# AOC Proposal
# Systolic Tensor Array

## Group｜把家齊高中還給家齊

### Member

| | | |
|---|---|---|
| 電機所114 | 張育誠 | N26121892 |
| 工科系113 | 劉子齊 | E94096209 |
| 電機系113 | 洪翊碩 | E24091027 |
| 電機系113 | 陳喬雅 | E14093140 |
| 電機系113 | 鄭喆嚴 | E14096724 |

# TABLE OF CONTENTS

**01** Model, Quantization and Pruning

**02** PE Architecture and Dataflow

**03** Our Accelerator Architecture and SPEC
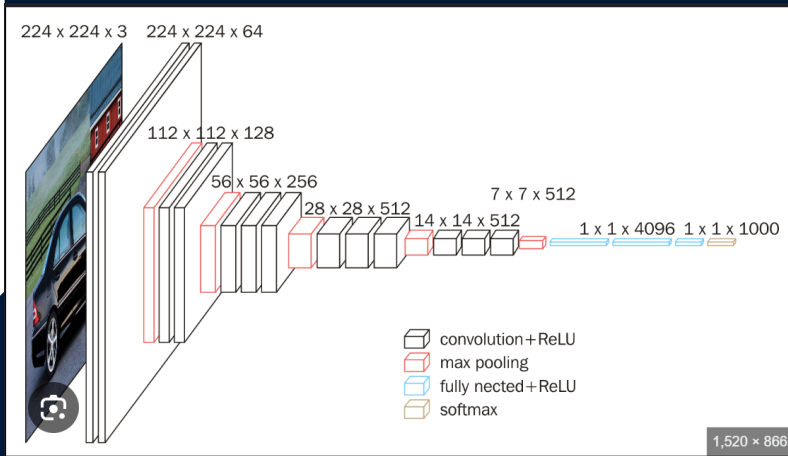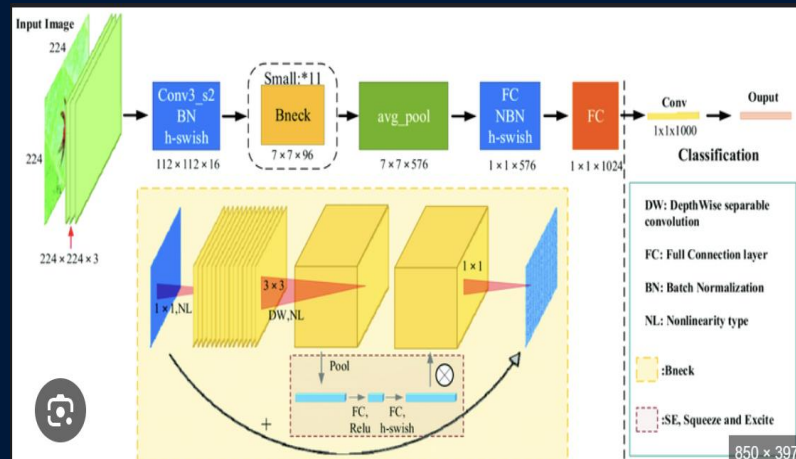
**04** Analysis and Question

# 01

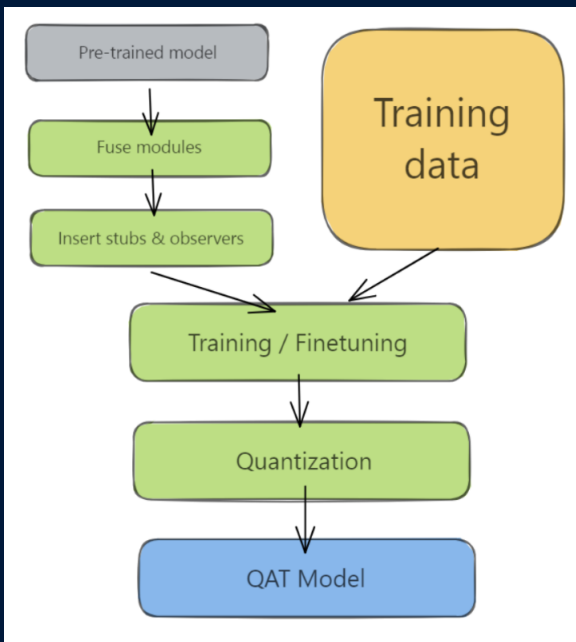## Model, Quantization and Pruning

# Model

## VGG16:



224 x 224 x 3  224 x 224 x 64

112 x 112 x 128

56 x 56 x 256

28 x 28 x 512

14 x 14 x 512

7 x 7 x 512

1 x 1 x 4096  1 x 1 x 1000

- convolution+ReLU
- max pooling
- fully nected+ReLU
- softmax

1,520 × 866

## Mobile NetV3:



Input Image

224

224

224 × 224 × 3

Conv3_s2 BN h-swish

112×112×16

Small:*11

Bneck

7 × 7 × 96

avg_pool

7 × 7 × 576

FC NBN h-swish

1 × 1 × 576

FC

1 × 1 × 1024

Conv 1x1x1000

Ouput

Classification

1 × 1,NL

3 × 3

DW,NL

1 × 1

Pool

FC, Relu  FC, h-swish

+

DW: DepthWise separable convolution

FC: Full Connection layer

BN: Batch Normalization

NL: Nonlinearity type

:Bneck

:SE, Squeeze and Excite

850 × 397
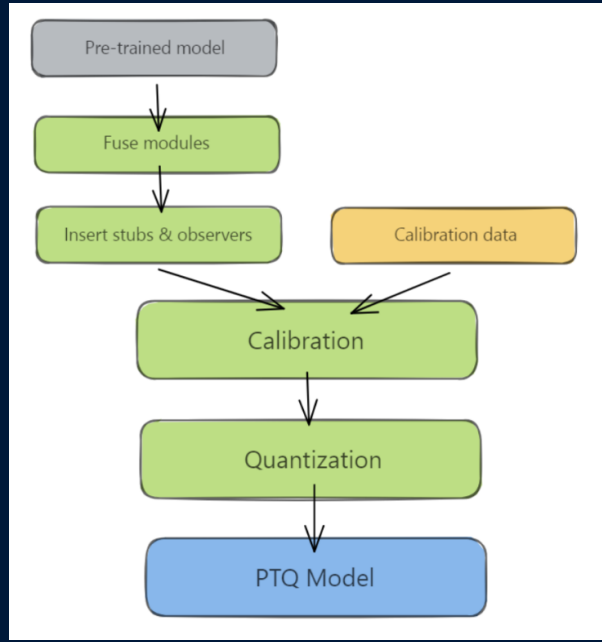
# Quantization

**QAT:**



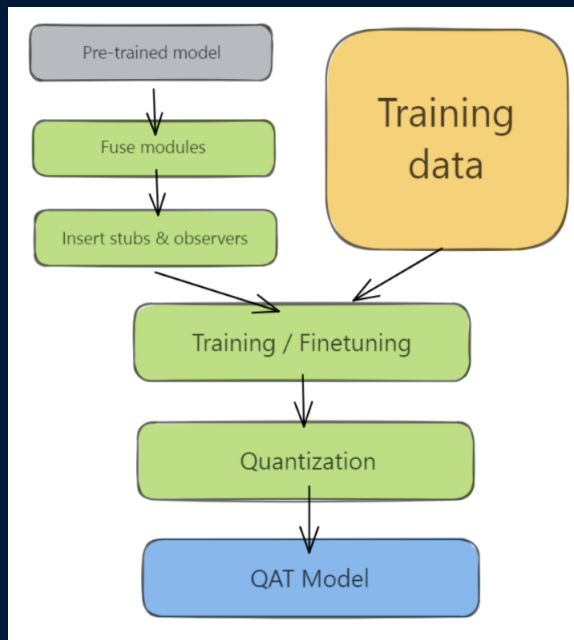**PTQ:**
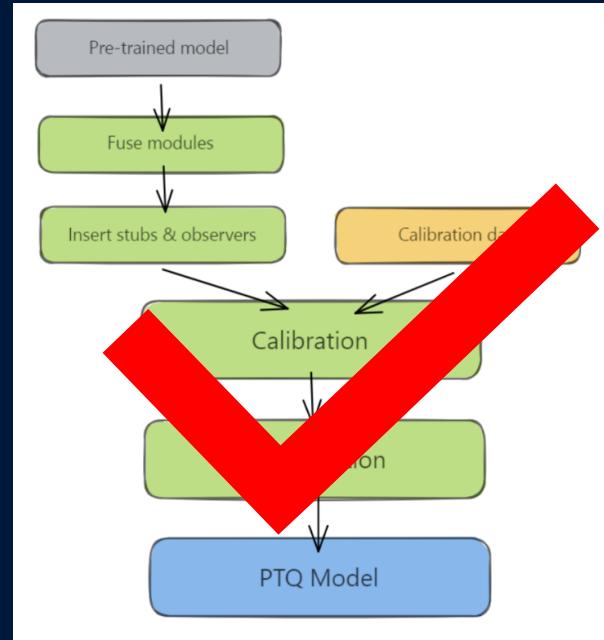
# Quantization

# Pruning

**Fix position pruning:**
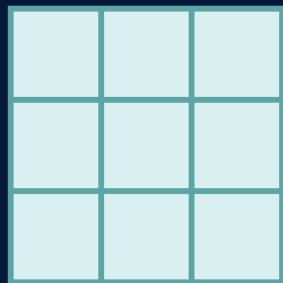
```python
if  SPARSITY:
        #  updateConv2D()
        #  train(1)
        test(sparsity_model)
        maskConv2D()
        test(model)
        test(sparsity_model)
```

# Pruning

## Original

```
          [[-1.2915e-04, -9.6068e-04, -1.1363e-03],
           [ 2.1713e-02,  7.5816e-03,  5.7231e-03],
           [ 2.7600e-03,  2.2870e-02,  3.2581e-03]],

          [[ 1.3669e-04,  3.5387e-02,  3.6681e-04],
           [-8.6679e-04,  3.6982e-03,  3.2041e-02],
           [ 5.4913e-03, -5.7838e-03,  1.4623e-03]],

          [[-9.4173e-04,  2.0093e-04,  6.6950e-04],
           [ 4.0101e-03, -1.5172e-02, -9.6157e-03],
           [ 5.9197e-03, -1.0201e-03, -1.6133e-02]]]], device='cuda:0')

Test set: Average loss: 0.3102, Accuracy: 9279/10000 (92.8%)


Test set: Average loss: 0.3081, Accuracy: 9283/10000 (92.8%)
```

# Pruning

## Saliency-based pruning:

```
if  SPARSITY:
        updateConv2D()
        train(1)
        test(sparsity_model)
        maskConv2D()
        test(model)
        test(sparsity_model)
```

# Pruning

## Saliency-based pruning:

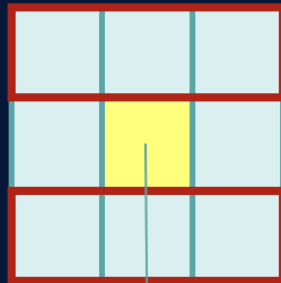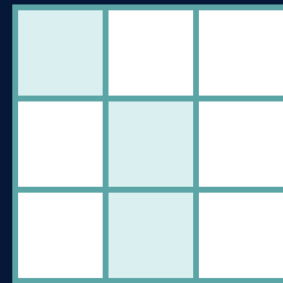According to scaling factor
Only keep the largest value

Sparsity

Prune

According to PatDNN
Center value is important, so it can't be pruned
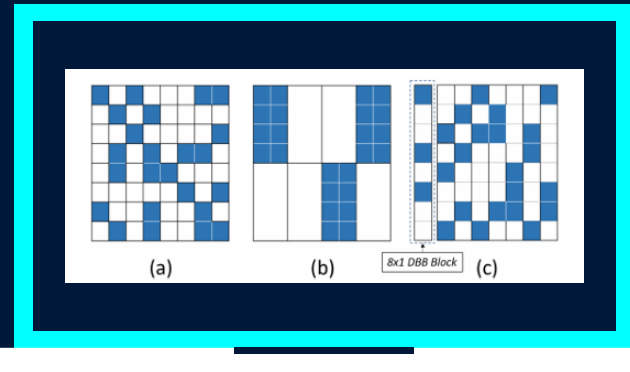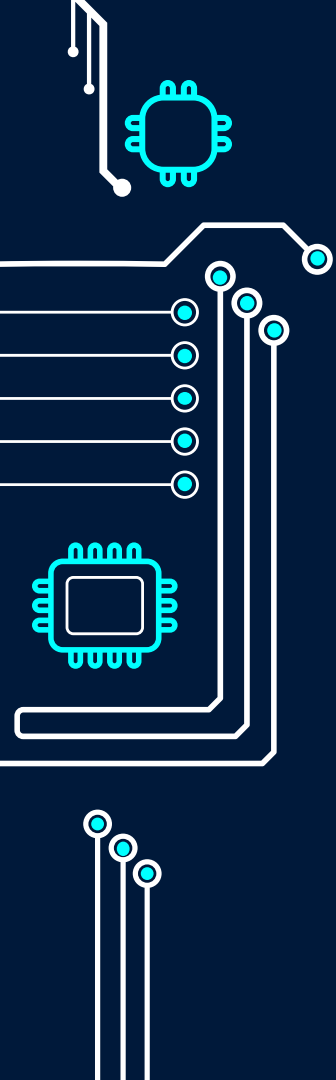
# 02

## PE Architecture and Dataflow

# GEMM Accelerator

## For Dense Matrix



(a) Conventional Systolic Array (SA)



**Systolic Tensor Array**

# GEMM Accelerator

## For Sparse Matrix

### DBB │ Density-Bound Block

Constraint mechanism for

the number of non-zero values

permissible in each block of a matrix.

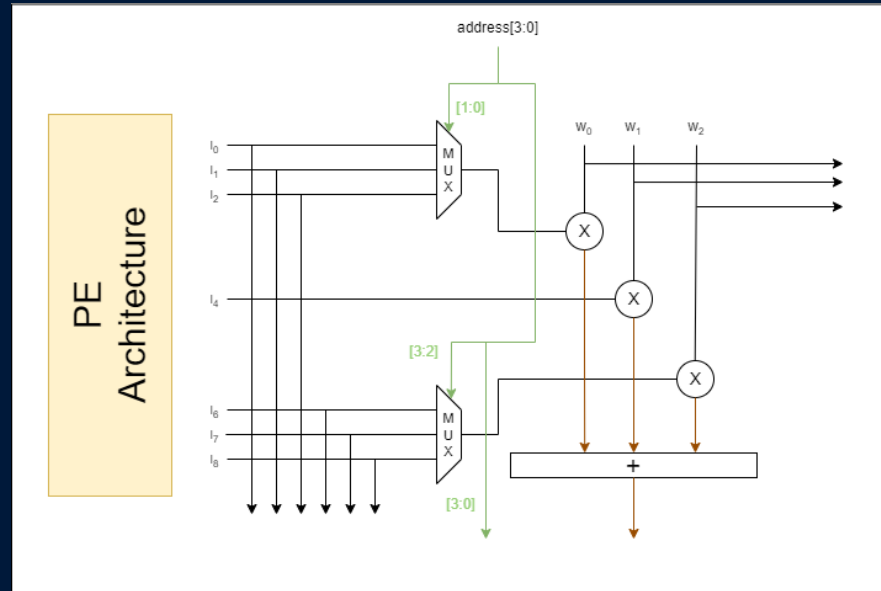# GEMM Accelerator

## For Sparse Matrix

### Our DBB Method



According to scaling factor
Only keep the largest value

Sparsity

Prune

According to PatDNN
Center value is important, so it can't be pruned

# GEMM Accelerator

## For Sparse Matrix



8x1 DBB Block

(c)



STA-DBB 2x4x2_2x2

Buffer#/MUL#=3/4
ACC#/MUL#=1/4

(c) Systolic Tensor Array for DBB (STA-DBB)

# PE Architecture

# PE I/O Ports

| Signal | I/O | Bits | Function |
|--------|-----|------|----------|
| ifmap_i | Input | 56 bits | $I_0$ -> [7:0]<br>$I_0$, $I_1$, $I_2$, $I_4$, $I_6$, $I_7$, $I_8$ |
| weight_i | Input | 24 bits | $W_0$ -> [7:0]<br>$W_0$, $W_1$, $W_2$ |
| address_i | Input | 4 bits | [1:0] for upper weight<br>[3:2] for lower weight |
| ifmap_o | Output | 56 bits | |
| weight_o | Output | 24 bits | |
| address_o | Output | 4 bits | |
| result | Output | 18 bits | Product and Add |

# PE Array Architecture

A for 3*3 Ifmap Sliding Window
B for 3*3 Weight
(num) num :which clock to send

# PE Array Architecture

A for 3*3 Ifmap Sliding Window
B for 3*3 Weight
(num) num :which clock to send
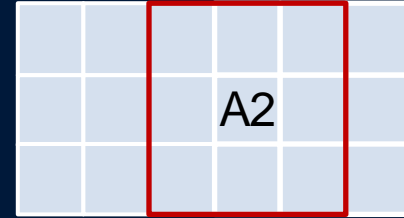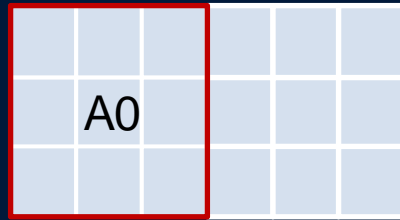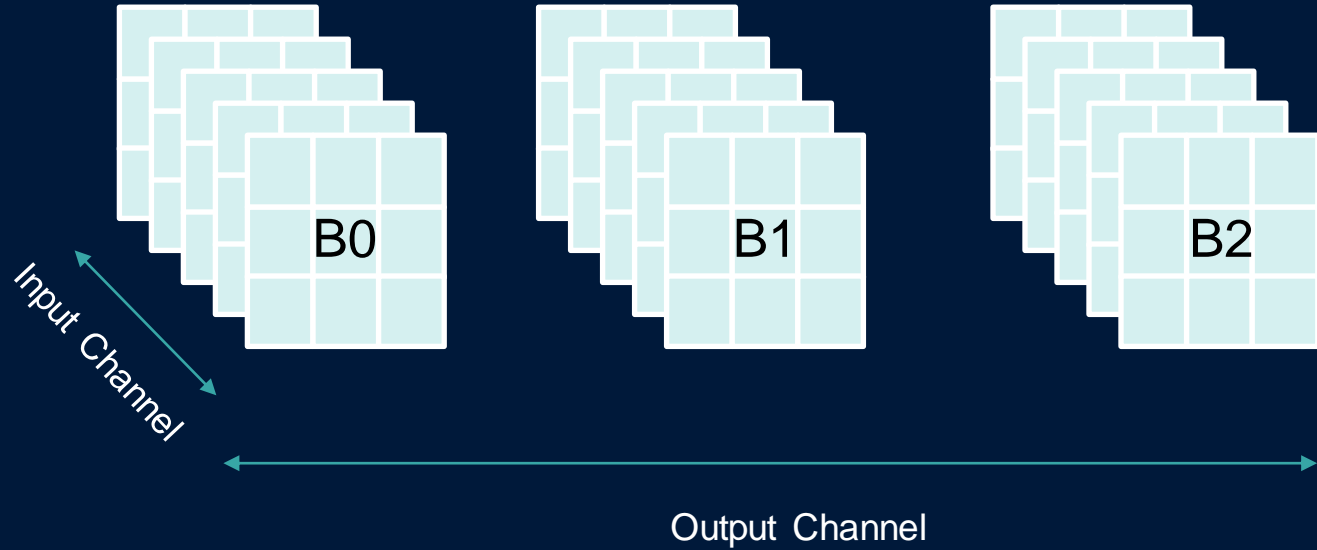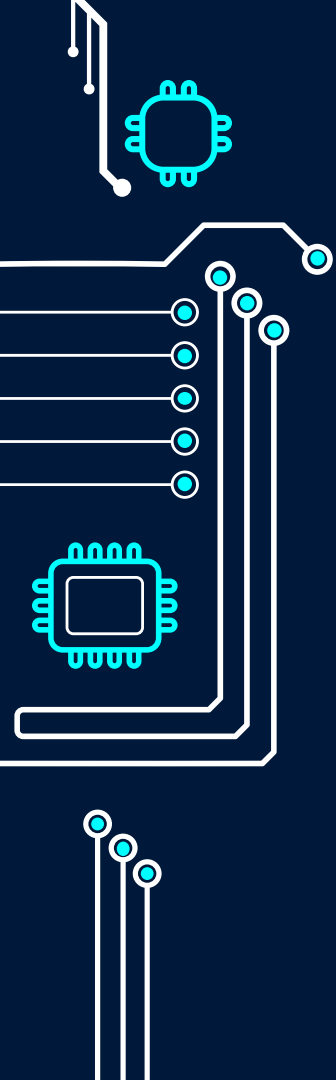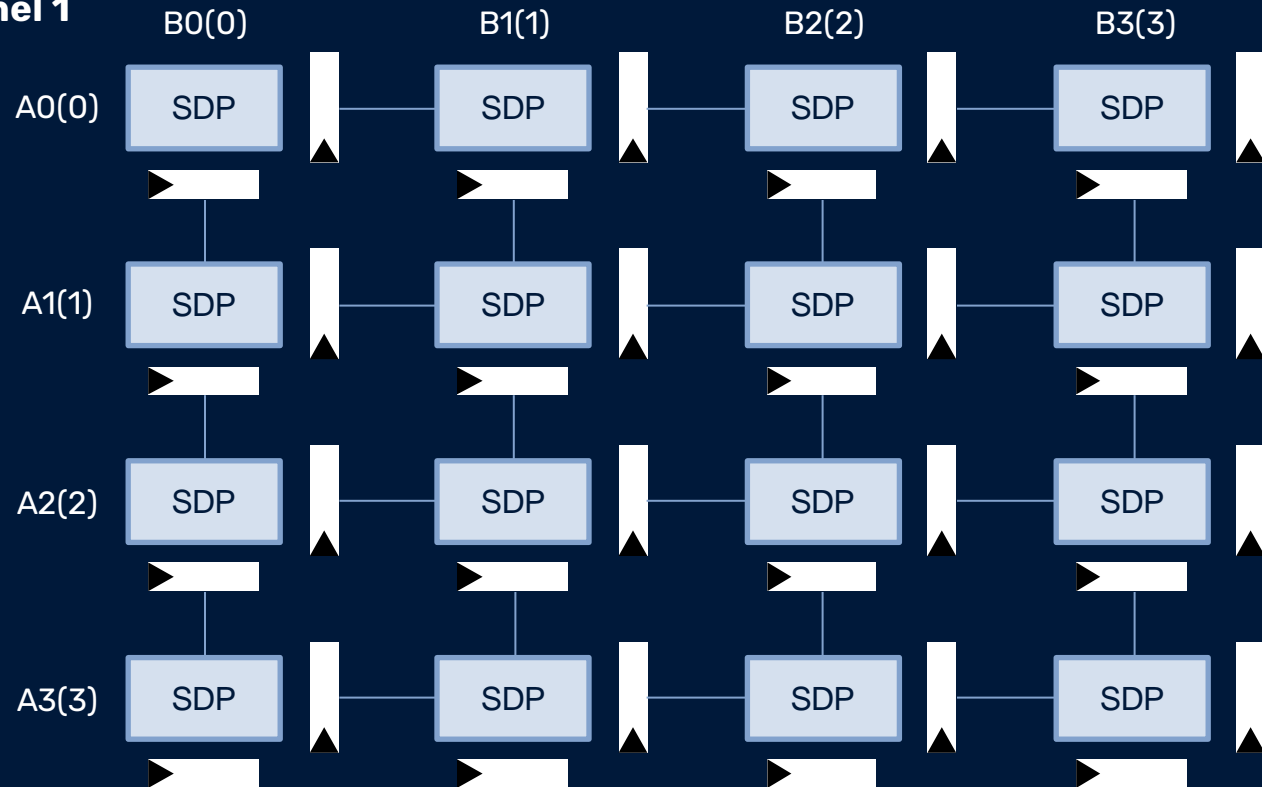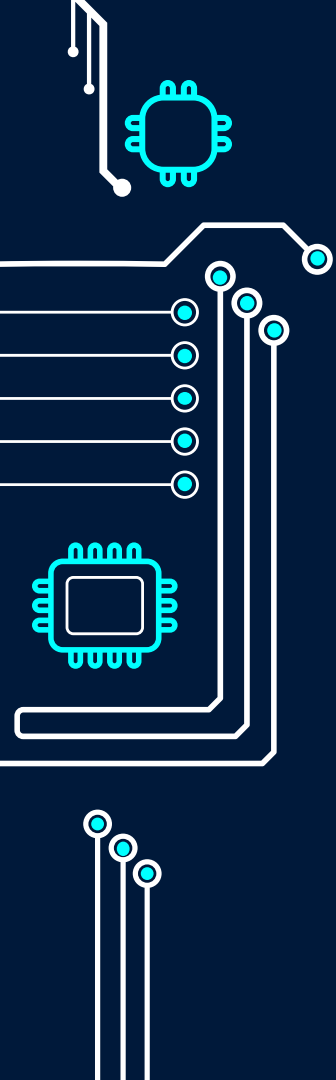
B0

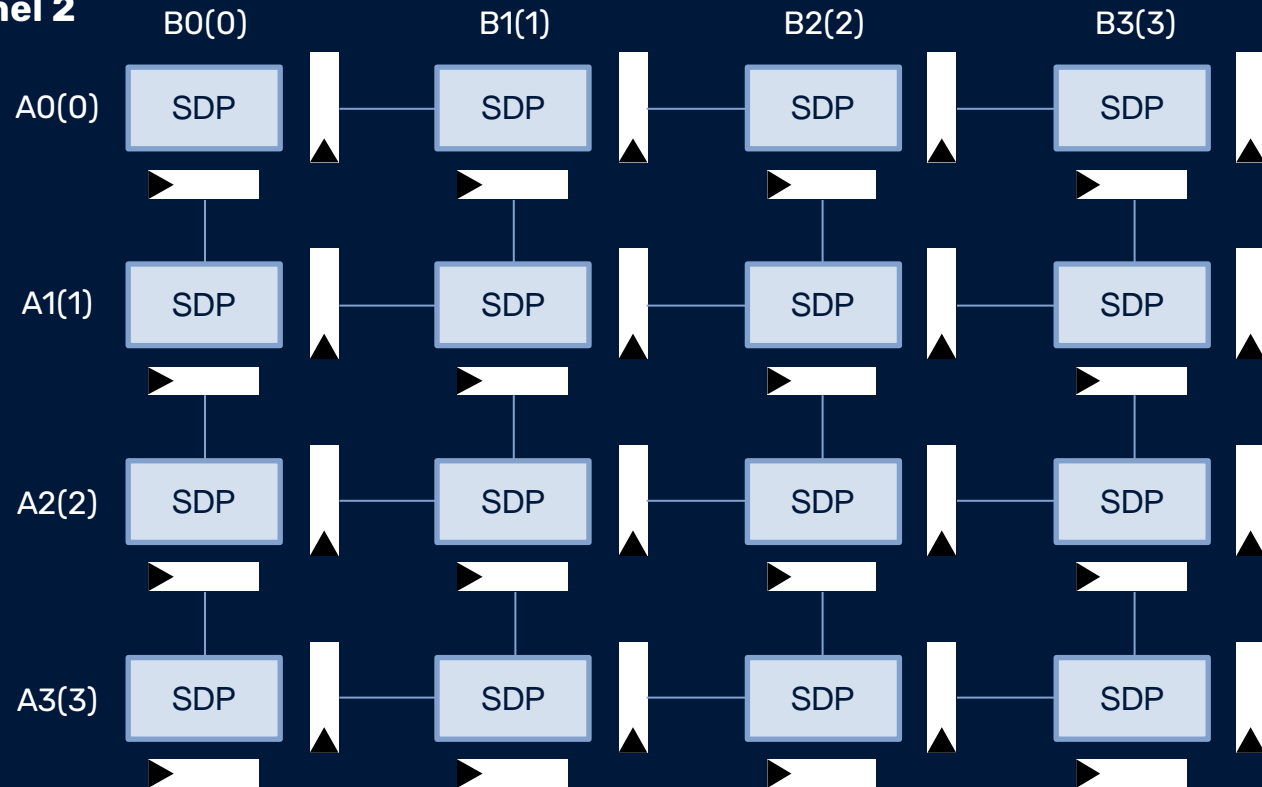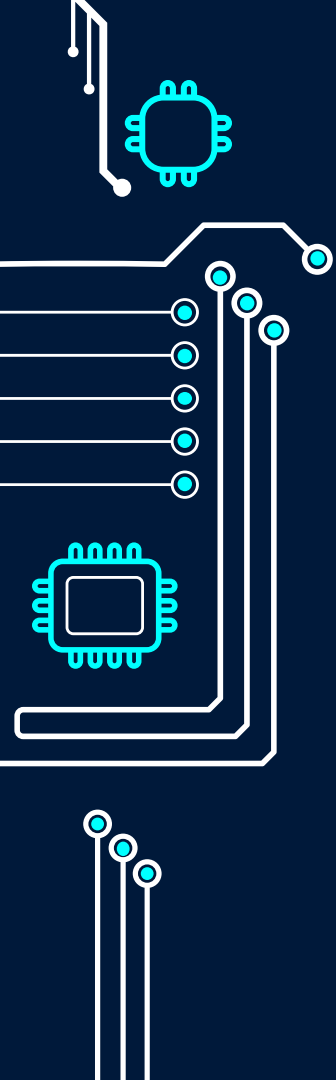B1

B2

Input Channel

Output Channel

# PE Array Architecture

A for 3*3 Ifmap Sliding Window
B for 3*3 Weight
(num) num :which clock to send

**Channel 2**

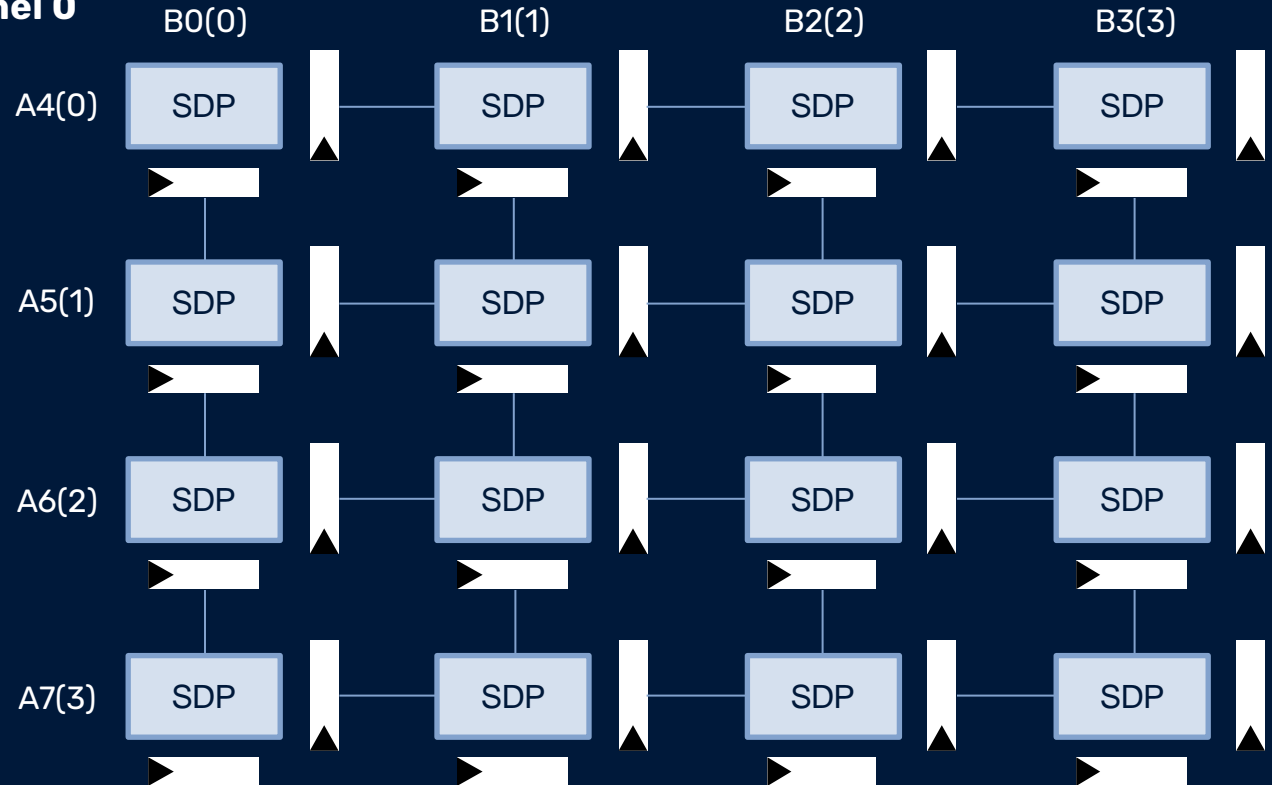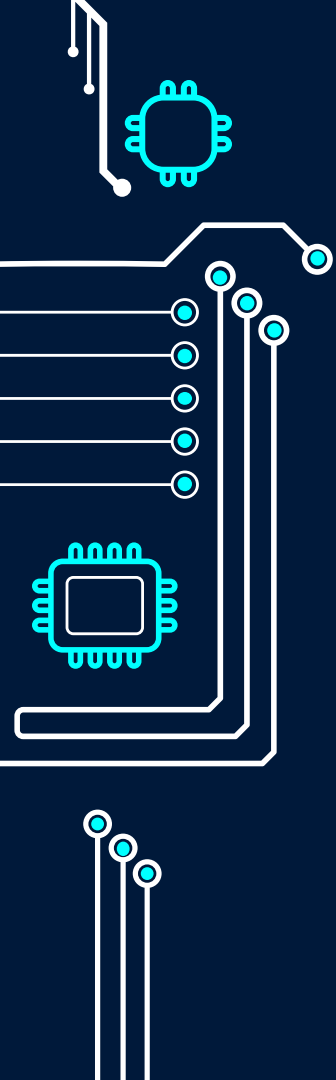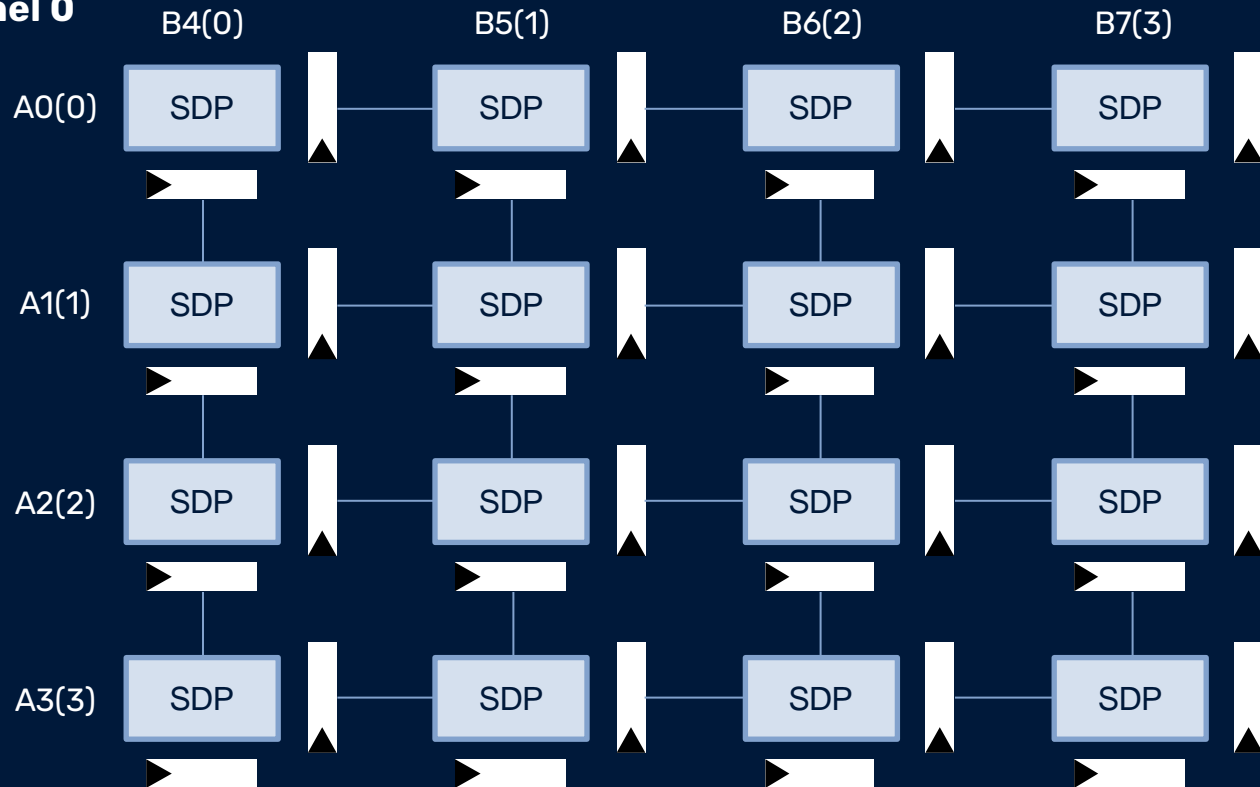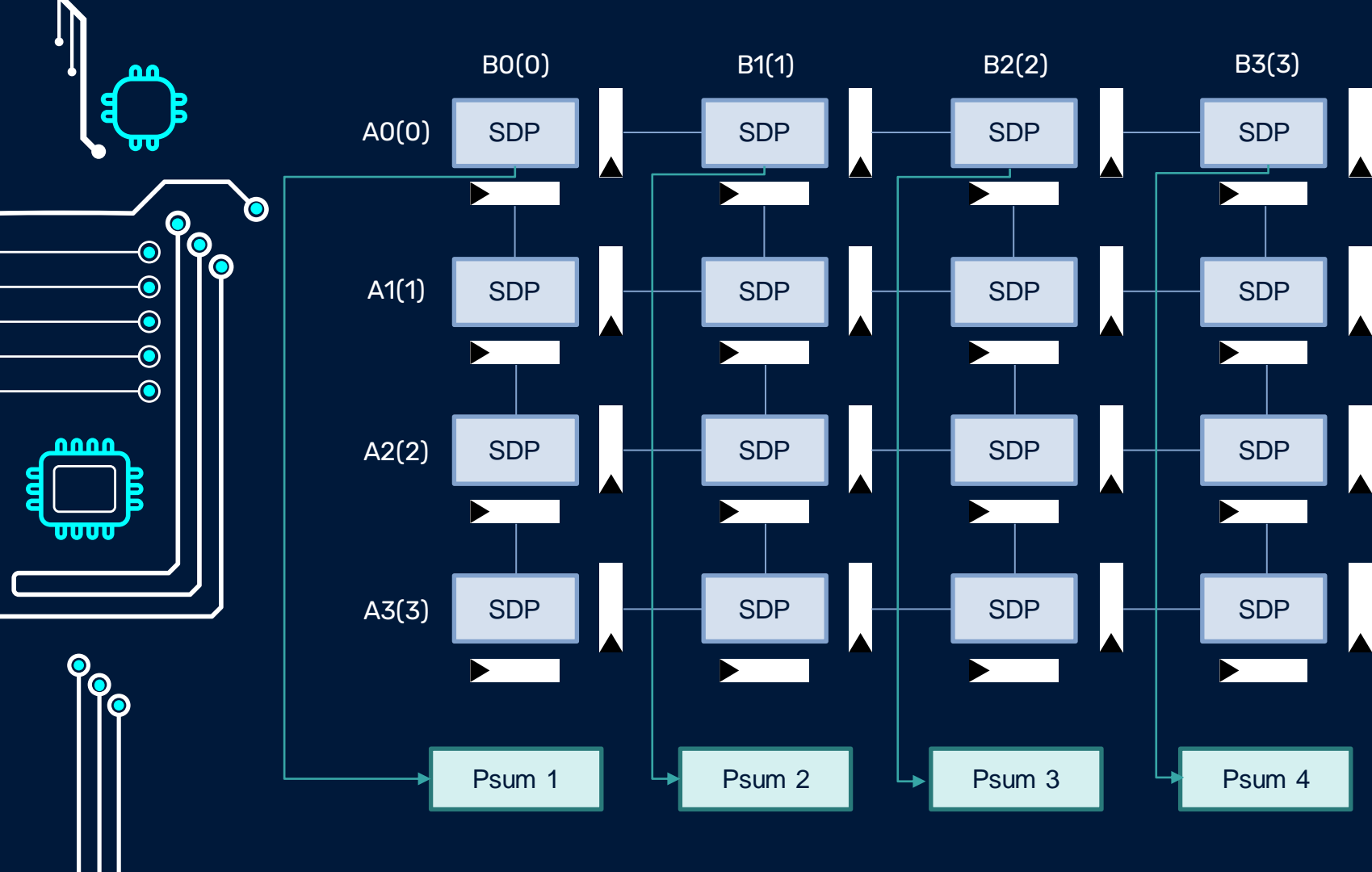|        | B0(0) | B1(1) | B2(2) | B3(3) |
|--------|-------|-------|-------|-------|
| A0(0)  | SDP   | SDP   | SDP   | SDP   |
| A1(1)  | SDP   | SDP   | SDP   | SDP   |
| A2(2)  | SDP   | SDP   | SDP   | SDP   |
| A3(3)  | SDP   | SDP   | SDP   | SDP   |

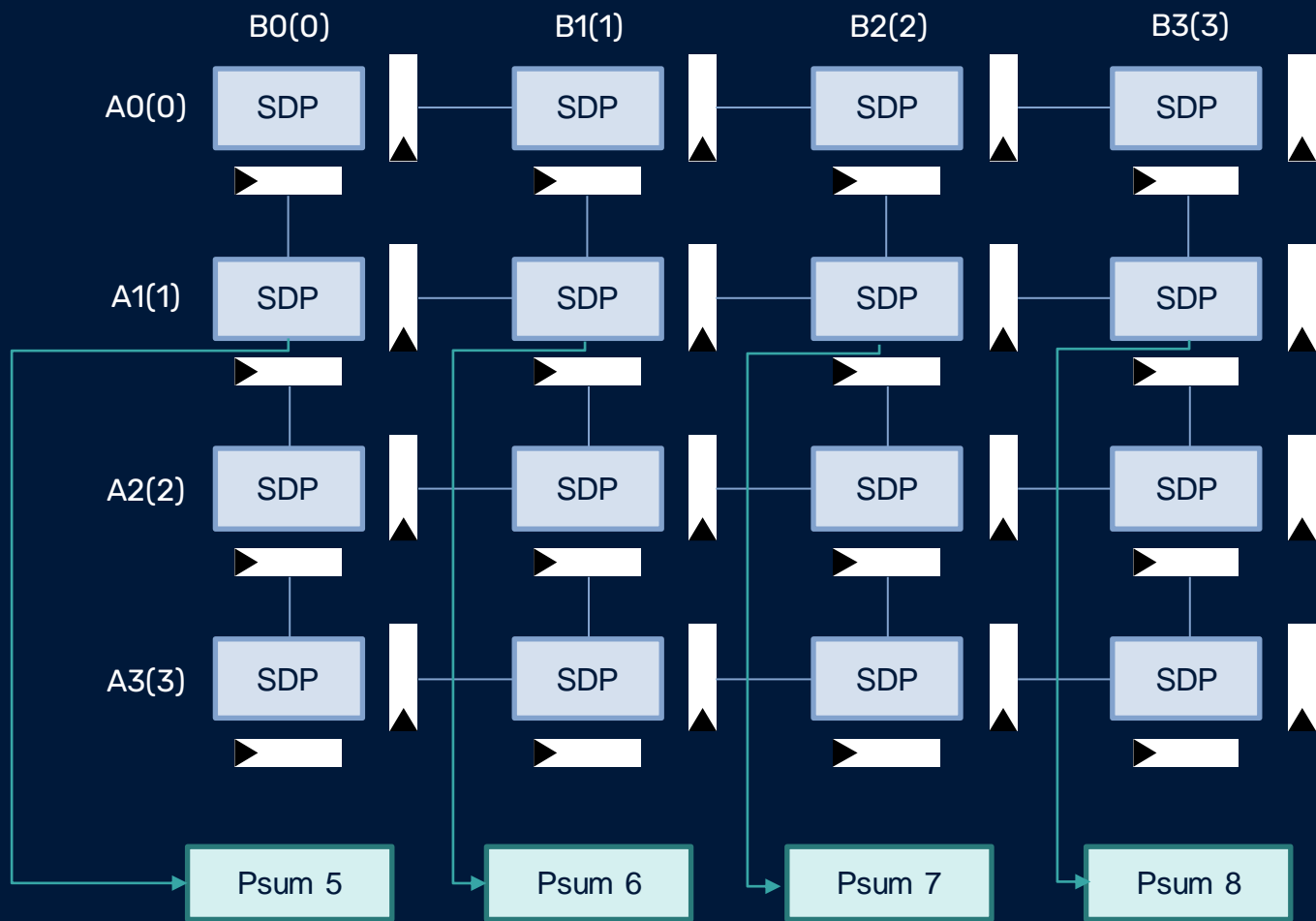# PE Array Architecture

A for 3*3 Ifmap Sliding Window
B for 3*3 Weight
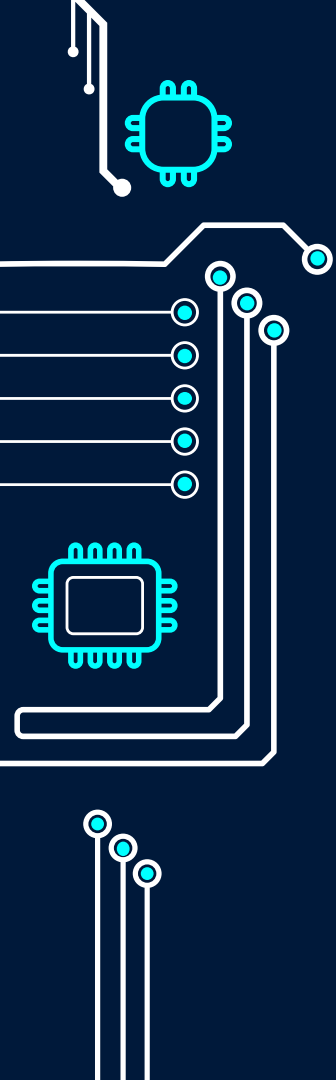(num) num :which clock to send

Channel 0

|  | B0(0) | B1(1) | B2(2) | B3(3) |
|---|---|---|---|---|
| A4(0) | SDP | SDP | SDP | SDP |
| A5(1) | SDP | SDP | SDP | SDP |
| A6(2) | SDP | SDP | SDP | SDP |
| A7(3) | SDP | SDP | SDP | SDP |

# 03

## Our Accelerator Architecture and SPEC

# Architecture

Off Chip DRAM

ifmap SRAM

weight SRAM

Controller

Weight buffer

ifmap buffer

PE Array

Accumulator buffer

Post ReLU

Pooling Engine

# SPEC

## SRAM

| Type | Size | Why we use this size? |
|------|------|----------------------|
| Ifmap SRAM | 16KB | 512(input channel) * 3(filter row) <br>• 10(6+4) * 1Byte = 15KB |
| Weight SRAM | 16KB | 512(input channel) * 8(output channel) * 3(NNZ) * 10(int8+2 bit address) = 15KB |

# SPEC

## On-Chip Buffer

| Type | Size | Why we use this size? |
|------|------|----------------------|
| Ifmap Buffer | 576B | 32(input channel) * 3(filter row) * 6(ifmap column) * 1Byte = 576B |
| Weight Buffer | 480B | 32(input channel) * 3(filter NNZ) * 4(PE a row) * 10bit (int8 + 2 bit address) = 480B |
| Accumulator Buffer | 64B | 32 Bits * 16 (PE) = 64B |

# SPEC

## Bandwidth

| Type | Size | Why we use this size? |
|------|------|----------------------|
| PE to Ifmap Buffer | 28B | 7(element)*4(PE a row)*1Byte = 28B |
| PE to Weight Buffer | 15B | 3(NNZ in a filter) * 4(PE a col) * 10 bits = 15B |
| PE to Accumulator Buffer | 64B | 4Byte * 16 = 64B |

# 04

## Analysis and Question

# Advantages

- Increased Efficiency

- Enhanced Area and Power Efficiency
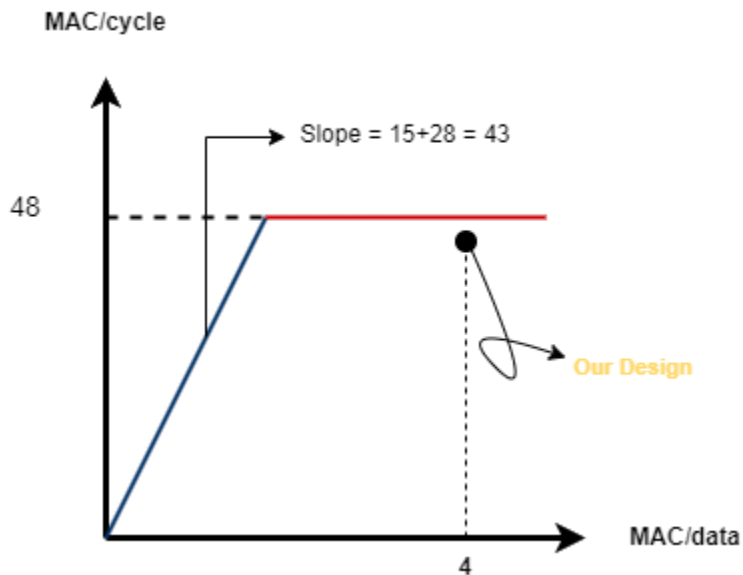
- Efficient Data Flow

# Roofline Model

# Problem Disscusion

Q1: Is the Accumulator on-chip bandwidth (64B) too large?

Q2: Is the number of times data is fetched from DRAM too many?

    (Up to 64 times per layer)

# RESOURCES

**Course**

- AI on Chip NCKU by Professor Chia-Chi Tsai

**Paper Reference**

- Liu, Z. G., Whatmough, P. N., & Mattina, M. (2020). Systolic tensor array: An efficient structured-sparse GEMM accelerator for mobile CNN inference. IEEE Computer Architecture Letters, 19(1), 34-37.

- Niu, W., Ma, X., Lin, S., Wang, S., Qian, X., Lin, X., ... & Ren, B. (2020, March). Patdnn: Achieving real-time dnn execution on mobile devices with pattern-based weight pruning. In Proceedings of the Twenty-Fifth International Conference on Architectural Support for Programming Languages and Operating Systems (pp. 907-922).

THANKS!