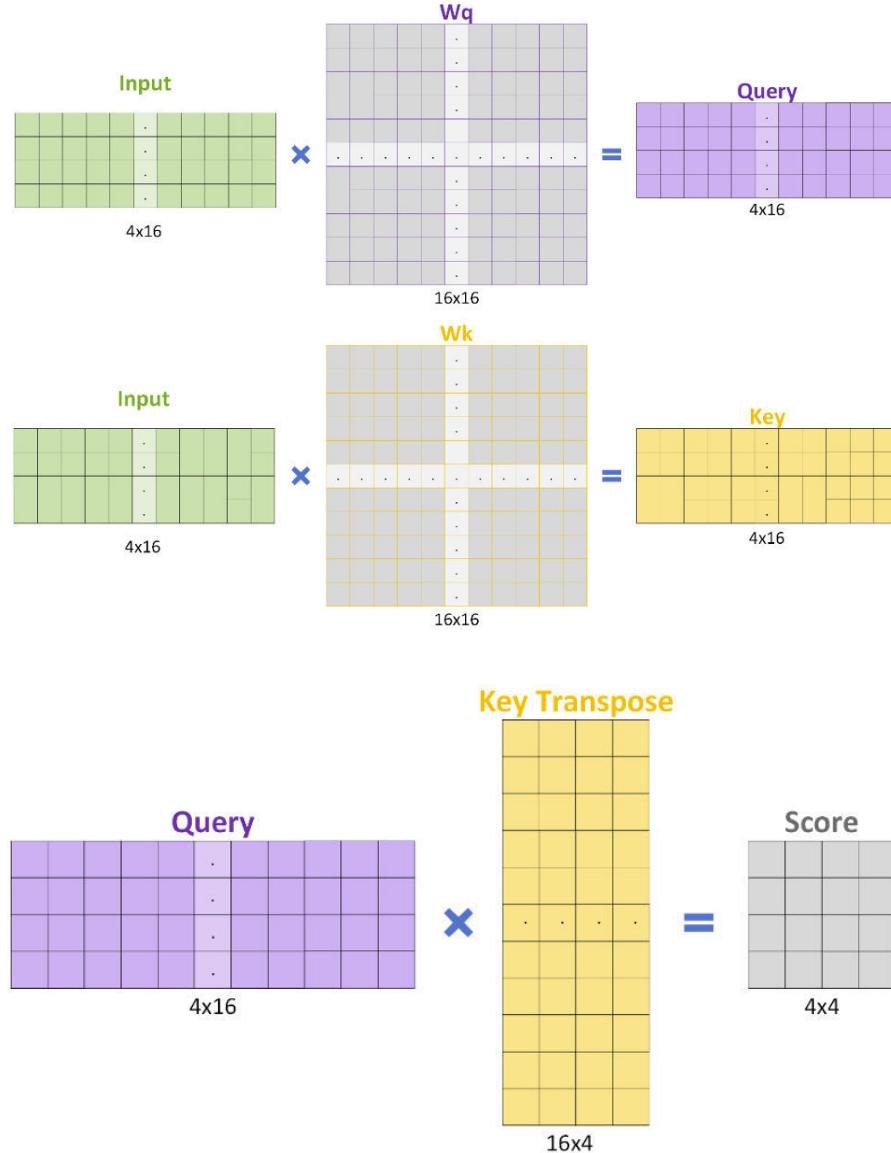


Figure 1: Calculation of Query, Key and Score matrices

To optimize performance, the design uses pipelining, allowing the overlapping of multiple computation stages to minimize latency. Additionally, parallel processing techniques accelerate matrix operations. Data is managed through an SRAM interface, ensuring efficient memory access and synchronization using control signals.

Key Features:

- **Pipelining:** Implements a multi-stage pipeline to enhance throughput and reduce overall computation time.
- **Parallelism:** Leverages parallel data paths to perform matrix operations efficiently.

2. Micro-Architecture:

2.1 Hardware “algorithmic” approach implemented:

2.1.1 Matrix Multiply:

- In order to perform matrix multiplication of two matrices, three counters are used, namely counter “i”, “j” and “k”.
- “Counter i” keeps track of no. of rows of input matrix. “Counter j” keeps track of no. of columns of weight matrix and “counter k” keeps track of no. of rows of weight matrix.
- All the three counter are down-counters. And they are initialized as: counter i = rows_of_input_matrix; counter j = col_of_weight_matrix; counter k = rows_of_weight_matrix.
- When the SRAM produces required matrix elements, the counters start decrementing.
- This approach is similar to matrix multiplication algorithm in high level languages like C++ i.e.

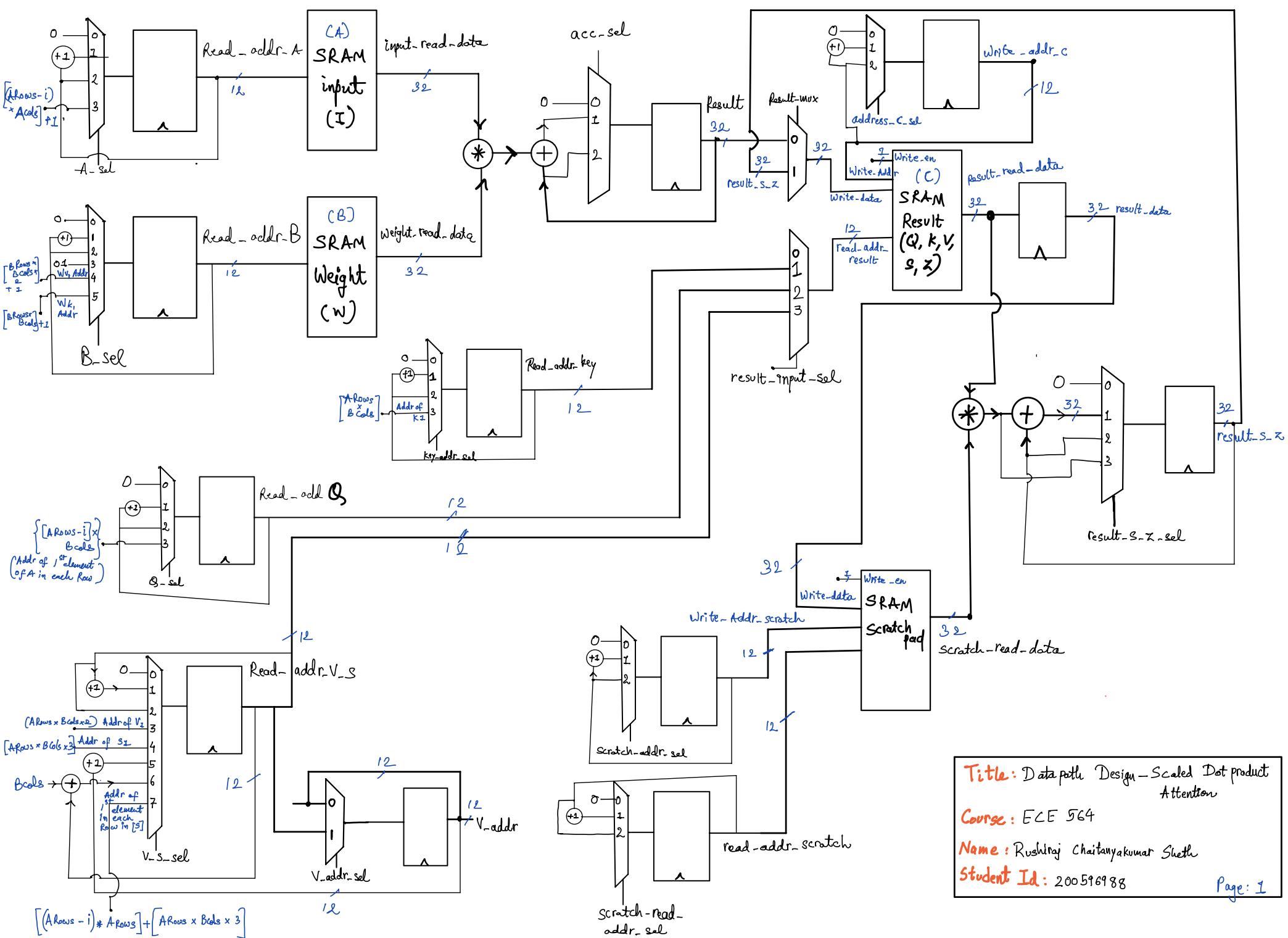
```
for(count i :1 to A_rows ){  
    for(count j:1 to B_cols ){  
        for(count k:1 to B_rows )  
            {      result[i][j] += (input[i][k]*weight[k][j]);  
    }}}
```

- Finally, since there are three weight matrices, the above multiplication is repeated three times i.e. computing Query(Q), Key(K) and value(V) matrix.
- “Counter c” is used to keep track of this.
- The computed elements of matrix Query, Key and Value are stored in result SRAM respectively.
- Further, to compute the score matrix, the elements of Key matrix are read from Result SRAM (in transpose order) and stored in scratchpad SRAM.
- Then, further the matrix multiplication is performed between elements of Query matrix (from result SRAM) and that of key matrix (read from scratchpad SRAM) and the convolution result as **Score matrix**, is stored in Result SRAM at address next to last stored value in it i.e. after elements of value matrix.
- Finally, the matrix V is stored in scratchpad SRAM in column major order and then convolution is performed between matrix S (read from result SRAM) and matrix V (read from Scratchpad SRAM in col. major order) and the result as **Attention matrix** is stored in result SRAM.

2.1.2 *Transpose of matrix:*

- In order to compute the transpose of matrix, the elements are stored in scratchpad SRAM in column major order and then referenced further.

2.1.3 *Datapath design:*



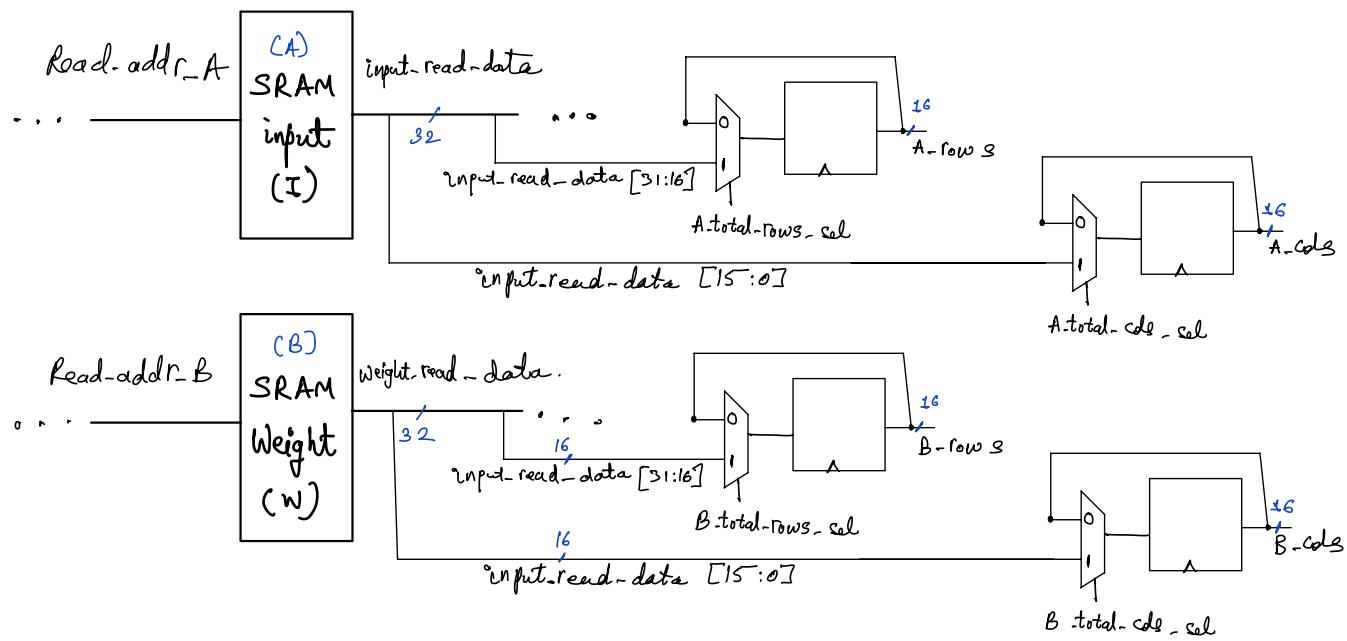
Title: Data path Design - Scaled Dot product Attention

Course: ECE 564

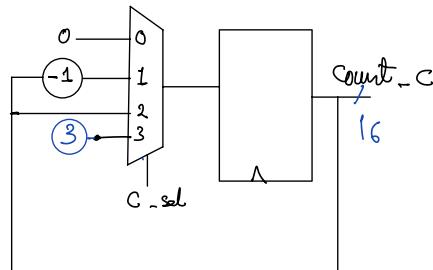
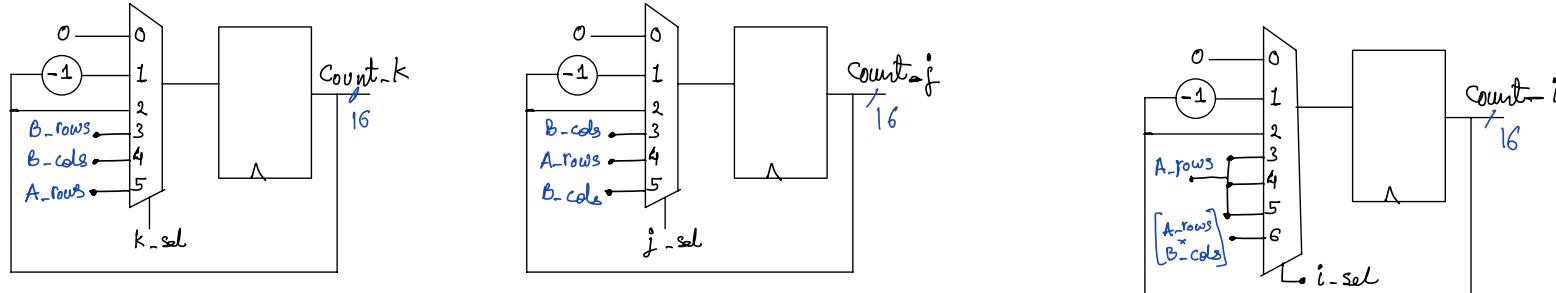
Name: Rushraj Chaitanya Kumar Sheth

Student Id: 200596988

Page: 1



* Counters:



Title: Data path Design - Scaled Dot product Attention
Course: ECE 564
Name: Rushikraj Chaitanya Kumar Shekhar
Student Id: 200596988
Page: 2

3. Interface Specification

3.1 List of signals, their widths and functions:

Sr. no.	Signal Name	Width(bits)	Function
1	A_sel	2	generates read address for input SRAM
2	B_sel	3	generates read address for weight SRAM
3	acc_sel	2	used for accumulating and storing the convolution result of input and weight matrices.
4	address_c_sel	2	generates write address for result SRAM
5	result_mux	1	0: store convolution data of input and weight matrices i.e. Q, K, V to result SRAM. 1: store convolution data of ($Q * K^T = S$) score and ($S * V = Z$) attention matrices to result SRAM.
6	key_addr_sel	2	generates read address for key matrix stored in Result SRAM.
7	Q_sel	2	Generates read address of Q matrix stored in result SRAM.
8	v_s_sel	3	Generates read address of matrix S and matrix V, both of which are stored in Result SRAM.
9	result_input_sel	2	Value 1: connects read address of key matrix to read_addr of Result SRAM. Value 2: connects read address of Query/Q matrix to read_addr of Result SRAM. Value 3: connects read address of value(V) and score(S) matrix to read_addr of Result SRAM.
10	scratch_addr_sel	2	Generates write address for scratchpad SRAM.
11	scratch_read_addr_sel	2	Generates read address for scratchpad SRAM.
12	result_s_z_sel	2	used for accumulating and storing the convolution result of score (s) and attention (z) matrices.

13	k_sel	3	Counter k as count down counter (usage described in 2.3)
14	j_sel	3	Counter j as count down counter (usage described in 2.3)
15	i_sel	3	Counter i as count down counter (usage described in 2.3)
16	c_sel	3	Counter c as count down counter (usage described in 2.3)
17	A_total_rows_sel	1	Used to store total rows of input matrix.
18	A_total_cols_sel	1	Used to store total columns of input matrix.
19	B_total_rows_sel	1	Used to store total rows of weight matrix.
20	B_total_cols_sel	1	Used to store total columns of weight matrix.

3.2 Timing Diagram:

- The timing diagram obtained after successful completion of all test cases is as below.
 - Below diagrams are related to testcase no. 1 only.

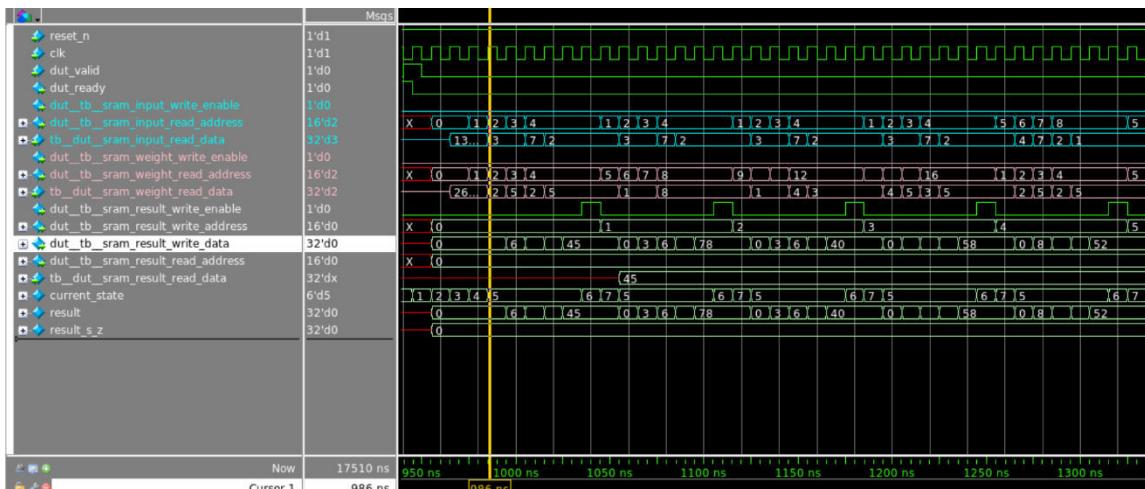


Figure 2: Test case 1 timing diagram

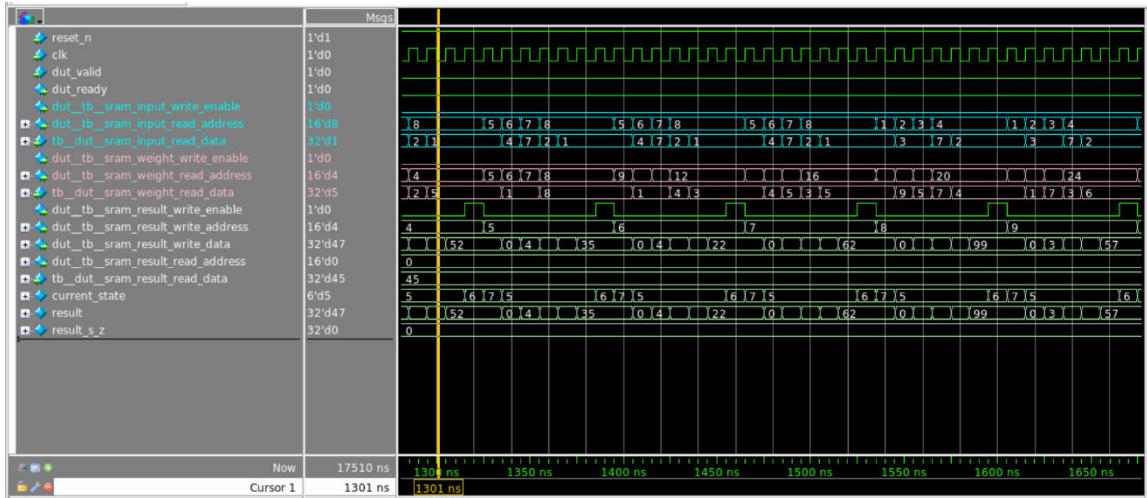


Figure 2: Test case 1 timing diagram

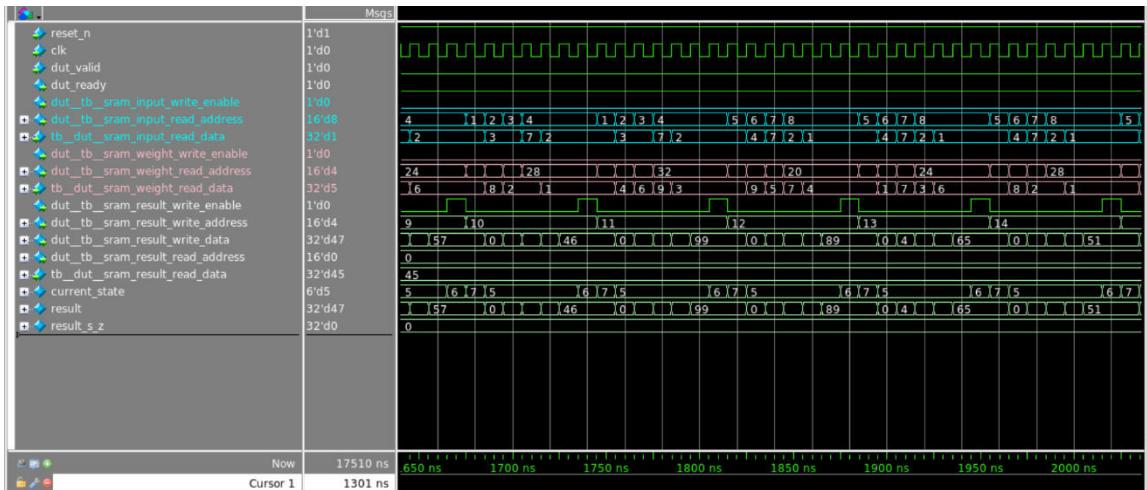


Figure 3: Test case 1 timing diagram

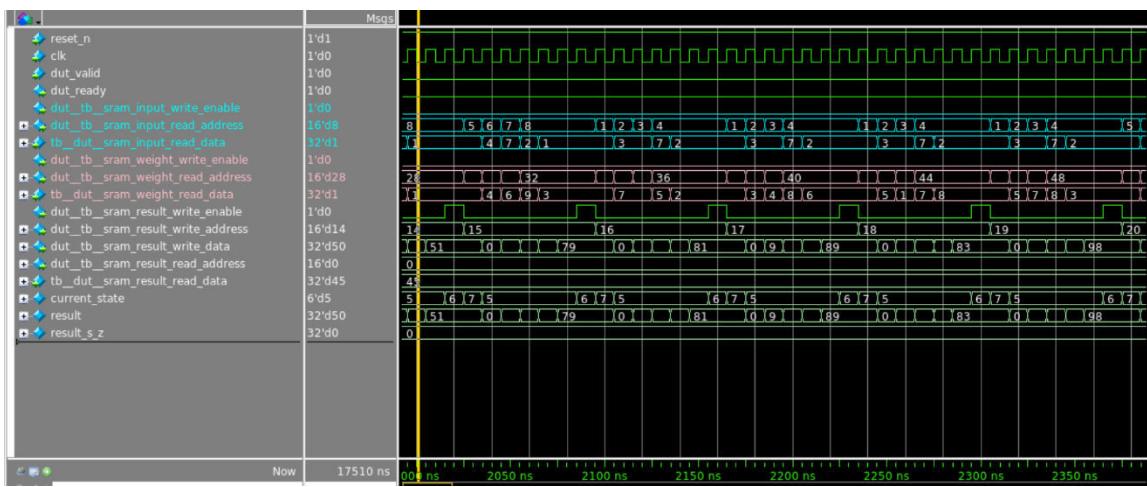


Figure 4: Test case 1 timing diagram

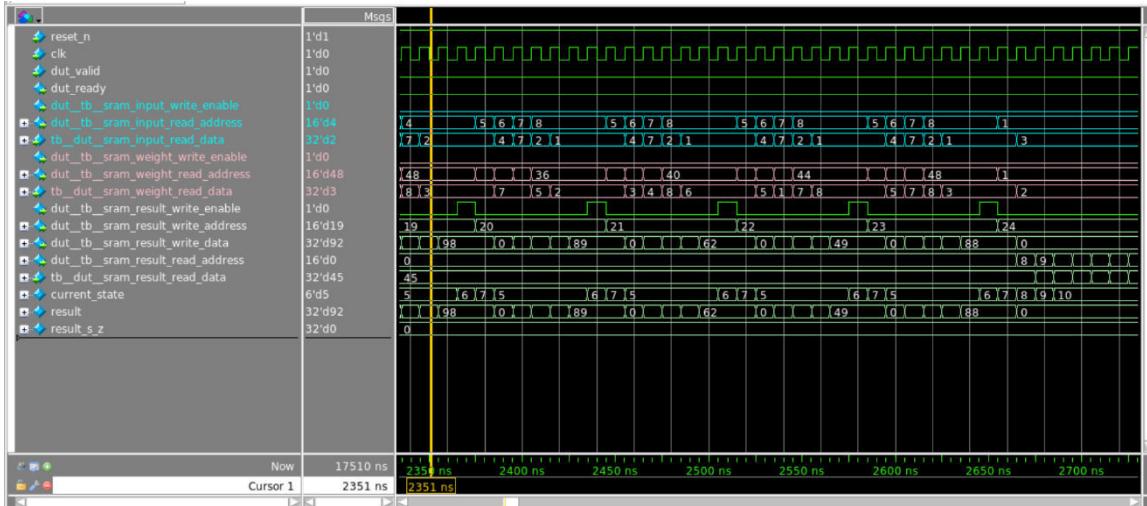


Figure 5: Test case 1 timing diagram

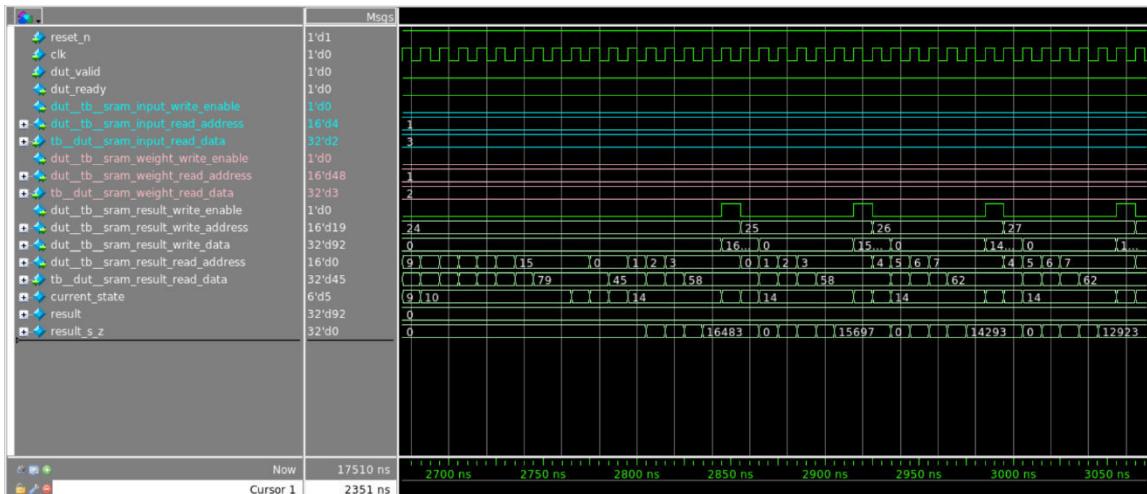


Figure 6: Test case 1 timing diagram

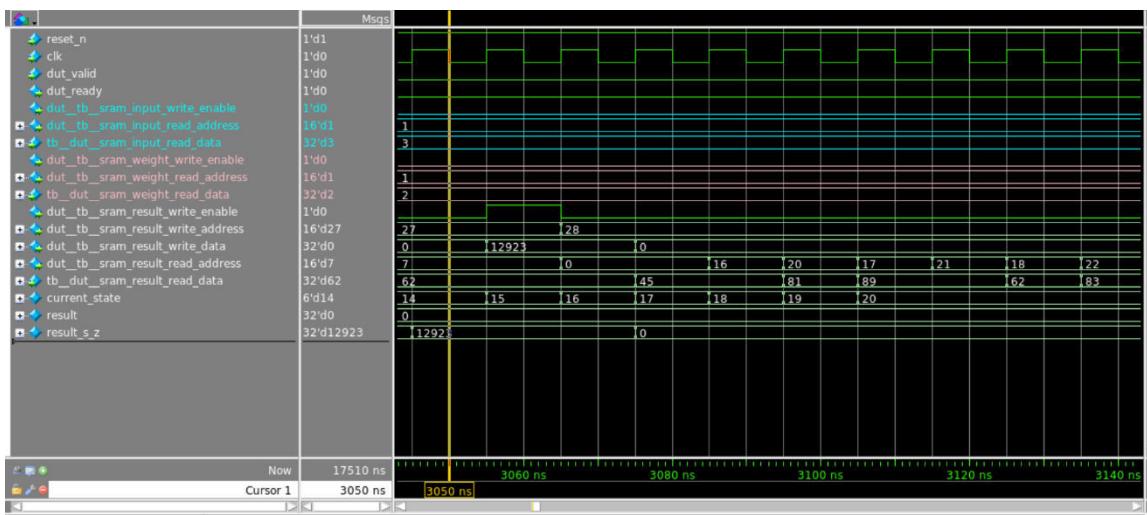


Figure 7: Test case 1 timing diagram

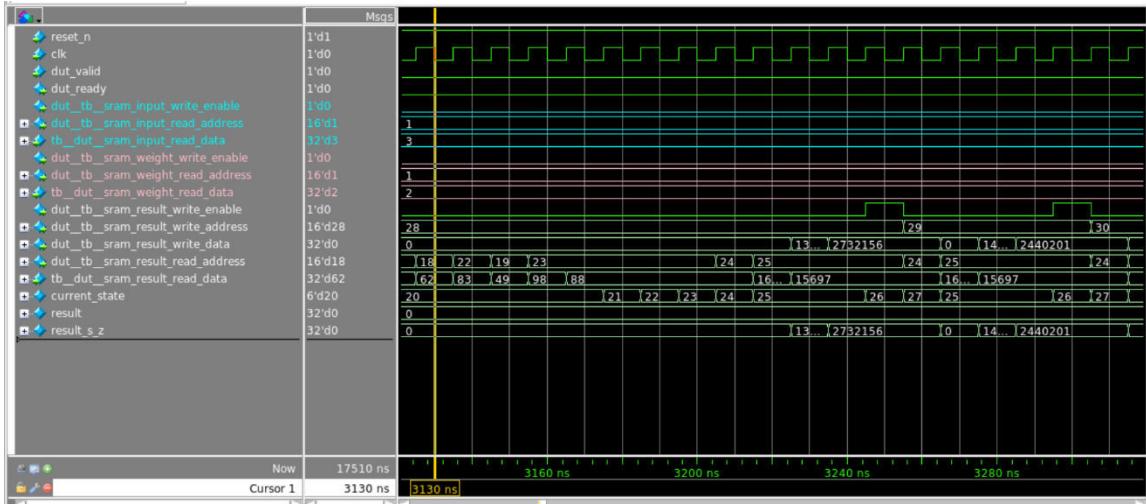


Figure 8: Test case 1 timing diagram

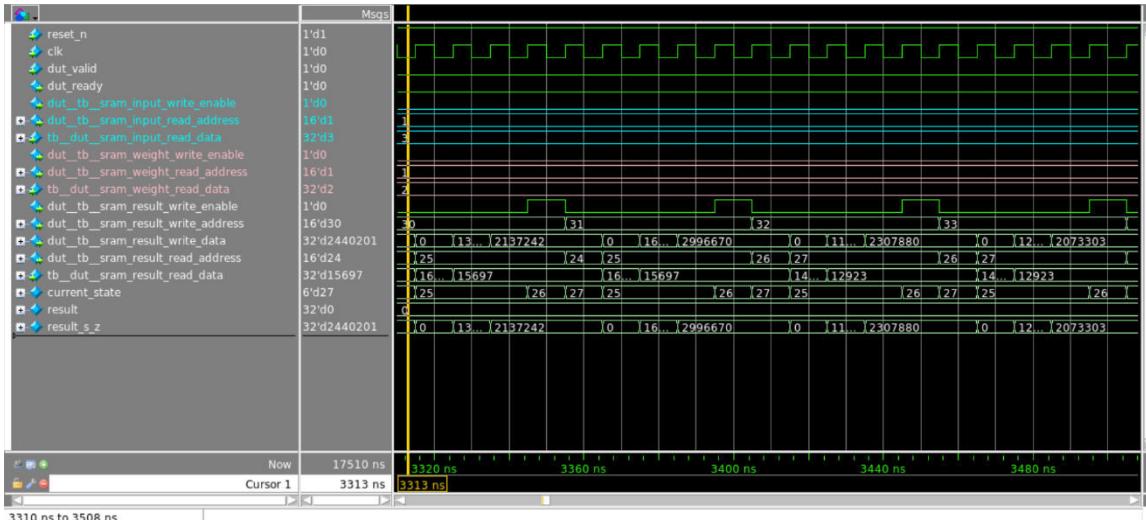


Figure 9: Test case 1 timing diagram

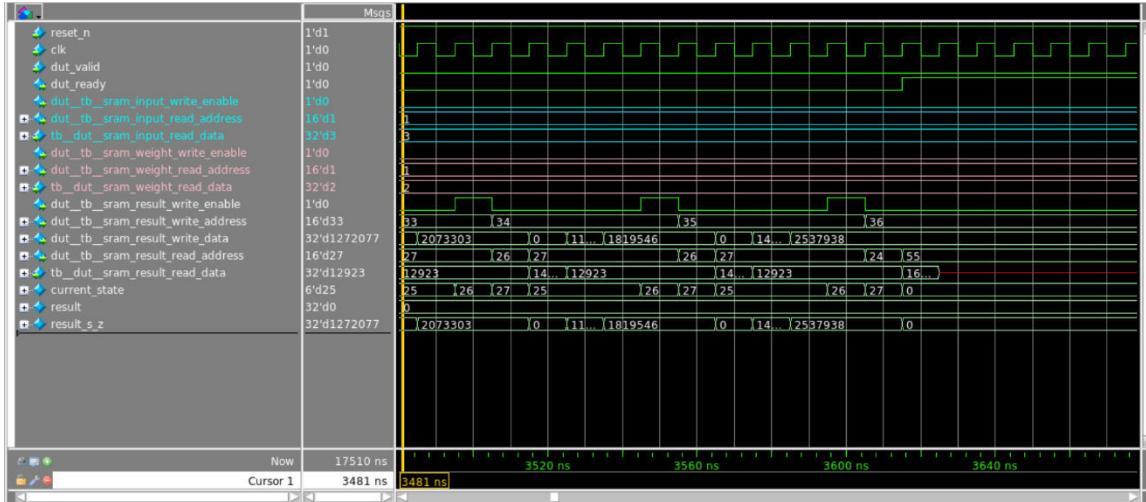
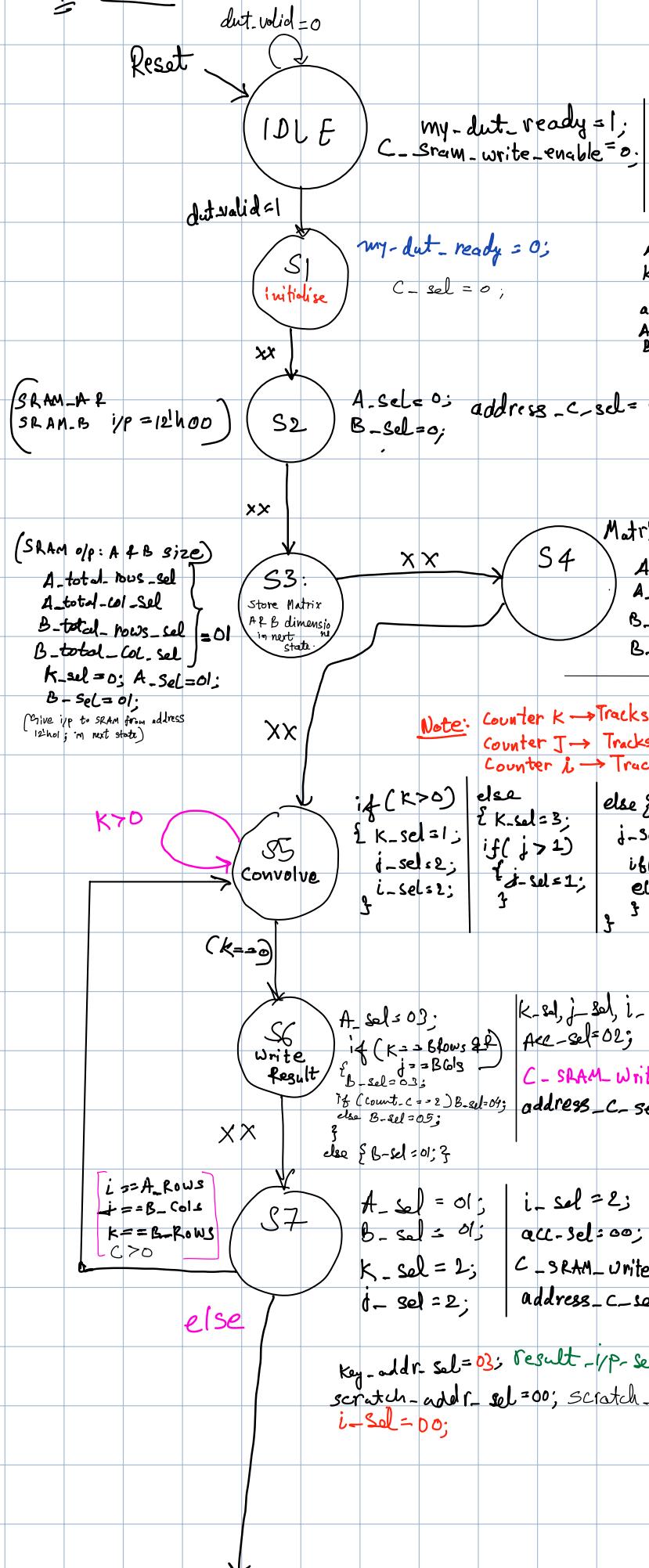


Figure 10: Test case 1 timing diagram

4. Finite state machine design:

- The designed FSM based on Moore machine is as mentioned below.
- It consists of 32 states.
- For optimizing area consumption, binary encoding is used to describe each state in Verilog.

FSM:



Name: Rushiraj Sheth

Student ID: 200596988

$$\left. \begin{array}{l} A_sel, B_sel, acc_sel, A_total_col_sel, \\ k_sel, j_sel, i_sel, A_total_col_sel \\ address_c_sel, \\ A_total_rows_sel, \\ B_total_rows_sel, \end{array} \right\} = 2^16xx;$$

$$\left. \begin{array}{l} A_sel, B_sel, acc_sel, A_total_col_sel, \\ k_sel, j_sel, i_sel, A_total_col_sel \\ address_c_sel, \\ A_total_rows_sel, \\ B_total_rows_sel, \end{array} \right\} = 2^{16}00;$$

Matrix dimension must be stored from now onwards, till complete.

Initialize & start counters from next clk cycle.

A-total-col-sel = 00, K-sel = 3
A-total-row-sel = 00, j-sel = 3
B-total-col-sel = 00, i-sel = 3
B-total-row-sel = 00, C-sel = 03;

Note: Counter K → Tracks B-Rows from #B Rows to 0. Counter C → Down Count from 3 to 1. (Since there are 3 weight matrices).

if (K > 0) { if (j > 1) { j-sel = 3, i-sel = 1; } else { acc-sel = 02; } } else { if (K == 0) { if (j == 1) { if (i == 1) { i-sel = 3, B-sel = 01; } else { i-sel = 2, B-sel = 02; } } else { i-sel = 01, B-sel = 01; } } }

K-sel, j-sel, i-sel = 02; C-sel = 02;
acc-sel = 02; result-MUX = 0;

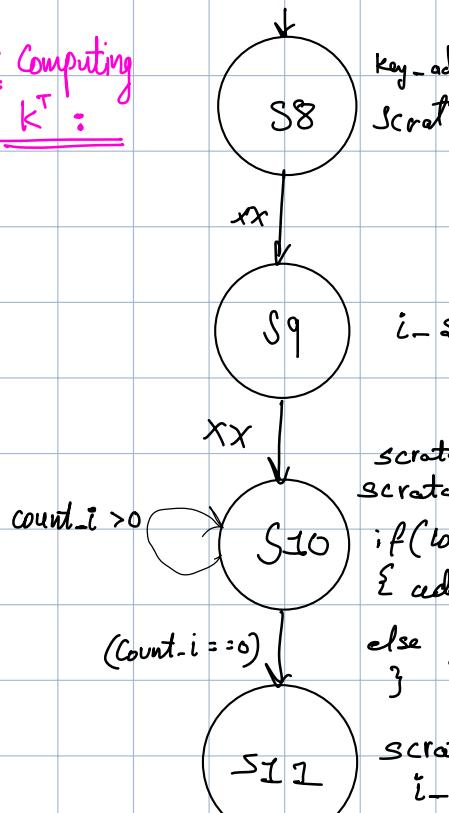
C-SRAM-write-en = 1;

address-c-sel = 1; (change address for next val. to be stored)

i-sel = 2; acc-sel = 00; C-SRAM-write-en = 0; address-c-sel = 2; key-addr-sel = 03; result-i/p-sel = 01; scratch-addr-sel = 00; scratch-write-en = 0; i-sel = 00;

if (i == A-rows & j == B-cols & C == 0) { A-sel = 02, B-sel = 02, C-sel = 02; acc-sel = 00; result-MUX = 0; C-SRAM-write-en = 0; address-c-sel = 1; } else { next-state = S5; }

Computing K^T :



`key_addr_sel = 01; i_SEL = 08;` // counter i = total elements of k
`scratch_addr_SEL = 02; scratch_Write_en = 0;` (Allows x reads)

count_i > 0

S_{10}

scratch_addr_sel = 01;
scratch_write_en = 1;
: 0C1 + i > ?)

if (count - i == 0)
 {
 i - sel = 06 ;
 scratch_addr . sel = 02 ;
 }
 else {
 i - sel = 2 ;
 }
 }

(Count - i == 0)

else
} {addr-keysel = 02;

 Transpose is Done here

10

Computing

Score Matrix:

$$S = Q \cdot K^T$$

S12

result - input - Sel = 01;

$$\text{Q-sel} = 0;$$

scratches - read addr - sel so; result - S - Z - Sel = 0;

scratch - read - addr r sel = 0 ;

K71

// Q_i & k_i are generated.
// counters are initialized.
// starting count down.

```

if( $k > 2$ )
{
     $g\_sel = 0;$ 
    scratch.read_addr.sel = 0;
}
else {
     $g\_sel = 2;$ 
    scratch.read_addr.sel = 2;
}

```

1

$\{ \begin{array}{l} f(k > 0) \\ K_{-Sel} = 1; \\ j-Sel = 2; \\ i-Sel = 2; \\ \text{else } \{ k_{-sel} = 4; \end{array} \}$

```

if( k > 0 )
{
    result_S-2 sel = 1;
}
else { result_S-2 sel = 2; }

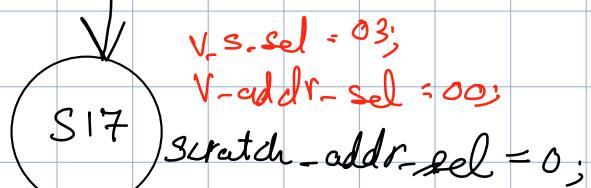
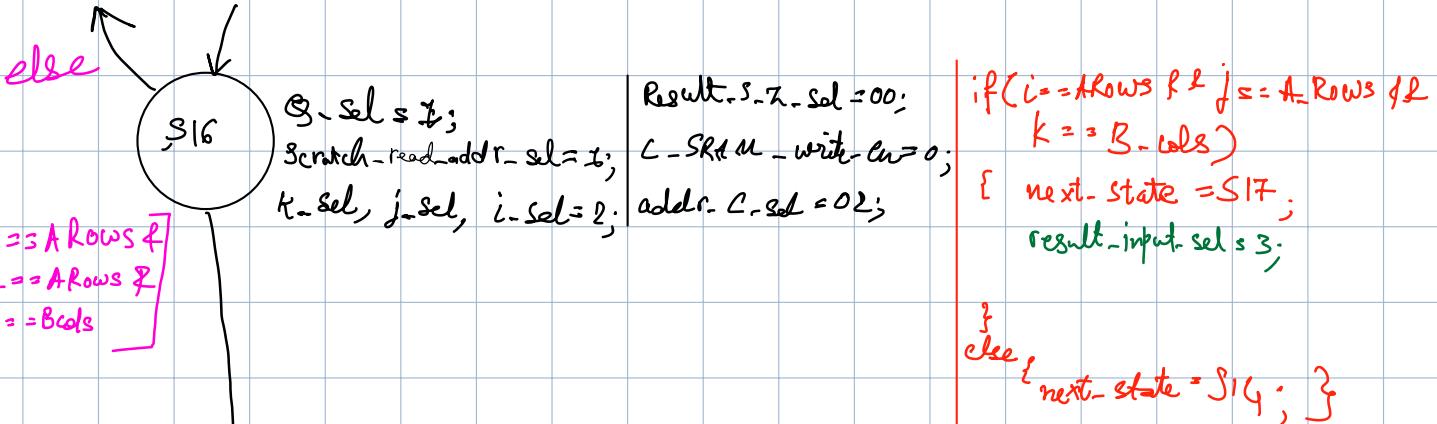
```

1

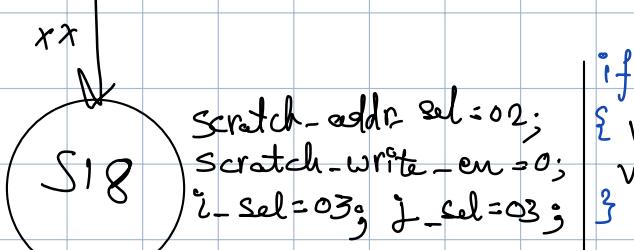
$\text{g_Sel} = 3;$
 $\text{if } C_L = -\beta \text{ cold ff}$

$\{ \text{scratch_read_addr_sel} = 0 \}$
 for $\{ \text{scratch_read_addr_sel} = 0 \}$
 $\}$

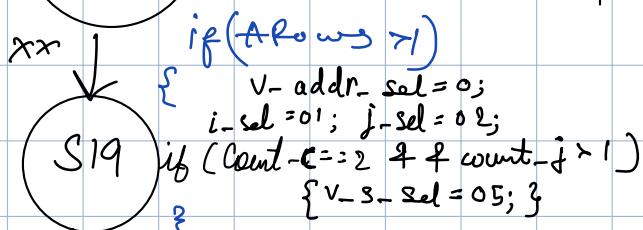
K-sel, J-sel, I-sel = 0;
Result-Sel \leq 0;
result-max \leq 1;
C-SRAM-Write-en = 1;
addr-C-sel = 1;



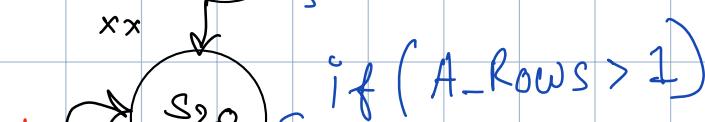
← Copying Matrix V to scratchpad SRAM
in column Major order.



if (A_Rows > 1)
{ V_addr_Sel = 01;
V_S_Sel = 06;
*}
*else { V_S_Sel = 01; }**



else { j_Sel = 01; i_Sel = 01; }

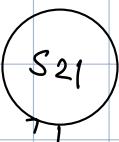


if (Count_i == 2 & Count_j > 1)
{ V_S_Sel = 05; }
if (Count_i == 1 & Count_j > 1)
{ V_addr_Sel = 1;
V_S_Sel = 06;
i_Sel = 03;
j_Sel = 01;
}

else { scratch_Write_en = 1;
scratch_addr_Sel = 1;
j_Sel = 01;
V_S_Sel = 01;
*}
if (Count_j == 1)
{ j_Sel = 03;
V_S_Sel = 08;
*}**

if (j <= 1)
{ next_state = S11; }
else { next_state = S10; }

(j == r & i == 1)



$V.S.Sel = 02;$ scratch.addr.sel = 2;
 $i.Sel = 02;$ $j.Sel = 02;$ scratch.write.en = 1;

scratch.addr.sel = 0; // copying of Matrix V is completed here.

xx

S23

result.input.sel = 03;
 $V.S.Sel = 4;$ scratch.read.addr.sel = 0;
Result.S.Z Sel = 0;

(A.Rows = 1)

S28

⋮

S24

$K.Sel, j.Sel, i.Sel = 05;$
 $V.S.Sel = 2;$ scratch.read.address.sel = 01;

xx

$K > 0$

if ($K \geq 0$)
{ $k.Sel = 01;$
 $j.Sel = 02;$
 $i.Sel = 02;$

else { $K.Sel = 05;$

if ($j > 1$)
{ $j.Sel = 1;$

else
 $j.Sel = 05;$

if ($i > 1$)
{ $i.Sel = 1;$

else { $i.Sel = 05;$

if ($K \geq 2$)
{ $V.S.Sel = 1;$
scratch.read.addr.sel = 1;

else { $V.S.Sel = 2;$
scratch.read.addr.sel = 2;

if ($K \geq 0$)
{ result.S.Z Sel = 1;

else { result.S.Z Sel = 2;

}

$(K \geq 0)$

}

↓

Starting convolution of Matrix S & V.
Computing Attention Matrix (Z).

S22

xx

xx

S23

xx

xx

S24

xx

xx

S25

xx

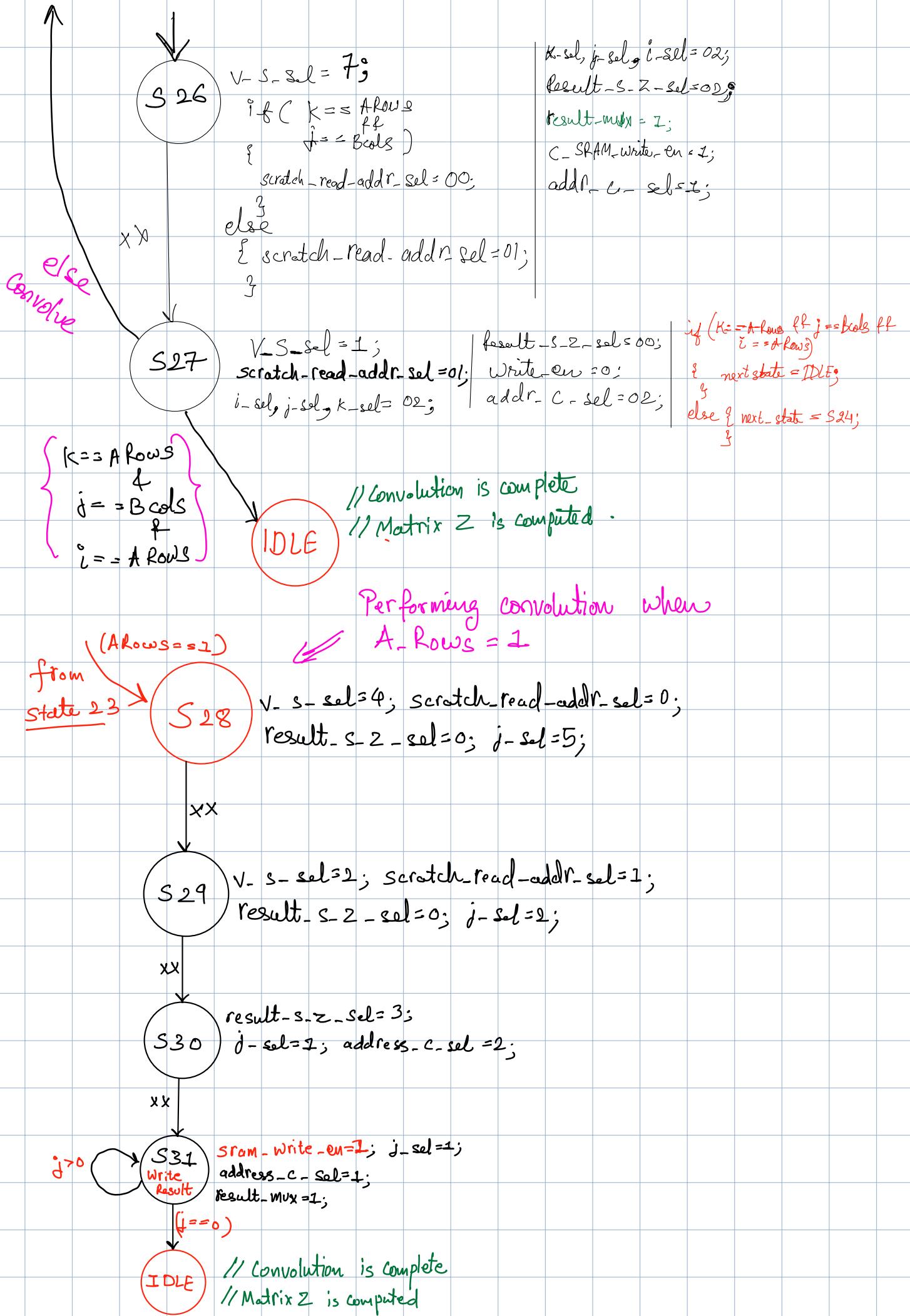
xx

⋮

S28

xx

xx



5. Results Achieved

- **Clock period:** 5.80 ns
- **# cycles":** 3354
- **Delay (clock period * #cycles):** 19453.2 ns
- **Area (um²):** 12866.6859 um²
- **1/(delay*area) ns⁻¹.um⁻² :** 4×10^{-9} ns⁻¹.um⁻²
- **Timing constraints:** All constraints MET.

6. Conclusions

- In this project, I understood the concept of self attention in LLM transformers and also implemented the RTL design using Verilog (2001) to realize the scaled dot product attention.
- The Finite state machine working as per Moore machine is implemented using binary encoding and consists of 32 states.
- Concept of pipelining is implemented for minimizing delay and efficient computation.
- Successfully synthesized the design using Synopsys with no unintentional latches, meeting all the timing constraints, optimal clock period and optimized area.