

# Project 2 - Matrix Multiplication Unit

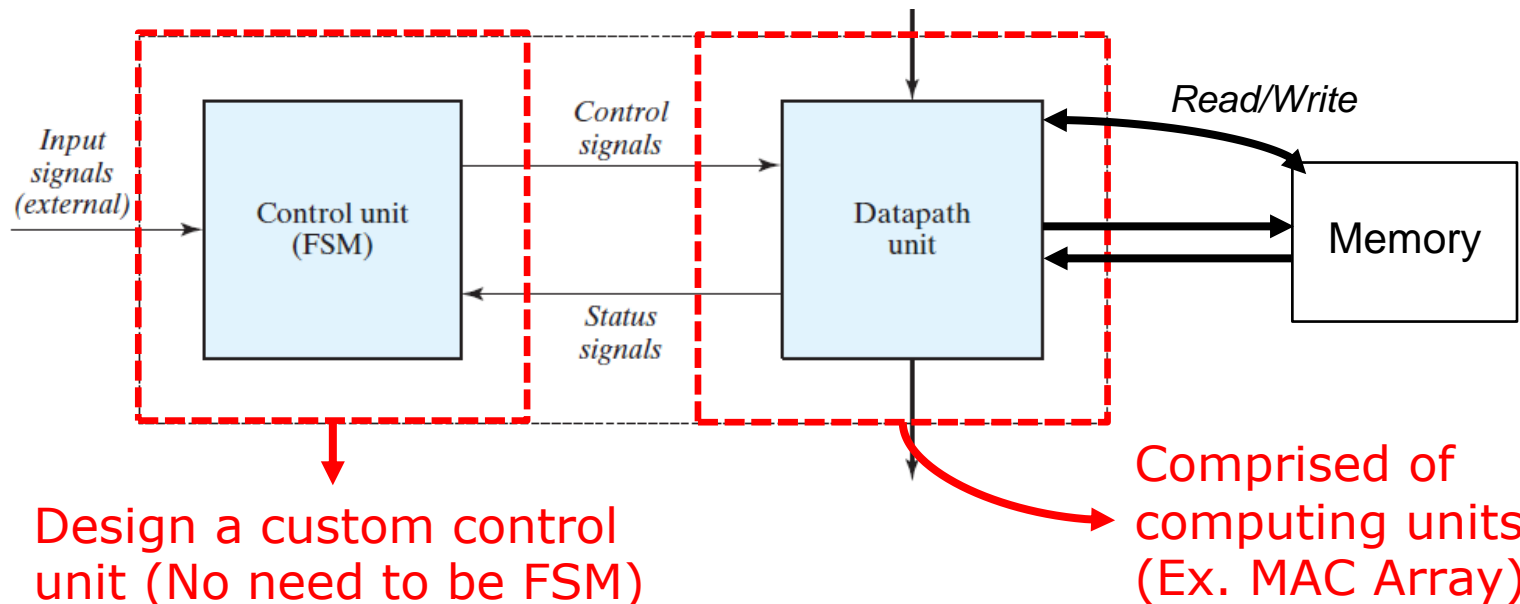
EE361 – Digital Circuit Design and Language

Kyung Hee University  
Electrical Engineering

Spring 2025  
Seungkyu Choi (seungkc@khu.ac.kr)

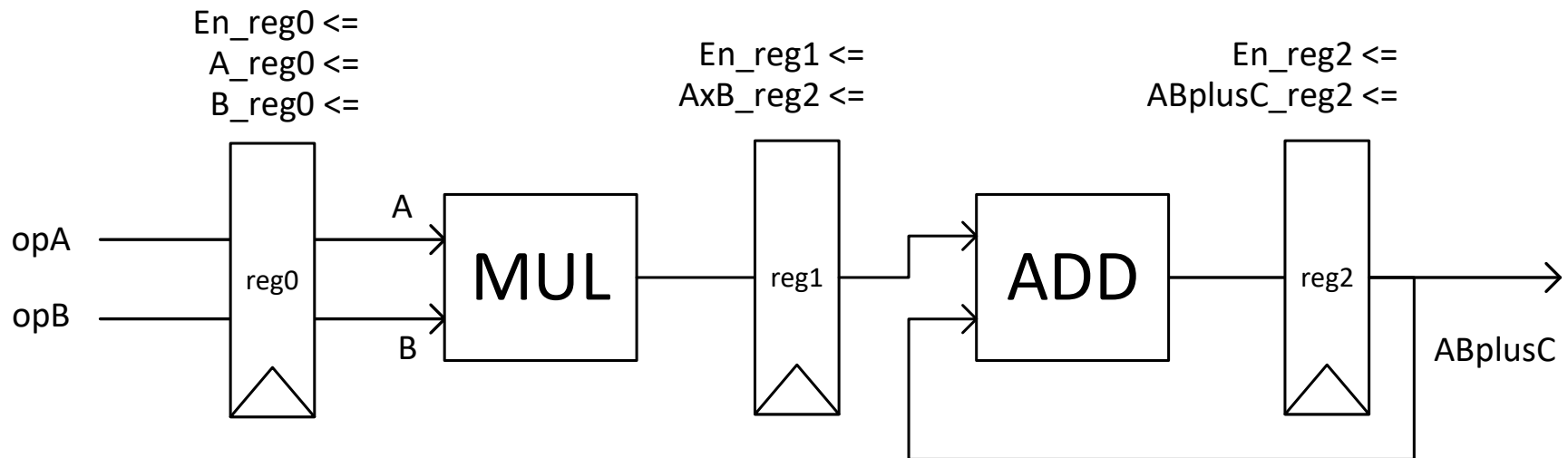
# System Design

- Control unit + datapath + memory structure
  - ✓ The datapath will be the processing part, which is comprised of multiple computing units.
  - ✓ The control unit can be implemented in any forms.
  - ✓ Input memory for read, output memory for write



# Utilizing a MAC Unit Design

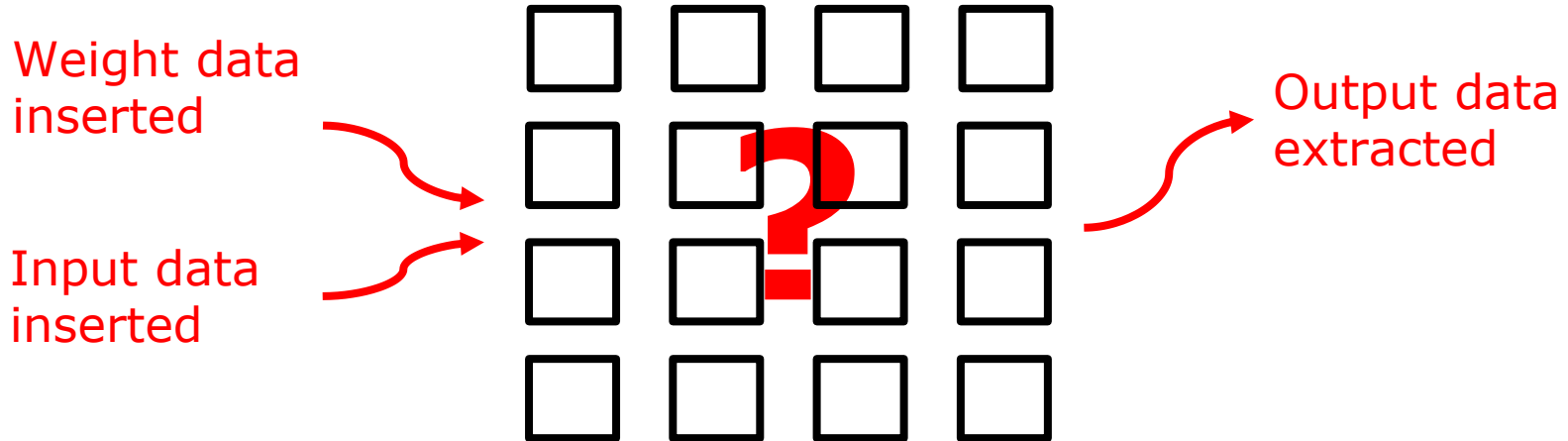
- A Pipelined Multiply-and-Accumulate (MAC) Unit
  - ✓ The MAC unit must be implemented with a pipelined system. (+MUL and ADD must be pipelined)
  - ✓ How to integrate the control system?



# Prj 2. 4x4 MAC Array Accelerator

- **$I (T \times N) \times W (N \times M) = O (T \times M)$  Matrix Multiplier**  
 **$(M, N, T = 1 \sim 8)$**

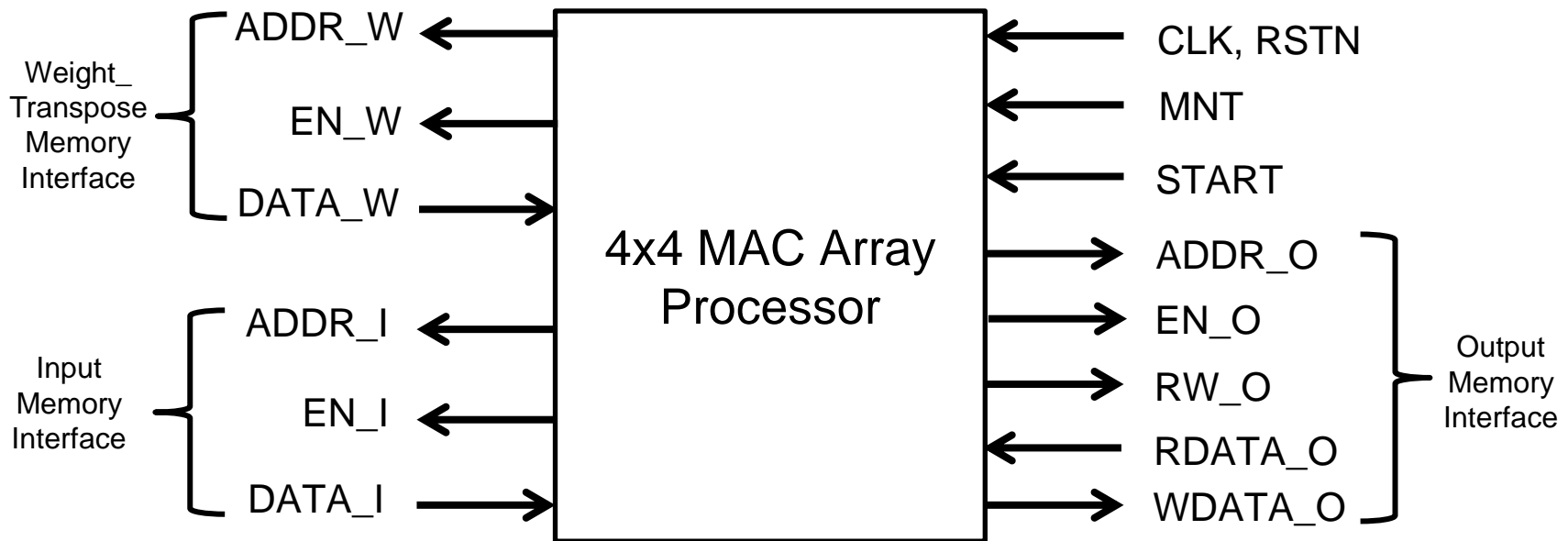
**Ex.**



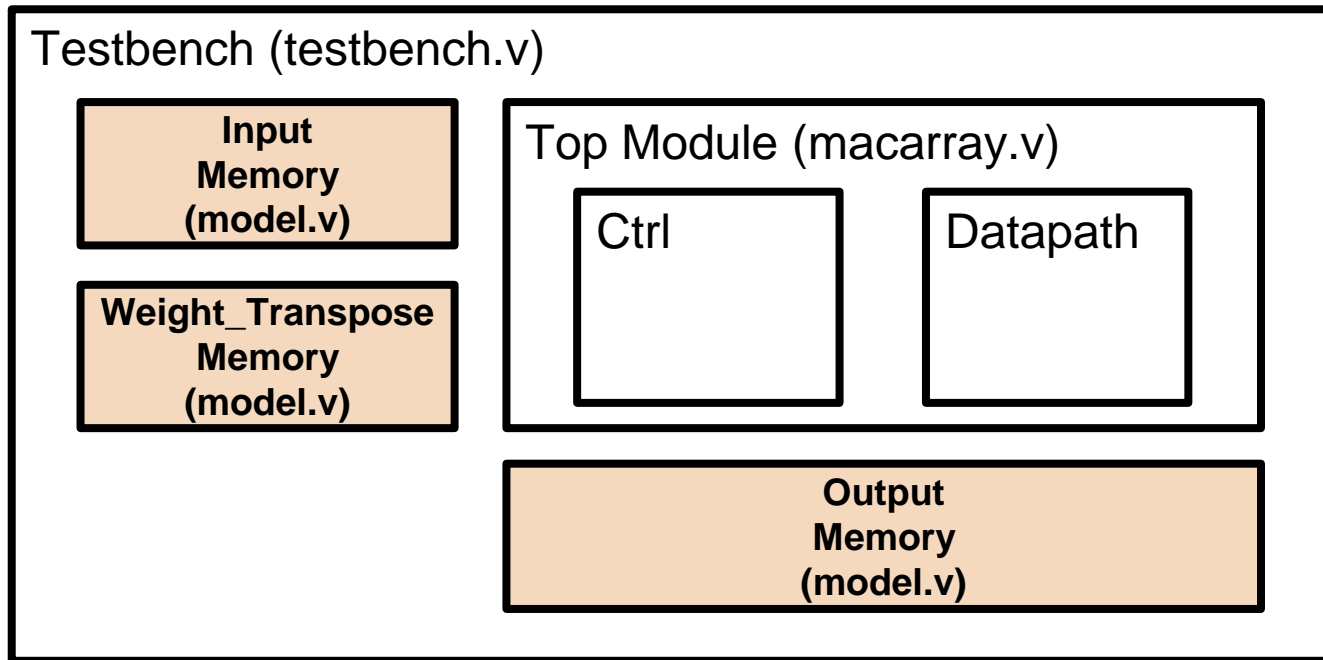
- Prepared  $(T \times N)$  input matrix and  $(N \times M)$  weight matrix are loaded from the input and weight memory, respectively.
- Output should be stored into the output memory in groups of 4: 4 output data are stored in the same address.
- Detailed load and store processes of the data will be explained.
- $M, N, T$  is a variable that can be any value from 1 to 8.

# Prj 2. Top Module Interface

<b>In CLK, RSTN</b>	Clock	<b>Out EN_W</b>	Weight_Tr Mem. Enable
<b>In [11:0] MNT</b>	4-bit M,N,T values	<b>In [31:0] RDATA_W</b>	Weight_Tr Data Bus
<b>In START</b>	Start Signal	<b>Out [3:0] ADDR_O</b>	Out Mem. Address
<b>Out [3:0] ADDR_I</b>	Input Mem. Address	<b>Out EN_O</b>	Out Mem. Enable
<b>Out EN_I</b>	Input Mem. Enable	<b>Out RW_O</b>	Out Mem. Read(0)/Write(1)
<b>In [31:0] RDATA_I</b>	Input Data Bus	<b>In [63:0] RDATA_O</b>	Out Mem. Read Data Bus
<b>Out [3:0] ADDR_W</b>	Weight_Tr Mem. Address	<b>Out [63:0] WDATA_O</b>	Out Mem. Write Data Bus



# Prj 2. Testbench I/O



Type	Width	Name	
Input	1	CLK	Clock
Input	1	RSTN	Asynchronous reset negative
Input	1	Start	Start signal
Input	12	MNT	4-bit concatenated M, N and T values ranging from 1~8 (Ex. 12'h568 means M=5, N=6, T=8)

# Prj 2. Memory Allocation Rule

- Read and Write data should be Input, Weight, Output participating the matrix multiplication process shown below.

$$\begin{array}{c}
 \begin{array}{c} \text{N=1~8} \\ \left[ \begin{array}{ccccc} x_{11} & x_{12} & x_{13} & \dots & x_{1N} \\ x_{21} & x_{22} & x_{23} & \dots & x_{2N} \\ x_{31} & x_{32} & x_{33} & \dots & x_{3N} \\ \dots & \dots & \dots & \dots & \dots \\ x_{T1} & x_{T2} & x_{T3} & \dots & x_{TN} \end{array} \right] \end{array} \\
 \begin{array}{c} \text{T=1~8} \end{array}
 \end{array}
 \times
 \begin{array}{c}
 \begin{array}{c} \text{M=1~8} \\ \left[ \begin{array}{ccccc} \omega_{11} & \omega_{12} & \omega_{13} & \dots & \omega_{1M} \\ \omega_{21} & \omega_{22} & \omega_{23} & \dots & \omega_{2M} \\ \omega_{31} & \omega_{32} & \omega_{33} & \dots & \omega_{3M} \\ \dots & \dots & \dots & \dots & \dots \\ \omega_{N1} & \omega_{N2} & \omega_{N3} & \dots & \omega_{NM} \end{array} \right] \end{array} \\
 \begin{array}{c} \text{N=1~8} \end{array}
 \end{array}
 \\
 \\
 \begin{array}{c}
 \begin{array}{c} \text{M=1~8} \\ \left[ \begin{array}{ccccc} o_{11} & o_{12} & o_{13} & \dots & o_{1M} \\ o_{21} & o_{22} & o_{23} & \dots & o_{2M} \\ o_{31} & o_{32} & o_{33} & \dots & o_{3M} \\ \dots & \dots & \dots & \dots & \dots \\ o_{T1} & o_{T2} & o_{T3} & \dots & o_{TM} \end{array} \right] \end{array} \\
 \begin{array}{c} \text{T=1~8} \end{array}
 \end{array}
 \end{array}$$

# Prj 2. Input and Weight Utilization

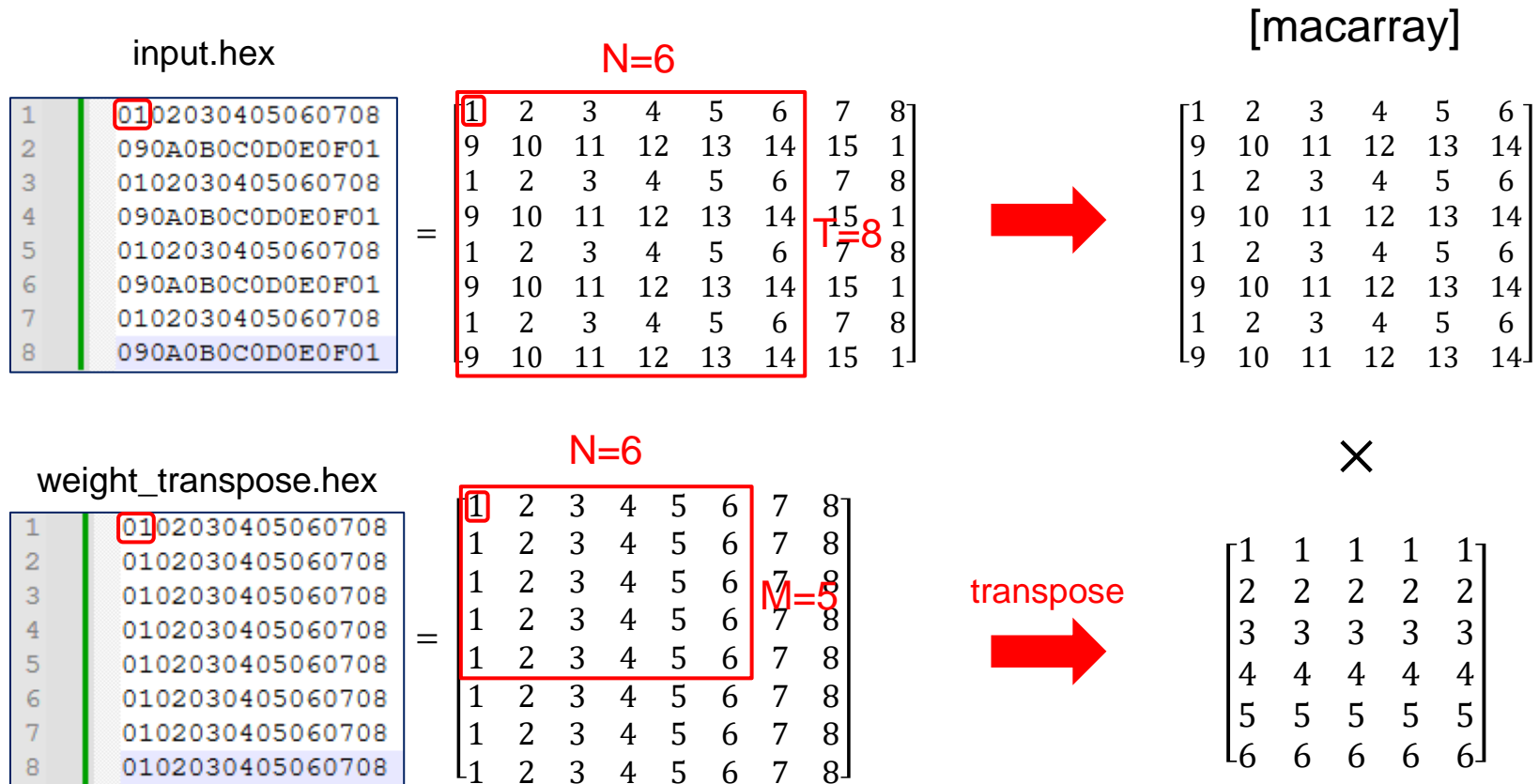
- Input and weight should be prepared with hex files as 8x8 matrices. The file names 'input.hex' and 'weight\_transpose.hex' should not be changed.  
(Only M,N,T signals can control the sizes of input and weight.)
- A single data is 8-bit 2's complement data (2 hex num).
- Be careful of the weight since it should be stored and inserted in a transposed version.

```
// -----  
// Input (T X N) and Weight (N X M) test matrices should be stored as.  
// -----  
// Caution : Assumption : input files have hex data like below.  
//      Input      : (1,1) (1,2) ... (1,N)  
//      (T x N)     (2,1) (2,2) ... (2,N)  
//                  (3,1) (3,2) ... (3,N)  
//                  ...   ...   ...   ...  
//                  (T,1) (T,2) ... (T,N)  
//  
//      Weight_transpose : (1,1) (1,2) ... (1,N)  
//      (M x N)           (2,1) (2,2) ... (2,N)  
//                        (3,1) (3,2) ... (3,N)  
//                        ...   ...   ...   ...  
//                        (M,1) (M,2) ... (M,N)  
//  
// hex files will be -> line1: (1,1) (1,2) (1,3) (1,4) (1,5) (1,6) (1,7) (1,8)  
//                    -> line2: (2,1) (2,2) (2,3) (2,4) (2,5) (2,6) (2,7) (2,8)  
//                    -> line3: (3,1) (3,2) (3,3) (3,4) (3,5) (3,6) (3,7) (3,8)  
//                    ...   ...   ...   ...   ...   ...   ...   ...  
//                    -> line8: (8,1) (8,2) (8,3) (8,4) (8,5) (8,6) (8,7) (8,8)
```



# Prj 2. Input and Weight Utilization

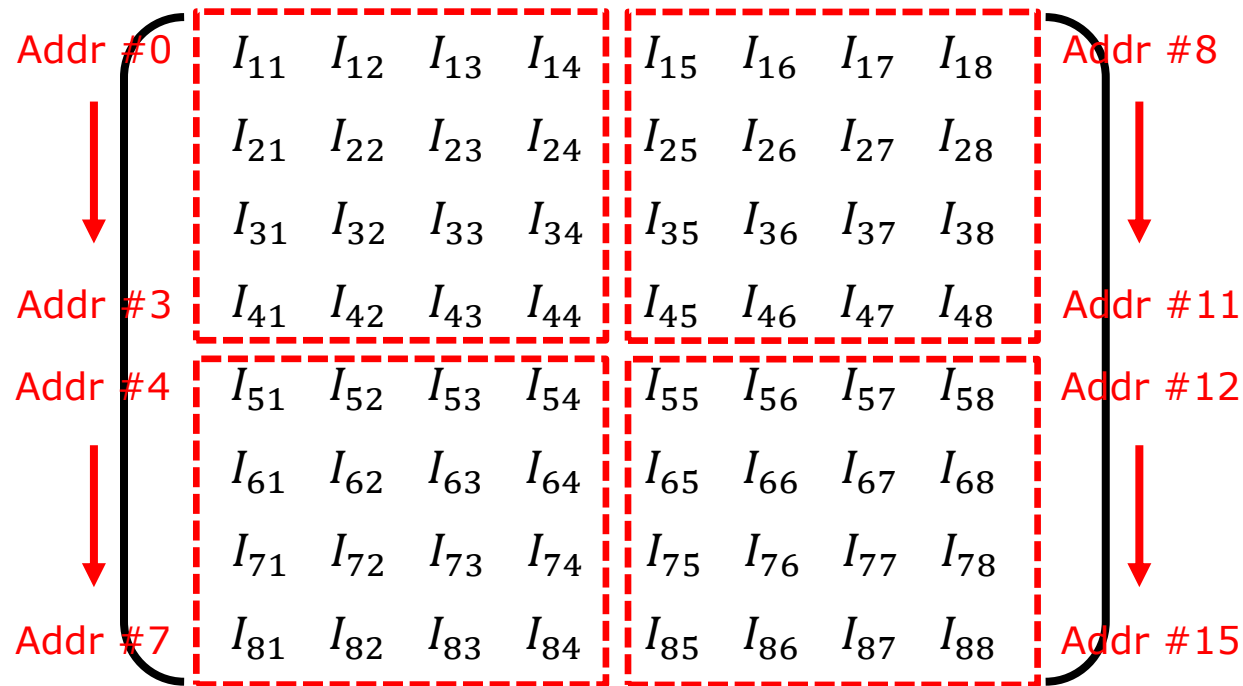
- Example (MNT=12'h568)



# Prj 2. In/W Memory Allocation Rule

- The data from hex file is automatically stored as below, having a 32-bit width per address.
- Input and weight have a fixed read bandwidth of 32 bits per cycle.

## Example) *Input*



[Same for *Weight*]

# Prj 2. In/W Memory Allocation Rule

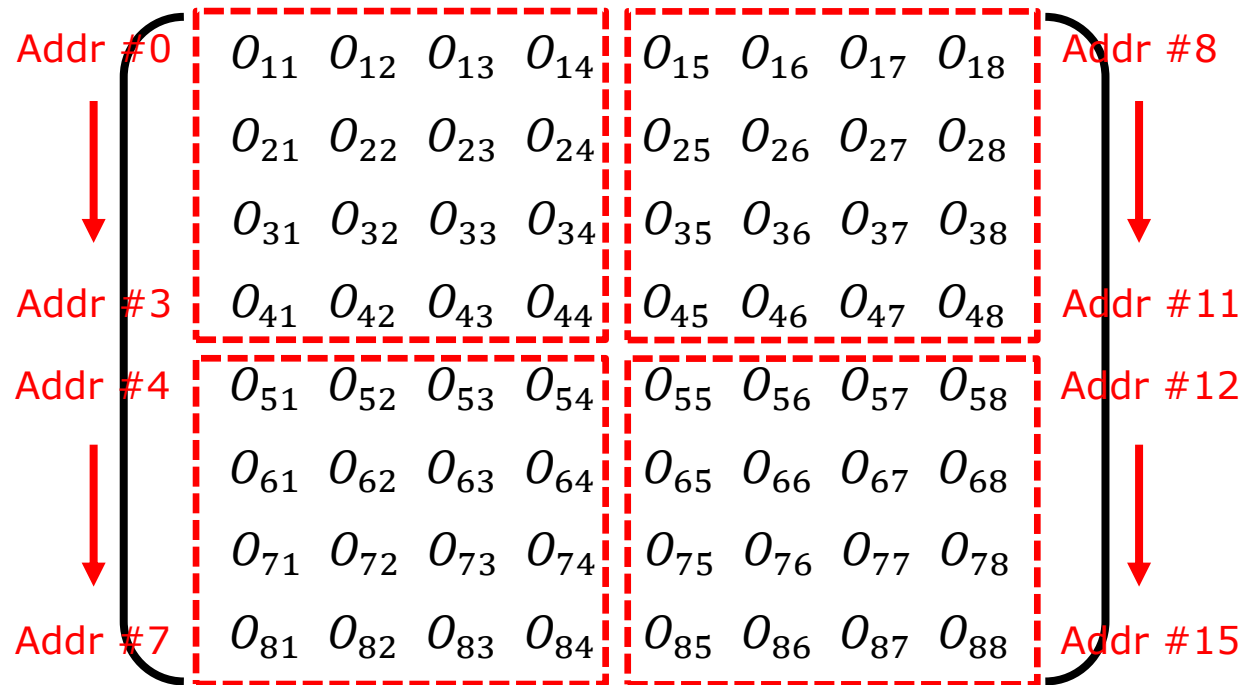
## ■ Example

Addr #0	01	23	45	67	89	5b	00	4f	Addr #8
	01	00	03	04	05	06	07	08	
	00	30	20	10	90	50	40	12	
Addr #3	3e	2c	1a	00	00	43	21	98	Addr #11
Addr #4	02	48	0c	51	35	79	4d	30	Addr #12
	75	53	00	59	85	52	45	56	
	01	02	03	04	05	06	07	08	
Addr #7	09	0a	0b	0c	0d	00	0f	00	Addr #15

/testbench/INPUT_MEM/ram	01...	01234567 01000304 00302010 3e2c1a00
[0]	01...	01234567
[1]	01...	01000304
[2]	00...	00302010
[3]	3e...	3e2c1a00
[4]	02...	02480c51
[5]	75...	75530059
[6]	01...	01020304
[7]	09...	090a0b0c
[8]	89...	895b004f
[9]	05...	05060708
[10]	90...	90504012
[11]	00...	00432198
[12]	35...	35794d30
[13]	85...	85524556
[14]	05...	05060708
[15]	0d...	0d000f00
/testbench/INPUT_MEM/readmem_inst	01...	01234567895b004f 0100030405060708 0
[0]	01...	01234567895b004f
[1]	01...	0100030405060708
[2]	00...	0030201090504012
[3]	3e...	3e2c1a0000432198
[4]	02...	02480c5135794d30
[5]	75...	7553005985524556
[6]	01...	0102030405060708
[7]	09...	090a0b0c0d000f00

# Prj 2. Output Memory Allocation Rule

- Output data is stored in the memory following the rule below, having a 64-bit width per address (16-bit width 4 data).
- Addressing is fixed, and unused row/column data should be set to zero.



# Prj 2. Output Monitoring

---

- OUT should be stored as the same manner stated in the memory allocation rule.
- A single data is 16-bit 2's complement data (4 hex num).
- You do not need to consider overflow cases.
- You should not modify any of the testbench file and model file to access output directly when submitting your source code.

# Prj 2. SDC table

Frequency (MHz)	Input Delay(ns)	Output Delay(ns)
~ 50	4.0 ~ 6.0	4.0 ~ 6.0
50 ~ 100	2.0 ~ 4.0	2.0 ~ 4.0
100 ~ 150	1.5 ~ 2.5	1.5 ~ 2.5
150 ~ 200	1.0 ~ 2.0	1.0 ~ 2.0

# Submission Rule (Due:6/16 11:59 PM)

- Submission file: Zip File - Verilog codes and reports

*File Name: prj2\_xx.zip (xx: team number 01~24)*

- Verilog Codes: **every files**  
(Do not change the name of the skeleton files.)
- Do not modify any of the testbench and model files when submitting your code.
- You can test by changing the values of the weight and input hex files. Be careful to follow the same rule indicated in the testbench code when making the weight and input matrix. (!!! Weight should be stored in transposed version.)
- Two reports should be prepared: An **8 page-limit word/hwp file** and a **6 min long PPT**
  - Reports should include any details of your design.
  - Screen captures of your Quartus results and analysis are mandatory.  
(Device: CycloneIII EP3C120F484C7)
  - Indicate the overview of the control system and key strengths of your design in the PPT file. Submitted version cannot be modified and will be the final version that will be used for the presentation.