



# Lecture 09: CNN Dataflow & Hardware Accelerators

# Recap

- Federated Learning
- Distributed DNN Training
- Distributed DNN Inference
- Speculative Decoding

# Notes

- Project Proposal due Thursday
  - Signup your team
- Mid-semester Course Feedback
- Midterm Exam Review

# Topics

- **Hardware accelerator: Overview**
- Convolutional operation conversion
- Hardware architecture of CNN accelerator
- Systolic array
- Popular accelerator design
  - Eyeriss
  - Diannao
  - Cnvlutin
  - EIE

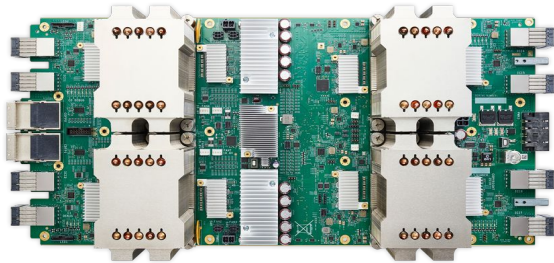
# Hardware Support for DNN

- GPU is better than CPU in terms of throughput for both Neural Network training and inference.
  - GPU leverages the highly parallelized architecture of its computing units to handle computational intensive operations.
  - GPU has 10x-20x higher throughput than CPU.
- However, GPU:
  - General purpose.
  - Power consumption and latency is high.
  - Does not support sophisticated pruning and quantization algorithms.



# Hardware Support for DNN

- ASIC-based implementations have been recently explored to accelerate the DNN inference.
  - Google's TPU, Apple's Neural Engine, Cerebras AI chip, ...
- FPGA-based accelerators for DNN inference have been recently developed.
  - Has good programmability and flexibility
  - Short development cycles
  - Can be used as a benchmark before implementing on ASIC

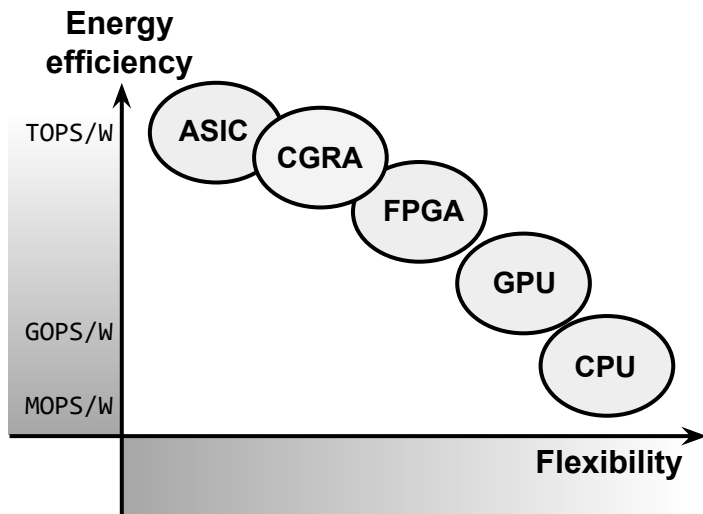


Tensor Processing Unit (Google)



Alveo Accelerator Card (Xilinx)

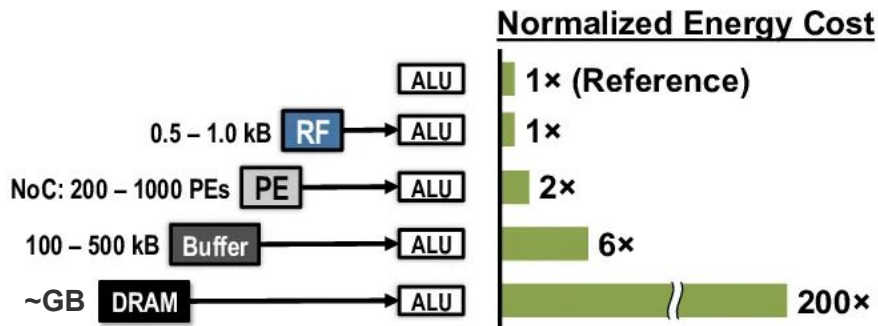
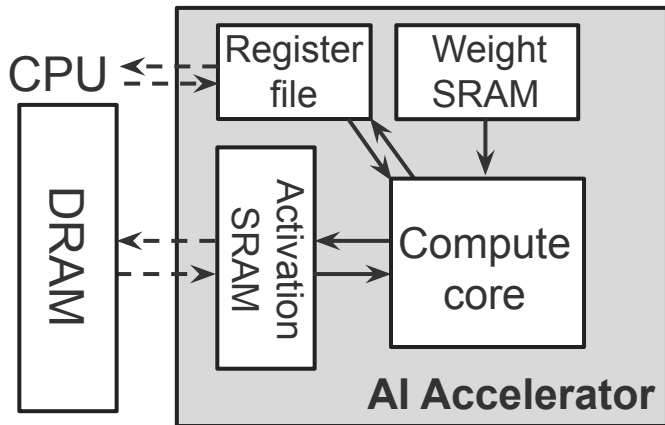
# Flexibility & Performance



- ASIC offers the highest energy efficiency but is only suitable for specific applications.
- The CPU is a general-purpose processor but has the lowest energy efficiency.

# AI Accelerator

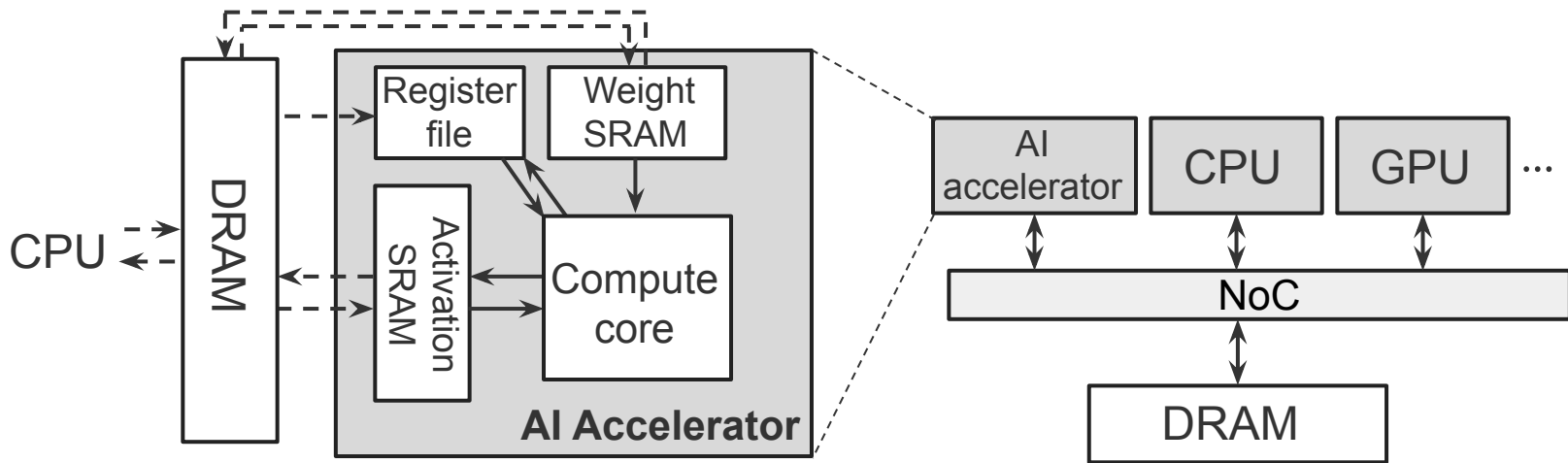
- Making any chip is a costly, difficult and lengthy process typically done by teams of 10 to 1000's of people depending on the size and complexity of the chip.



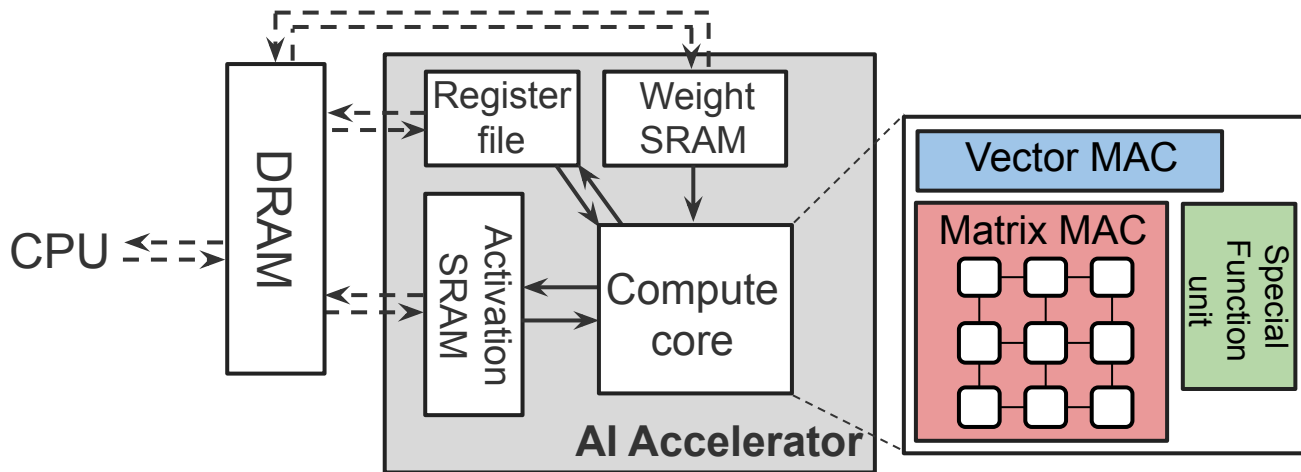


# AI Accelerator

- The AI accelerator can execute part of the machine code that is related to the AI workload.



# AI Accelerator



- The compute core consists of Multiply and accumulator (MAC) engine for 2D matrix multiplication.
- It also contains vector multiplier MAC as well as special function unit.

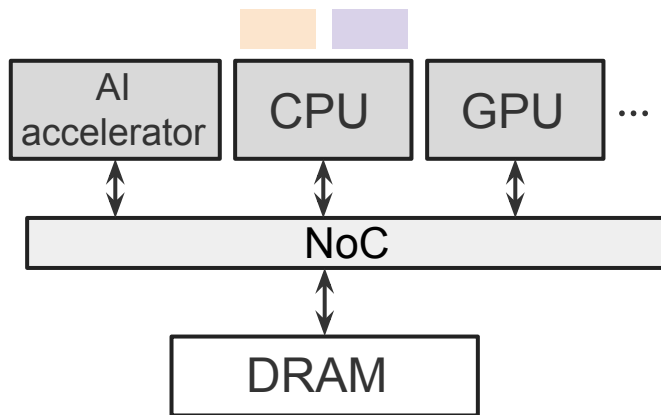
# AI Accelerator



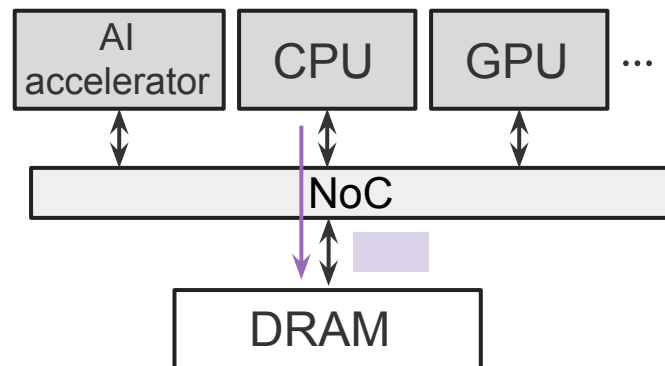
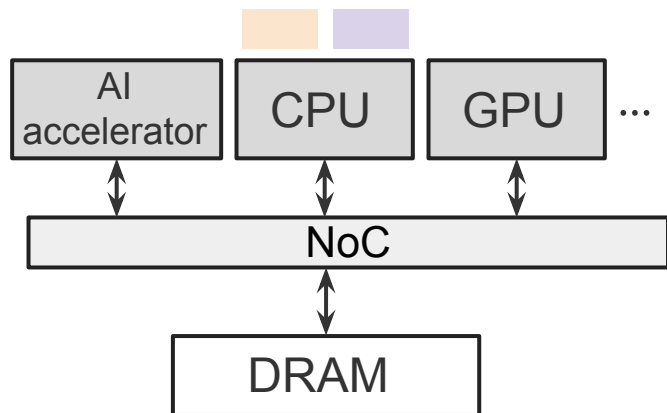
```
import torch
import torch.nn as nn
input = torch.randn(1, 1, 5, 5)
conv = nn.Conv2d(1, 1, kernel_size=3)
output = conv(input)
print(output)
```

Compiler

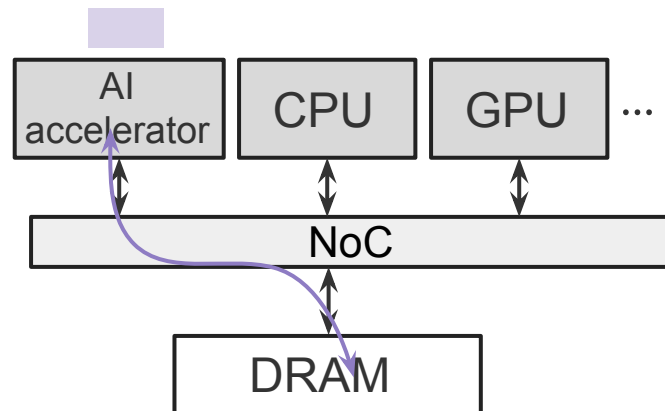
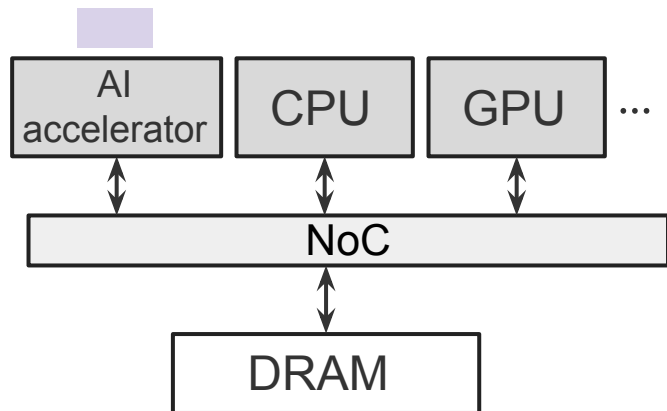
```
01001000 10001001 11011000
01001000 10000011 11000000 00001010
01001000 10000011 11101011 00000101
01001000 00111001 11011000
01110100 00000101
01001000 10001001 11000001
01001000 01101001 11001001 00000010
01001000 10001001 00001111
11101011 00000011
```



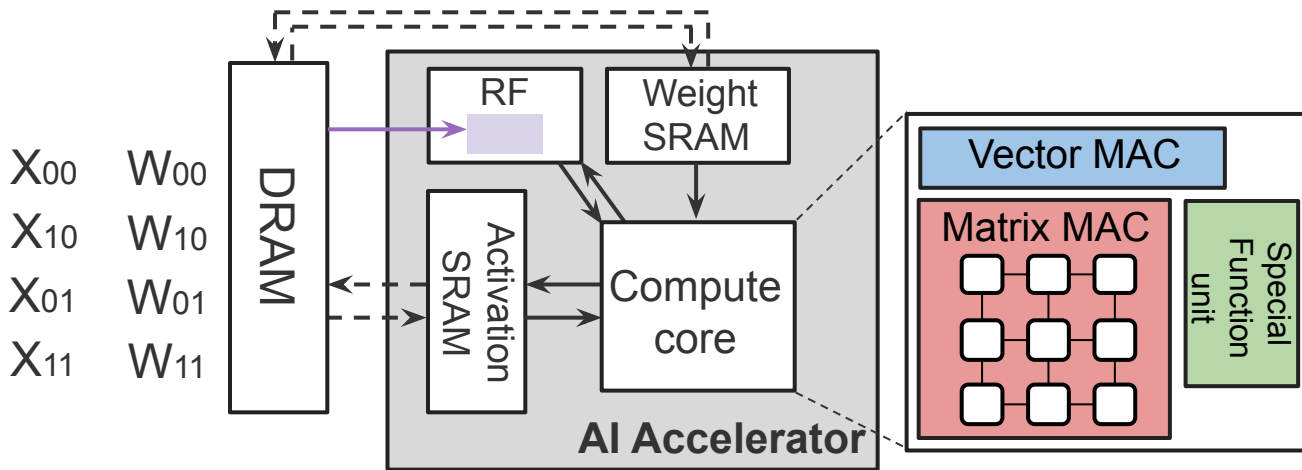
# AI Accelerator



# AI Accelerator

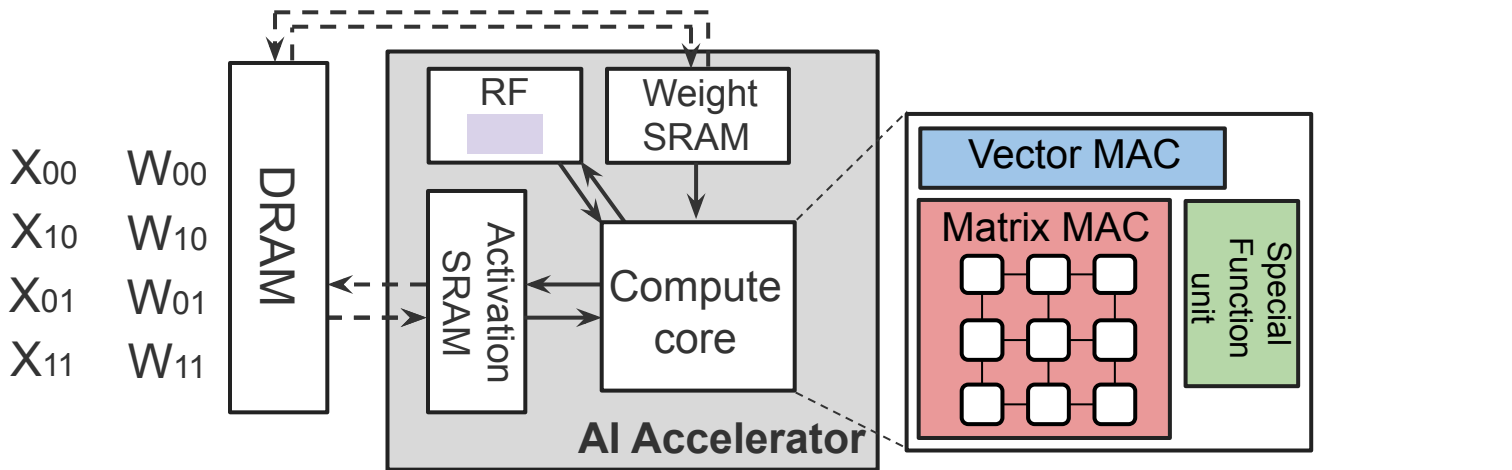


# AI Accelerator



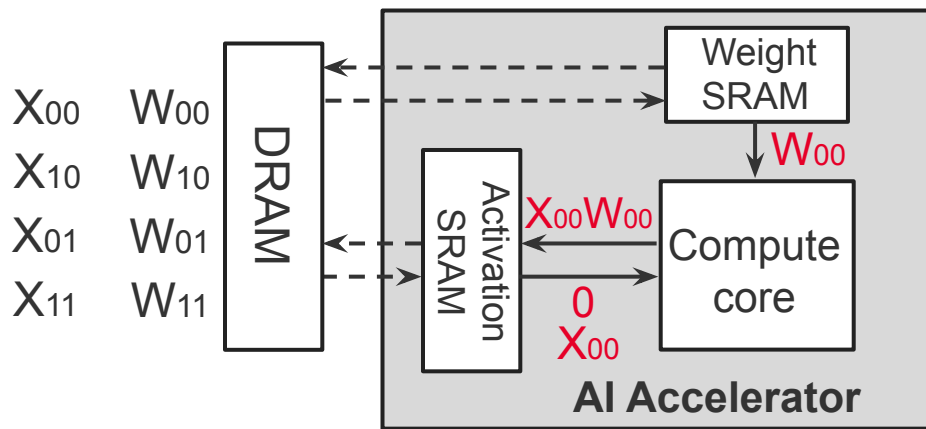
$$\begin{bmatrix} W_{00} & W_{01} \\ W_{10} & W_{11} \end{bmatrix} \times \begin{bmatrix} X_{00} & X_{01} \\ X_{10} & X_{11} \end{bmatrix} = \begin{bmatrix} W_{00}X_{00}+W_{01}X_{10} & W_{00}X_{01}+W_{01}X_{11} \\ W_{10}X_{00}+W_{11}X_{10} & W_{10}X_{01}+W_{11}X_{11} \end{bmatrix} = \begin{bmatrix} Y_{00} & Y_{01} \\ Y_{10} & Y_{11} \end{bmatrix}$$

# AI Accelerator



$$\begin{bmatrix} W_{00} & W_{01} \\ W_{10} & W_{11} \end{bmatrix} \times \begin{bmatrix} X_{00} & X_{01} \\ X_{10} & X_{11} \end{bmatrix} = \begin{bmatrix} W_{00}X_{00}+W_{01}X_{10} & W_{00}X_{01}+W_{01}X_{11} \\ W_{10}X_{00}+W_{11}X_{10} & W_{10}X_{01}+W_{11}X_{11} \end{bmatrix} = \begin{bmatrix} Y_{00} & Y_{01} \\ Y_{10} & Y_{11} \end{bmatrix}$$

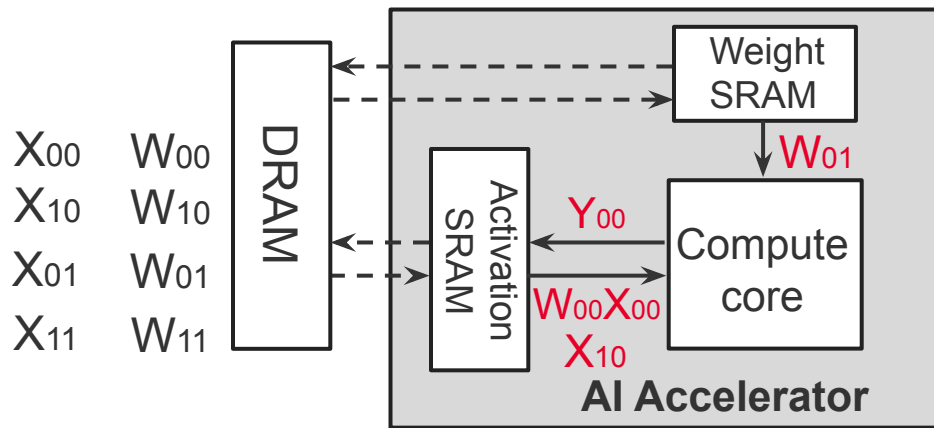
# AI Accelerator



$$\begin{bmatrix} W_{00} & W_{01} \\ W_{10} & W_{11} \end{bmatrix} \times \begin{bmatrix} X_{00} & X_{01} \\ X_{10} & X_{11} \end{bmatrix} = \begin{bmatrix} W_{00}X_{00} + W_{01}X_{10} & W_{00}X_{01} + W_{01}X_{11} \\ W_{10}X_{00} + W_{11}X_{10} & W_{10}X_{01} + W_{11}X_{11} \end{bmatrix} = \begin{bmatrix} Y_{00} & Y_{01} \\ Y_{10} & Y_{11} \end{bmatrix}$$

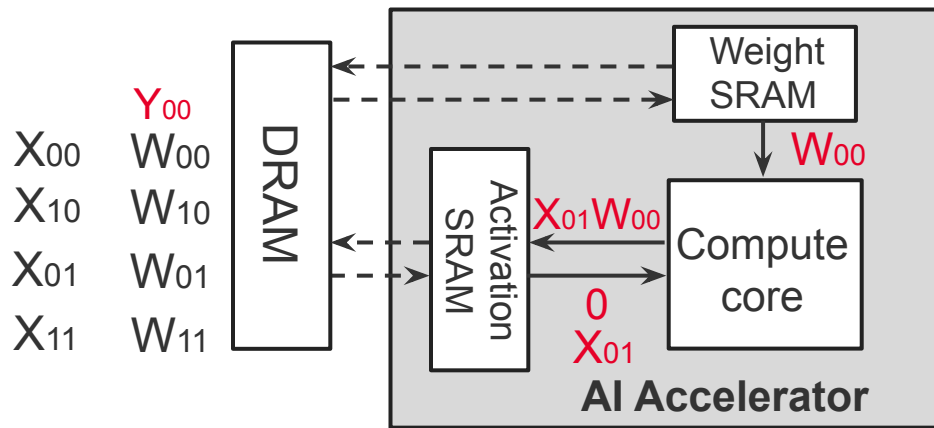


# AI Accelerator



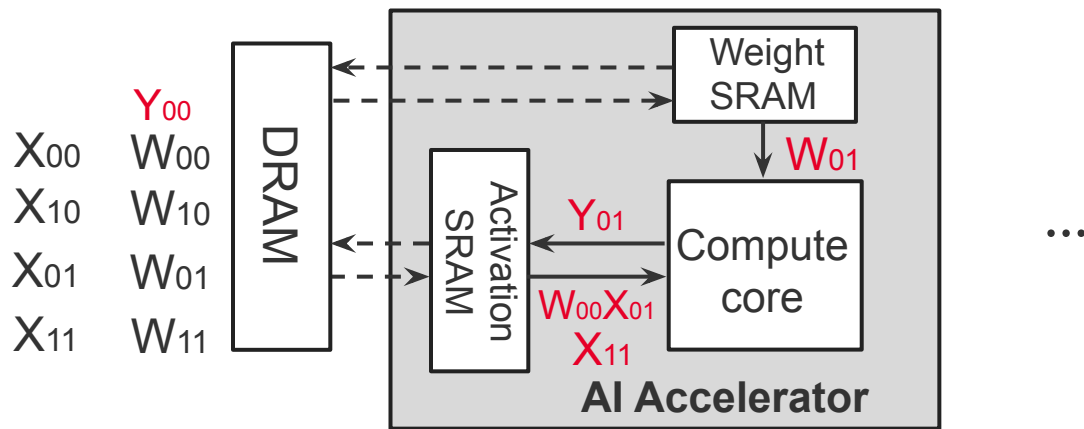
$$\begin{bmatrix} W_{00} & W_{01} \\ W_{10} & W_{11} \end{bmatrix} \times \begin{bmatrix} X_{00} & X_{01} \\ X_{10} & X_{11} \end{bmatrix} = \begin{bmatrix} W_{00}X_{00} + W_{01}X_{10} & W_{00}X_{01} + W_{01}X_{11} \\ W_{10}X_{00} + W_{11}X_{10} & W_{10}X_{01} + W_{11}X_{11} \end{bmatrix} = \begin{bmatrix} Y_{00} & Y_{01} \\ Y_{10} & Y_{11} \end{bmatrix}$$

# AI Accelerator



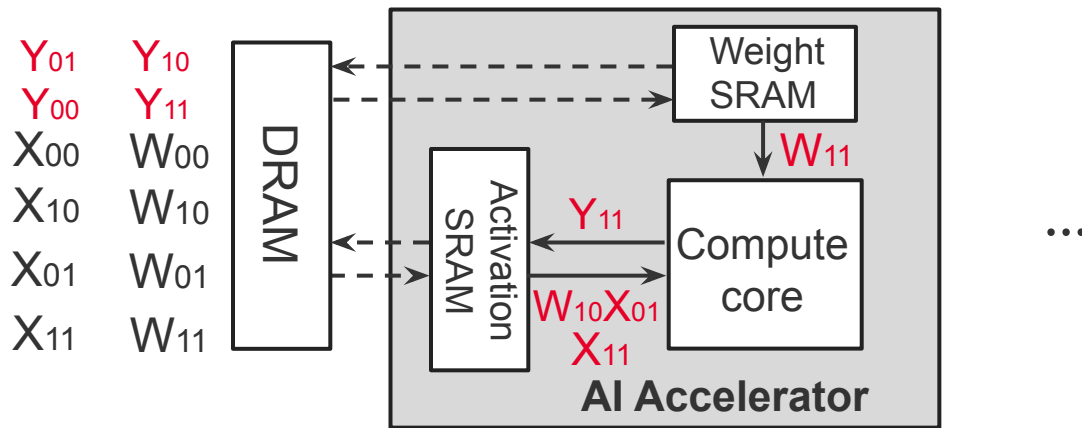
$$\begin{bmatrix} W_{00} & W_{01} \\ W_{10} & W_{11} \end{bmatrix} \times \begin{bmatrix} X_{00} & X_{01} \\ X_{10} & X_{11} \end{bmatrix} = \begin{bmatrix} W_{00}X_{00}+W_{01}X_{10} & \color{red}W_{00}X_{01}+W_{01}X_{11} \\ W_{10}X_{00}+W_{11}X_{10} & W_{10}X_{01}+W_{11}X_{11} \end{bmatrix} = \begin{bmatrix} Y_{00} & Y_{01} \\ Y_{10} & Y_{11} \end{bmatrix}$$

# AI Accelerator



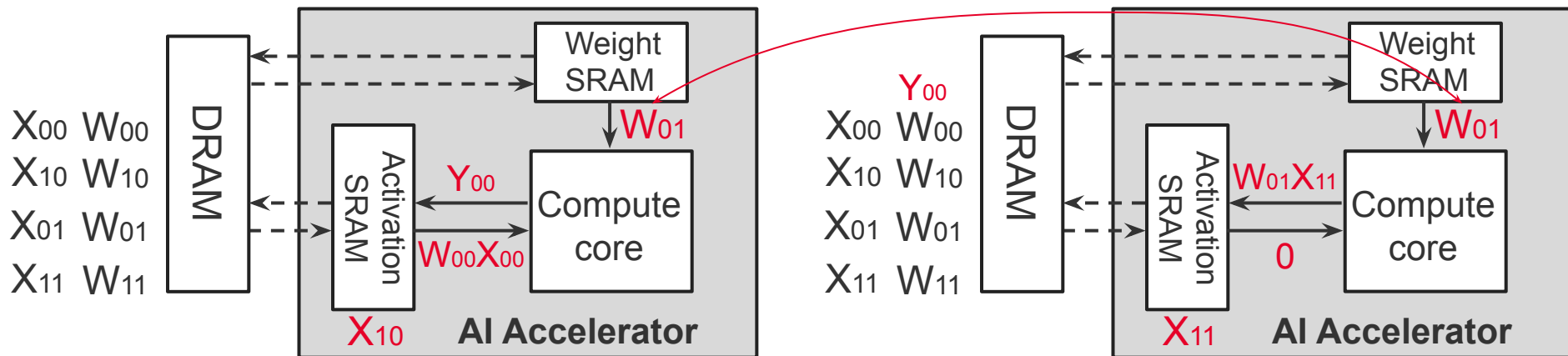
$$\begin{bmatrix} W_{00} & W_{01} \\ W_{10} & W_{11} \end{bmatrix} \times \begin{bmatrix} X_{00} & X_{01} \\ X_{10} & X_{11} \end{bmatrix} = \begin{bmatrix} W_{00}X_{00} + W_{01}X_{10} & W_{00}X_{01} + W_{01}X_{11} \\ W_{10}X_{00} + W_{11}X_{10} & W_{10}X_{01} + W_{11}X_{11} \end{bmatrix} = \begin{bmatrix} Y_{00} & Y_{01} \\ Y_{10} & Y_{11} \end{bmatrix}$$

# AI Accelerator



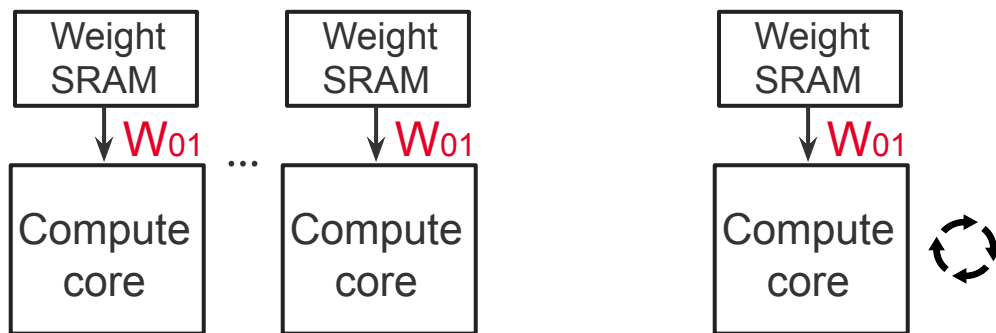
$$\begin{bmatrix} W_{00} & W_{01} \\ W_{10} & W_{11} \end{bmatrix} \times \begin{bmatrix} X_{00} & X_{01} \\ X_{10} & X_{11} \end{bmatrix} = \begin{bmatrix} W_{00}X_{00}+W_{01}X_{10} & W_{00}X_{01}+W_{01}X_{11} \\ W_{10}X_{00}+W_{11}X_{10} & \color{red}W_{10}X_{01}+\color{red}W_{11}X_{11} \end{bmatrix} = \begin{bmatrix} Y_{00} & Y_{01} \\ Y_{10} & \color{red}Y_{11} \end{bmatrix}$$

# Memory Access Reduction



- The computation and memory access pattern can be changed to minimize the computational cost without impacting the results.

# Memory Access Reduction

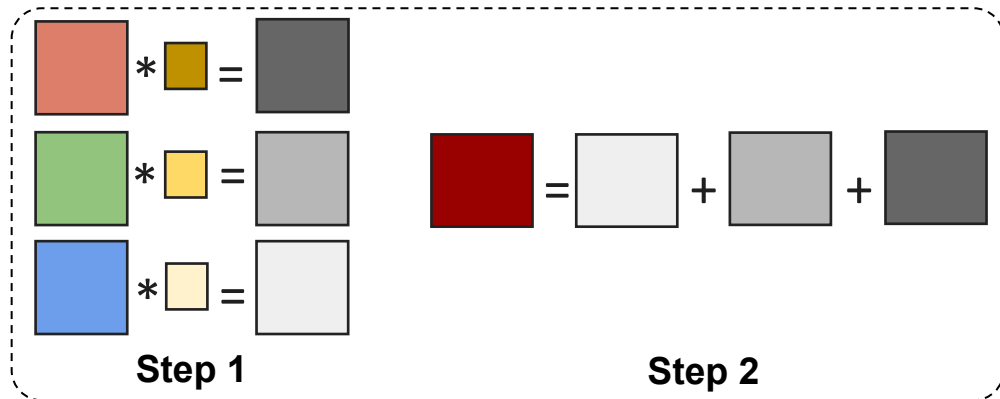
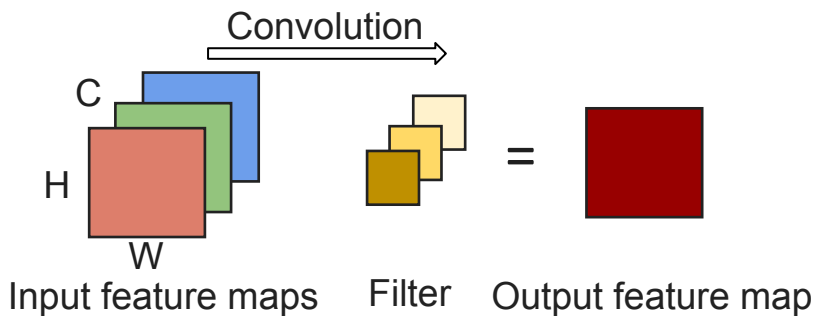


- It is preferable to minimize memory access by maximizing the reuse of loaded data.

# Topics

- Hardware accelerator: Overview
- Convolutional operation conversion
- Hardware architecture of CNN accelerator
- Systolic array
- Popular accelerator design
  - Eyeriss
  - Diannao
  - Cnvlutin
  - EIE

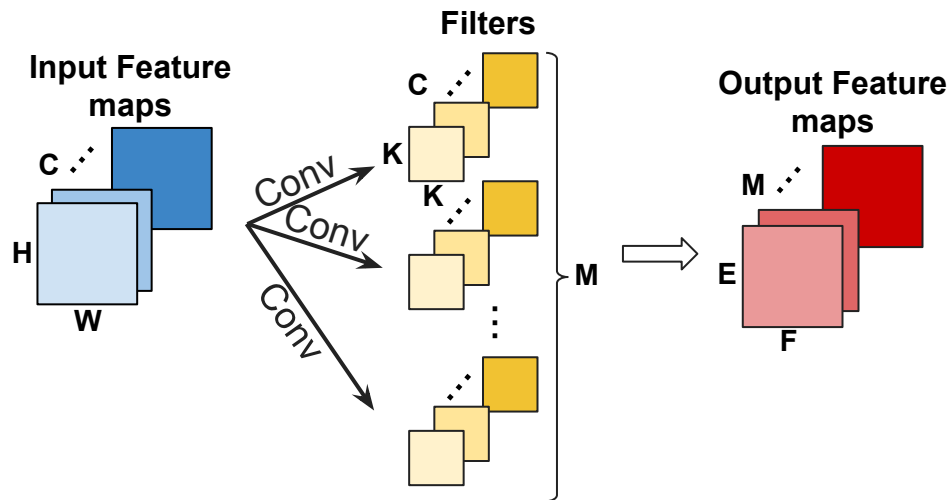
# Convolutional Layers



- Core building block of a CNN, it is also the most computational intensive layer.



# Convolution



- Number of MACs:  $M \times K \times K \times C \times E \times F$
- Storage cost:  
 $32 \times (M \times C \times K \times K + C \times H \times W + M \times E \times F)$

$C$ : number of input channels

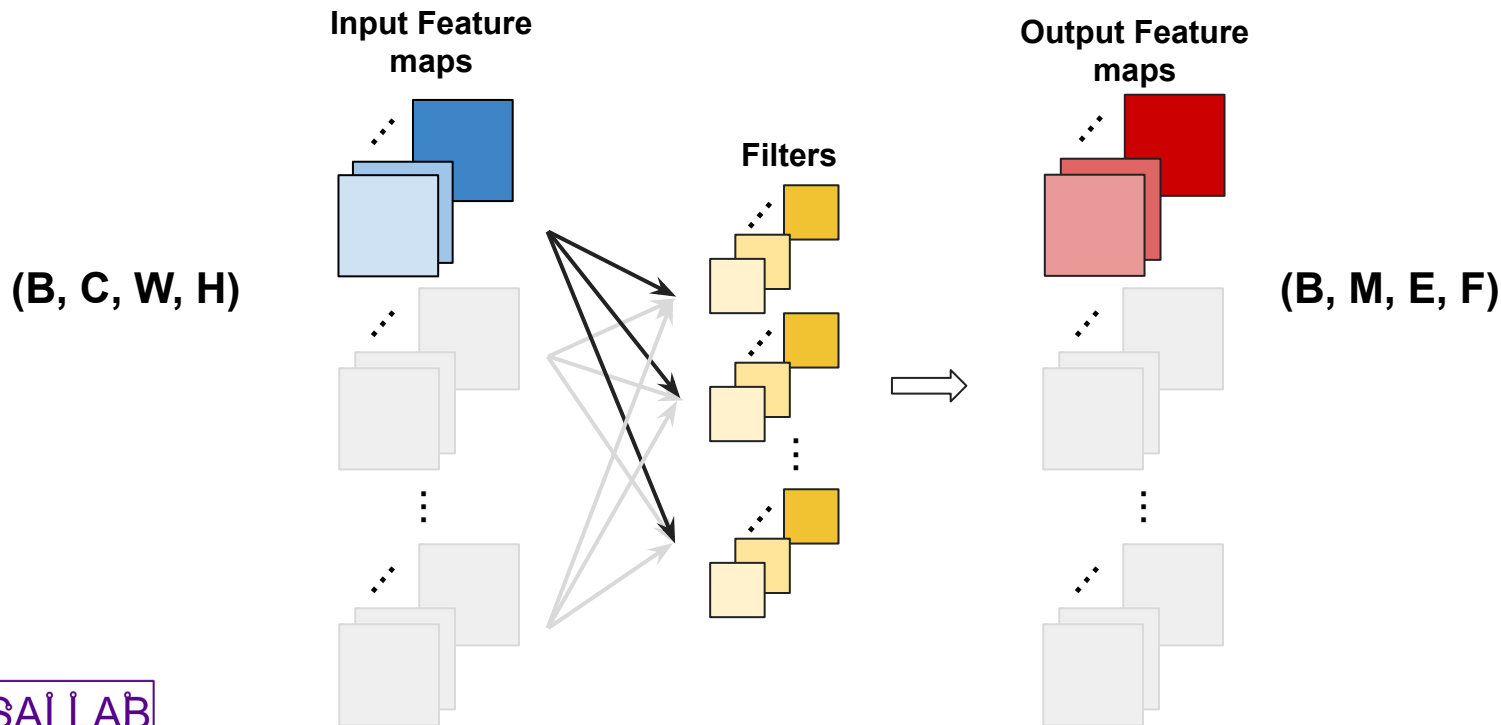
$H, W$ : size of the input feature maps

$M$ : number of weight filters

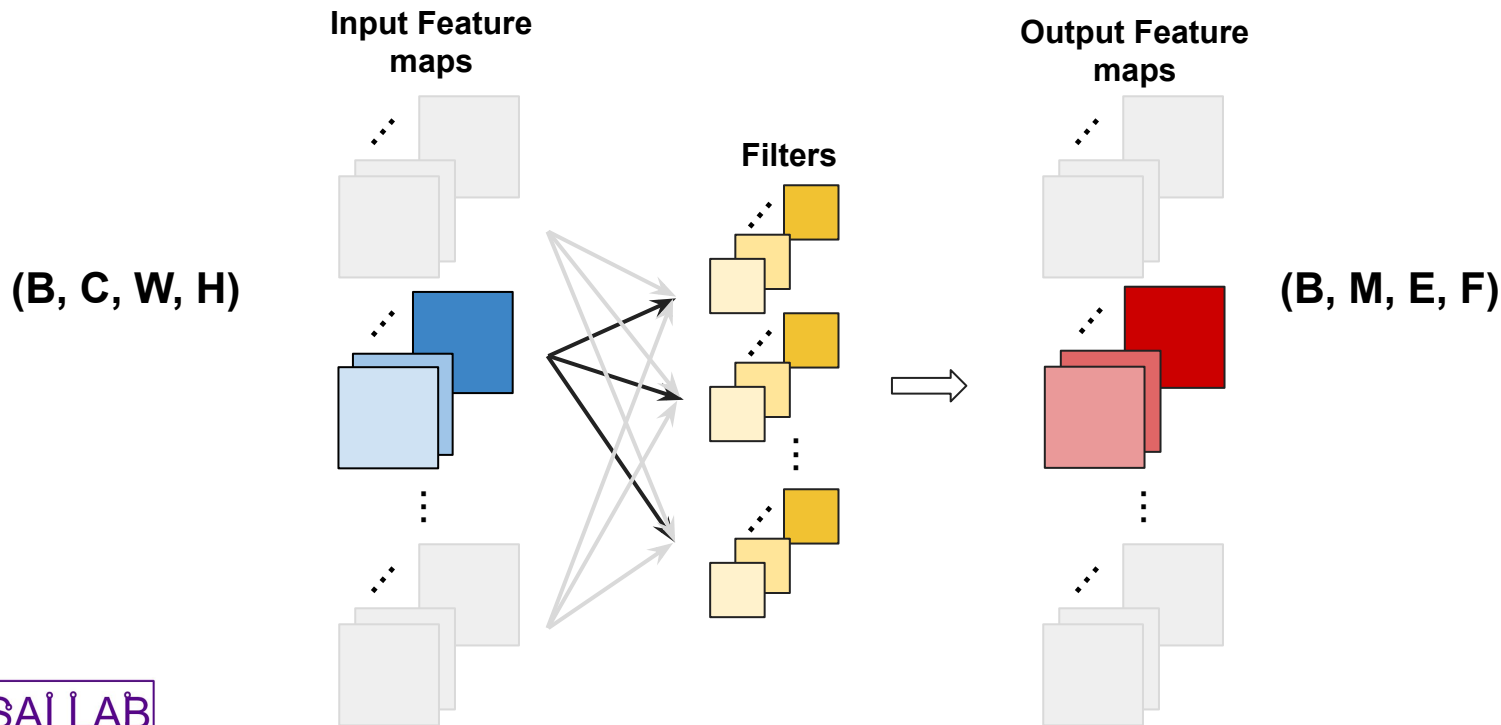
$K$ : weight kernel size

$E, F$ : size of the output feature maps

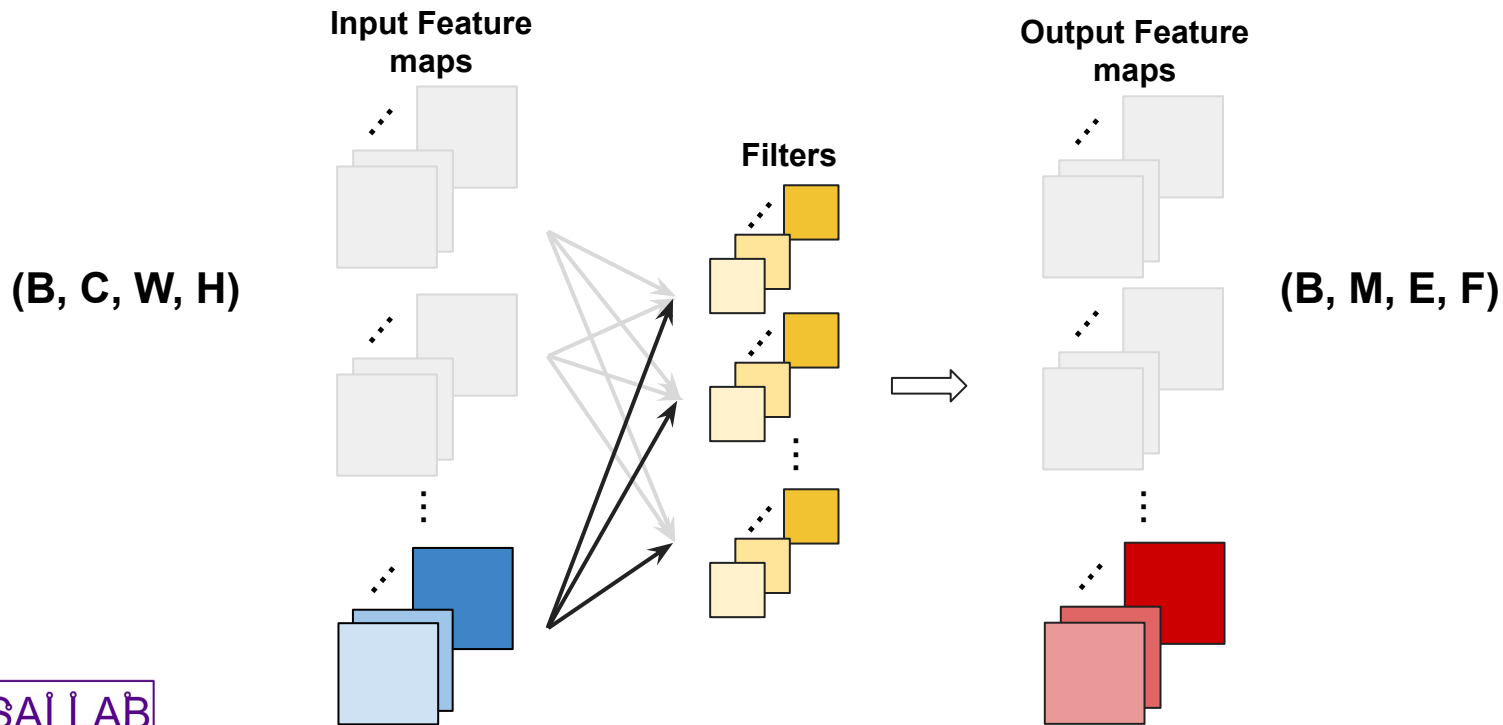
# Convolution



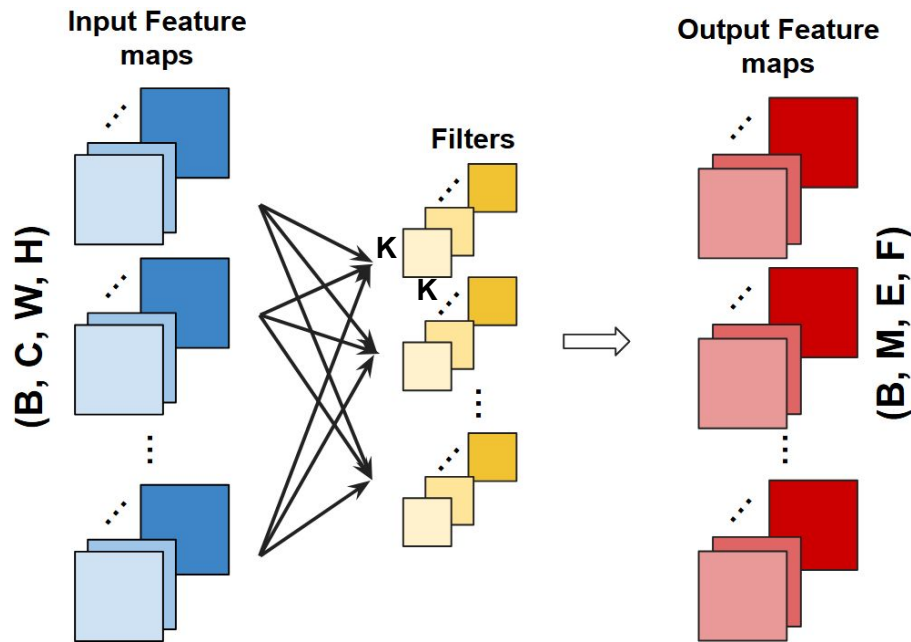
# Convolution



# Convolution



# Computational Cost: Standard Convolution



- Number of MACs:  $B \times M \times K \times K \times C \times E \times F$
- Storage cost:  
 $32 \times (M \times C \times K \times K + B \times C \times H \times W + B \times M \times E \times F)$

B: batch size

C: number of input channels

H,W: size of the input feature maps

M: number of weight filters

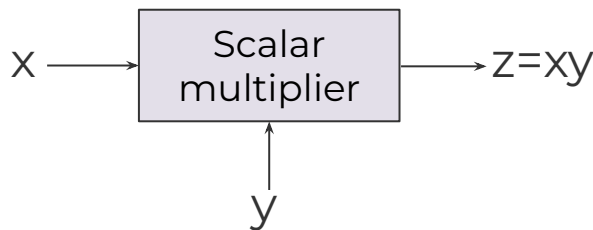
K: weight kernel size

E,F: size of the output feature maps

- We need to iterate over seven dimensions:
  - B, M, C, E, F, K(kernel width), K (kernel height)

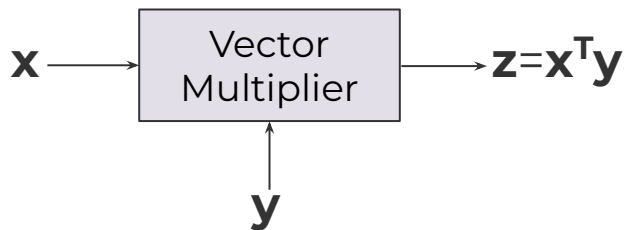
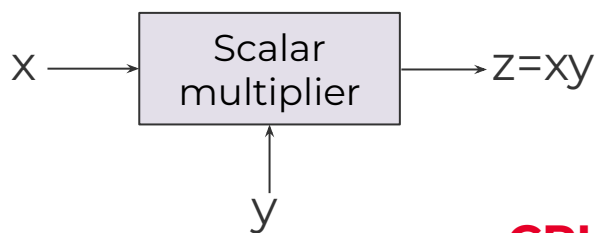
# Computational Dataflow for CNN

```
for b = 1 to B
  for m = 1 to M
    for c = 1 to C
      for w = 1 to E
        for h = 1 to F
          for k1 = 1 to K
            for k2 = 1 to K
              out[b][m][e][f] += in[b][c][e+k1-(K+1)/2][f+k2-(K+1)/2] * filter[m][c][k1][k2];
```

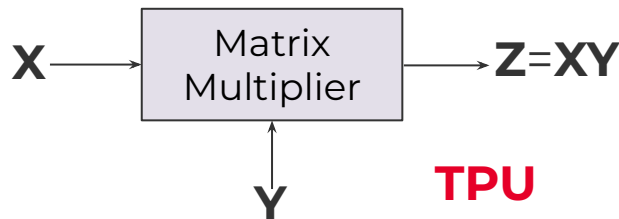


- This simple loop nest can be transformed in numerous ways to capture different reuse patterns of the activations and weights and to map the computation to a hardware accelerator implementation.
- A CNN's dataflow defines how the loops are ordered, partitioned, and parallelized
- We can use the scalar machine to compute the results of CNN using this for loop

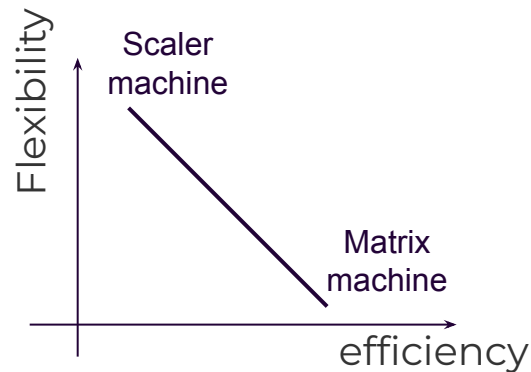
# Computational Dataflow for CNN



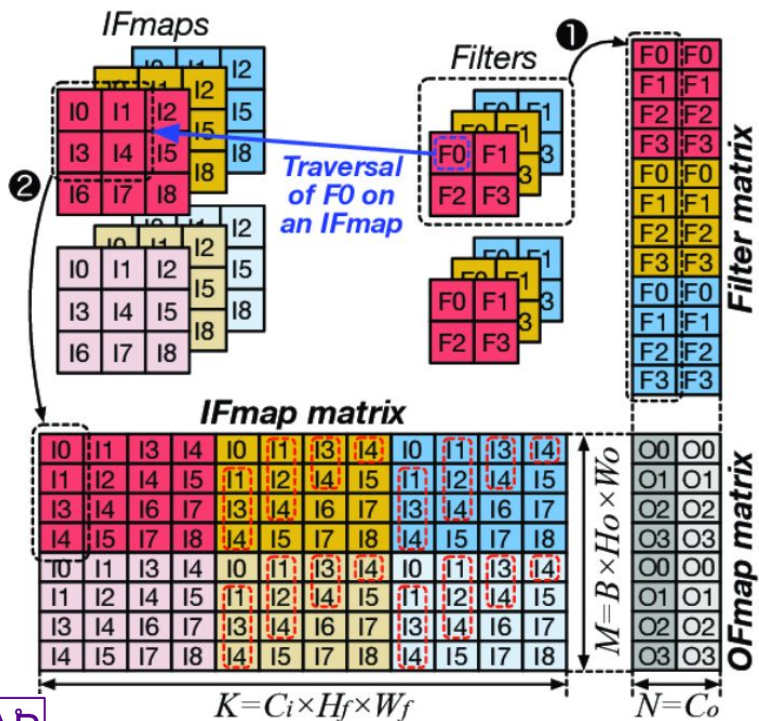
**CPU & GPU**



**TPU**



# How to Convert to Matrix Multiplication?



- A standard Convolutional operation can be converted to 2D matrix multiplication using **Im2Col** operations.



# How to Convert to Matrix Multiplication?

Input feature map

a <sub>00</sub>	a <sub>01</sub>	a <sub>02</sub>	a <sub>03</sub>
a <sub>10</sub>	a <sub>11</sub>	a <sub>12</sub>	a <sub>13</sub>
a <sub>20</sub>	a <sub>21</sub>	a <sub>22</sub>	a <sub>23</sub>
a <sub>30</sub>	a <sub>31</sub>	a <sub>32</sub>	a <sub>33</sub>

Weight filter

W <sub>00</sub>	W <sub>01</sub>	W <sub>02</sub>
W <sub>10</sub>	W <sub>11</sub>	W <sub>12</sub>
W <sub>20</sub>	W <sub>21</sub>	W <sub>22</sub>

\*



Input matrix

a <sub>00</sub>	a <sub>01</sub>	a <sub>02</sub>	a <sub>10</sub>	a <sub>11</sub>	a <sub>12</sub>	a <sub>20</sub>	a <sub>21</sub>	a <sub>22</sub>
a <sub>01</sub>	a <sub>02</sub>	a <sub>03</sub>	a <sub>11</sub>	a <sub>12</sub>	a <sub>13</sub>	a <sub>21</sub>	a <sub>22</sub>	a <sub>23</sub>
a <sub>10</sub>	a <sub>11</sub>	a <sub>12</sub>	a <sub>20</sub>	a <sub>21</sub>	a <sub>22</sub>	a <sub>30</sub>	a <sub>31</sub>	a <sub>32</sub>
a <sub>11</sub>	a <sub>12</sub>	a <sub>13</sub>	a <sub>21</sub>	a <sub>22</sub>	a <sub>23</sub>	a <sub>31</sub>	a <sub>32</sub>	a <sub>33</sub>

Weight matrix

W <sub>00</sub>
W <sub>01</sub>
W <sub>02</sub>
W <sub>10</sub>
W <sub>11</sub>
W <sub>12</sub>
W <sub>20</sub>
W <sub>21</sub>
W <sub>22</sub>

×

# How to Convert to Matrix Multiplication?

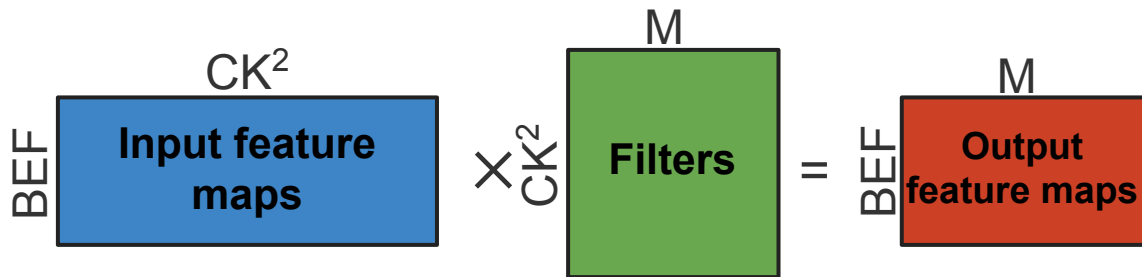
Convolution:

Filter		Input Fmap		Output Fmap																	
<table border="1"><tr><td>1</td><td>2</td></tr><tr><td>3</td><td>4</td></tr></table>	1	2	3	4	*	<table border="1"><tr><td>1</td><td>2</td><td>3</td></tr><tr><td>4</td><td>5</td><td>6</td></tr><tr><td>7</td><td>8</td><td>9</td></tr></table>	1	2	3	4	5	6	7	8	9	=	<table border="1"><tr><td>1</td><td>2</td></tr><tr><td>3</td><td>4</td></tr></table>	1	2	3	4
1	2																				
3	4																				
1	2	3																			
4	5	6																			
7	8	9																			
1	2																				
3	4																				



Matrix Mult:

<table border="1"><tr><td>1</td><td>2</td><td>3</td><td>4</td></tr></table>	1	2	3	4	x	<table border="1"><tr><td>1</td><td>2</td><td>4</td><td>5</td></tr><tr><td>2</td><td>3</td><td>5</td><td>6</td></tr><tr><td>4</td><td>5</td><td>7</td><td>8</td></tr><tr><td>5</td><td>6</td><td>8</td><td>9</td></tr></table>	1	2	4	5	2	3	5	6	4	5	7	8	5	6	8	9	=	<table border="1"><tr><td>1</td><td>2</td><td>3</td><td>4</td></tr></table>	1	2	3	4
1	2	3	4																									
1	2	4	5																									
2	3	5	6																									
4	5	7	8																									
5	6	8	9																									
1	2	3	4																									



# Tiling

- In order to handle matrix multiplication with large size, it is usually decomposed into tiles.

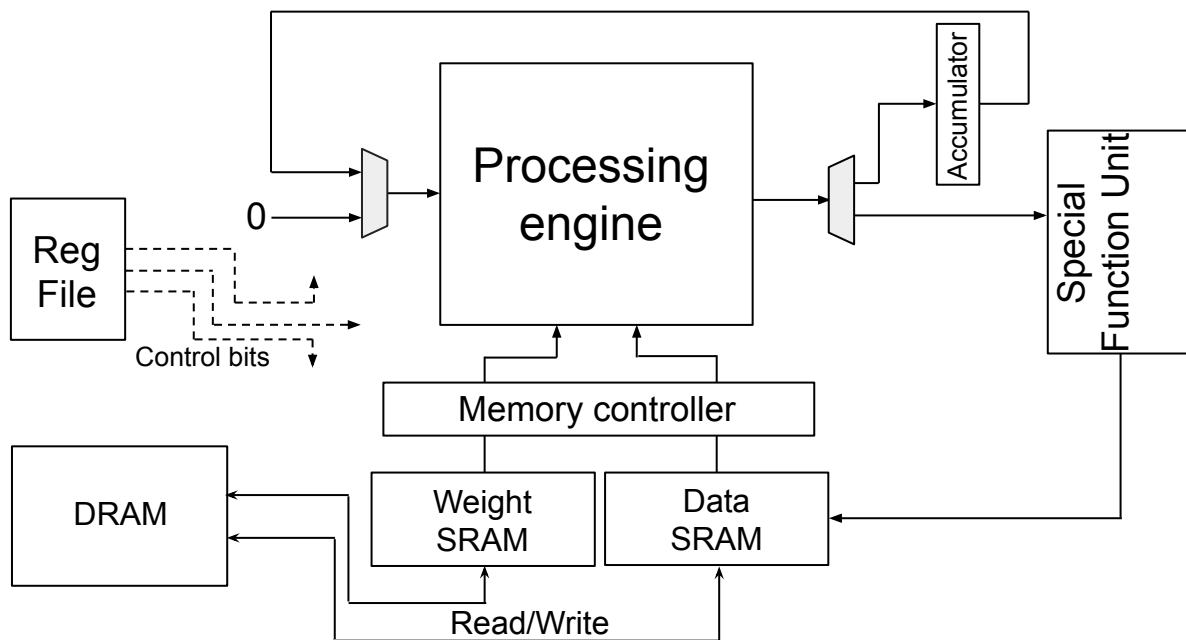
$$\begin{bmatrix} W_{00} & W_{01} \\ W_{10} & W_{11} \end{bmatrix} \times \begin{bmatrix} X_{00} & X_{01} \\ X_{10} & X_{11} \end{bmatrix} = \begin{bmatrix} W_{00}X_{00}+W_{01}X_{10} & W_{00}X_{01}+W_{01}X_{11} \\ W_{10}X_{00}+W_{11}X_{10} & W_{10}X_{01}+W_{11}X_{11} \end{bmatrix} = \begin{bmatrix} Y_{00} & Y_{01} \\ Y_{10} & Y_{11} \end{bmatrix}$$

- Each of  $W_{ij}$  and  $X_{ij}$  can be a sub-matrix.

# Topics

- Hardware accelerator: Overview
- Convolutional operation conversion
- **Hardware architecture of CNN accelerator**
- Systolic array
- Popular accelerator design
  - Eyeriss
  - Diannao
  - Cnvlutin
  - EIE

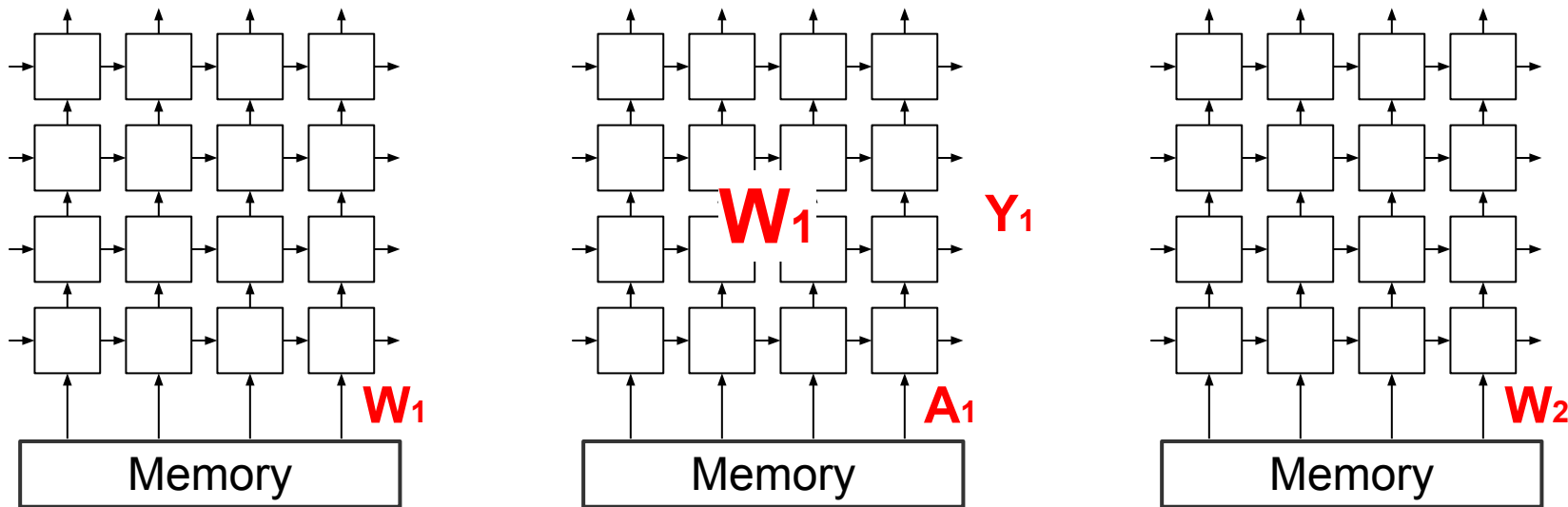
# Hardware Architectures for DNN Processing



Major building blocks:

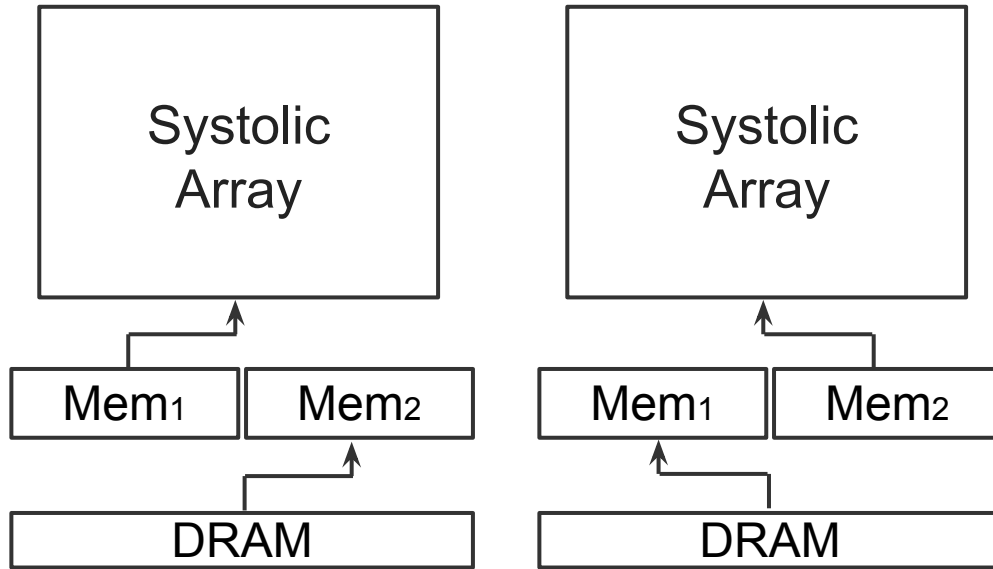
- Processing engine
- Accumulator
- Reg file
- Special function unit
- Memory subsystem
  - Weight SRAM
  - Data SRAM
  - DRAM

# Computing Paradigms



- Spatial architecture can achieve great reuse of the extracted content, leading to a reduced memory access cost.

# Double Buffering



- Double buffering in hardware design is a technique used to improve the efficiency and performance of data processing, especially in systems that require smooth and continuous data transfer.
- The idea is to overlap the data production and consumption processes to avoid delays.

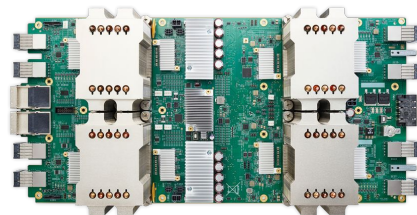
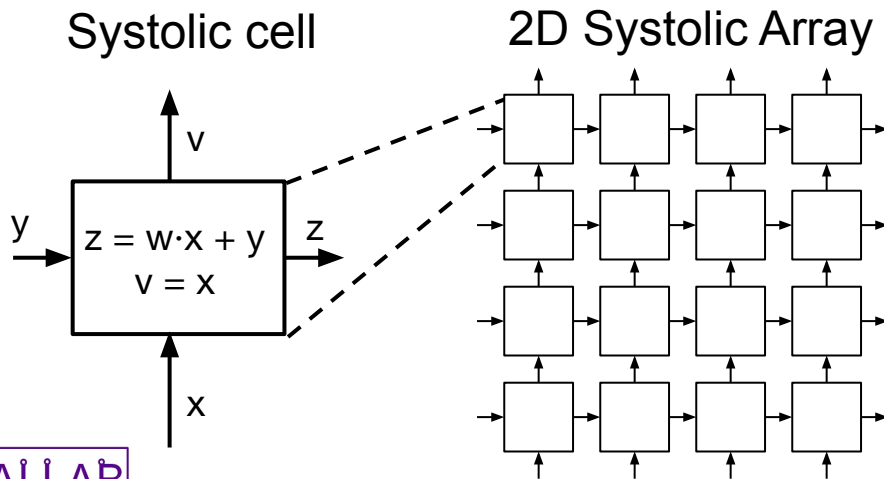
# Topics

- Hardware accelerator: Overview
- Convolutional operation conversion
- Hardware architecture of CNN accelerator
- **Systolic array**
- Popular accelerator design
  - Eyeriss
  - Diannao
  - Cnvlutin
  - EIE



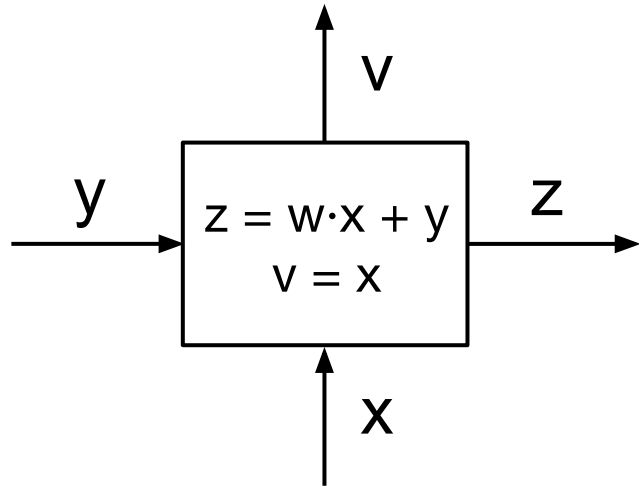
# Systolic Array (Weight Stationary Version)

- Kung and Leiserson, "Systolic Arrays for VLSI," 1978 and Kung, "Why systolic architectures?" 1982
- 2D grid of multiplier-accumulators (MACs) for matrix multiplication
- Used by Google TPU for deep learning (2017), etc



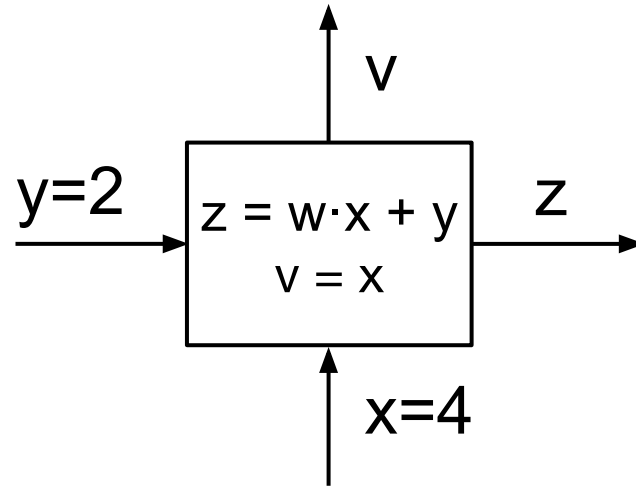
TPU (Google)

# Systolic Cell

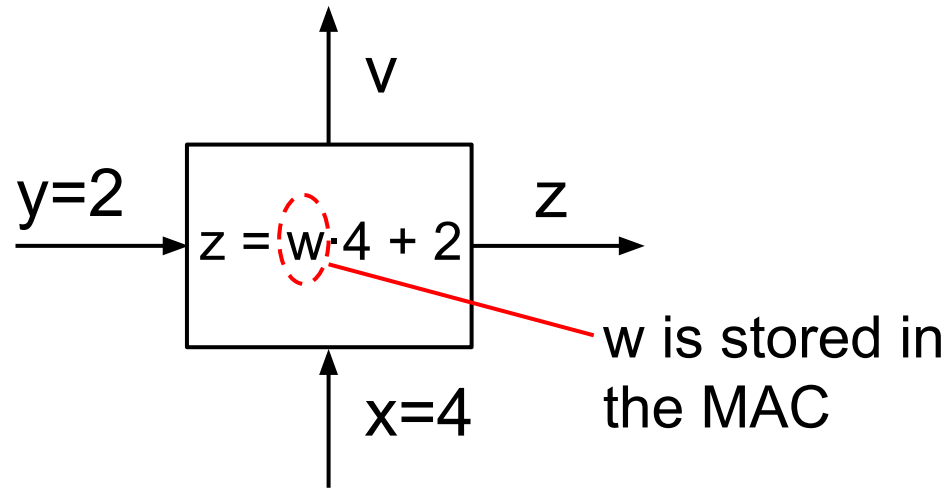


- Takes data ( $x$  and  $y$ ) as input
- $w$  stays in the systolic cell
- Performs a multiply-accumulate operation

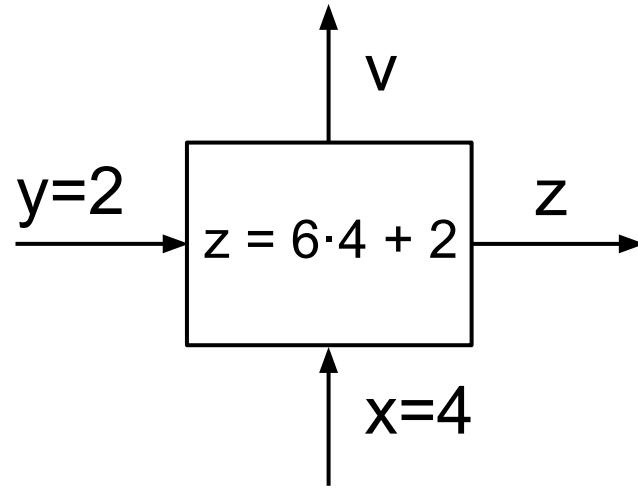
# Systolic Cell



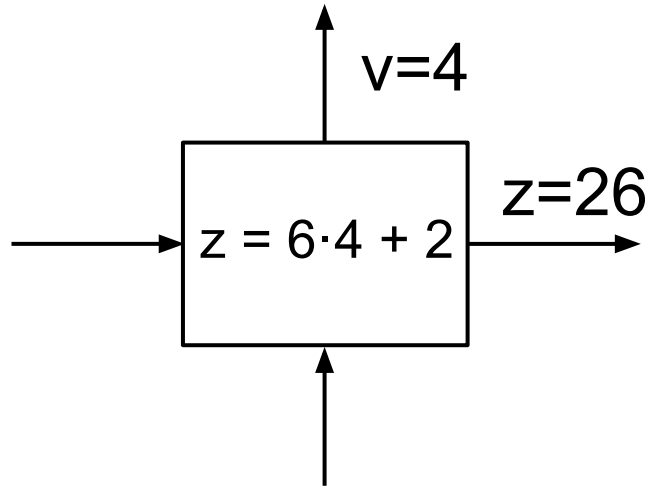
# Systolic Cell



# Systolic Cell



# Systolic Cell

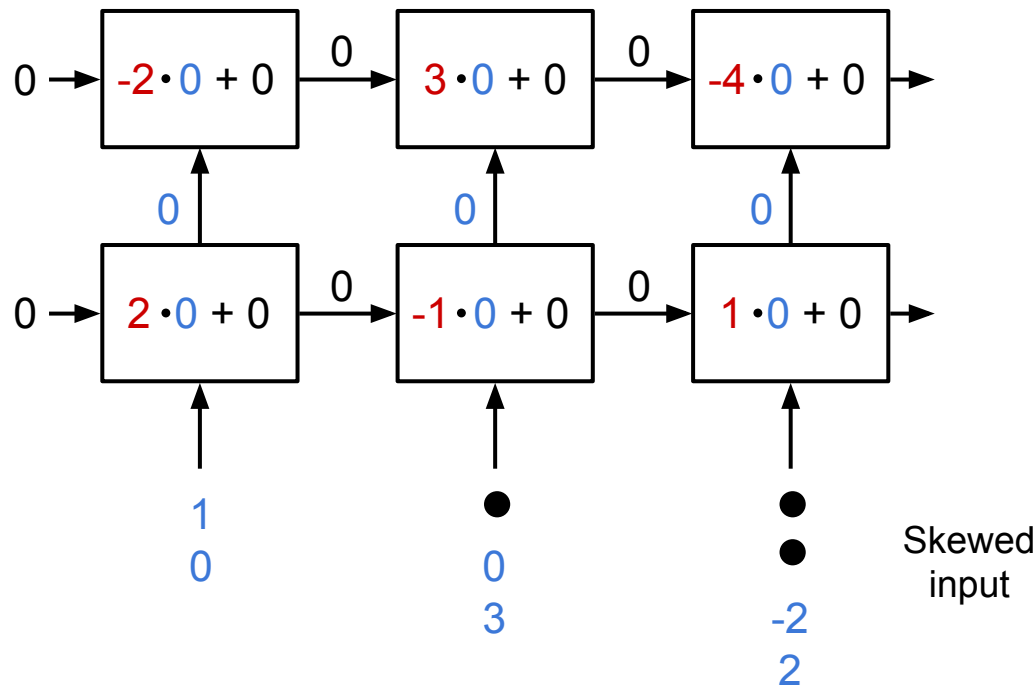


# Visualizing Systolic Array Multiplication

Weight Matrix      Data Matrix      Result Matrix

$$\begin{bmatrix} 2 & -1 & 1 \\ -2 & 3 & -4 \end{bmatrix} \times \begin{bmatrix} 1 & 0 \\ 0 & 3 \\ -2 & 2 \end{bmatrix} = \begin{bmatrix} 0 & -1 \\ 6 & 1 \end{bmatrix}$$

Weights in **red** are preloaded into the systolic array

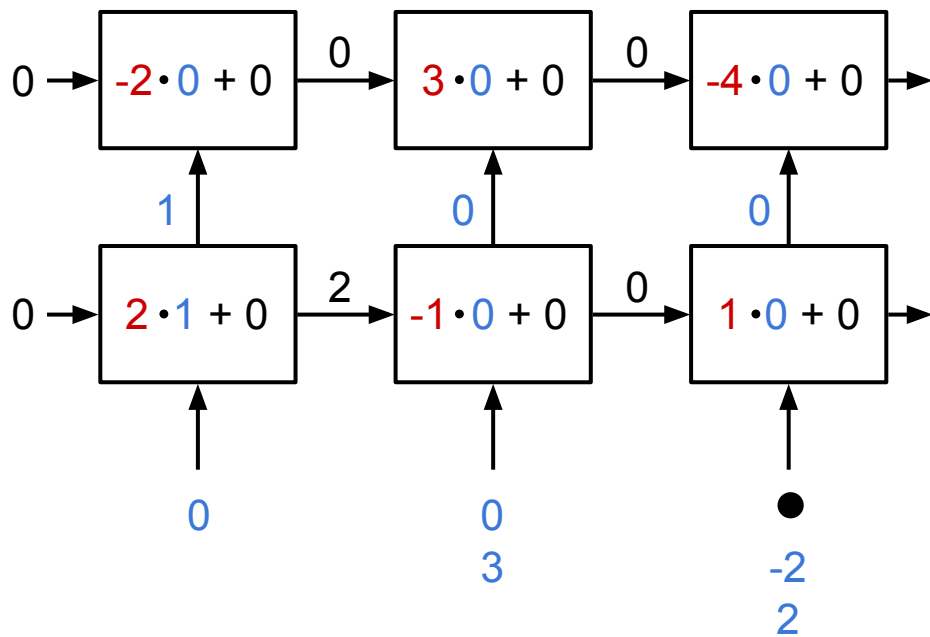


# Visualizing Systolic Array Multiplication

Weight Matrix      Data Matrix      Result Matrix

$$\begin{bmatrix} 2 & -1 & 1 \\ -2 & 3 & -4 \end{bmatrix} \times \begin{bmatrix} 1 & 0 \\ 0 & 3 \\ -2 & 2 \end{bmatrix} = \begin{bmatrix} 0 & -1 \\ 6 & 1 \end{bmatrix}$$

Weights in red are preloaded into the systolic array



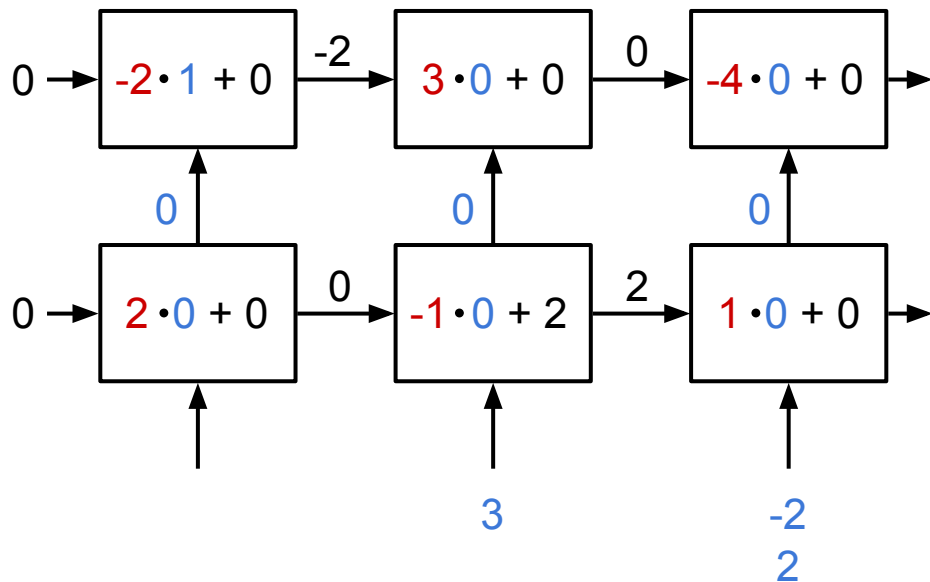


# Visualizing Systolic Array Multiplication

Weight Matrix      Data Matrix      Result Matrix

$$\begin{bmatrix} 2 & -1 & 1 \\ -2 & 3 & -4 \end{bmatrix} \times \begin{bmatrix} 1 & 0 \\ 0 & 3 \\ -2 & 2 \end{bmatrix} = \begin{bmatrix} 0 & -1 \\ 6 & 1 \end{bmatrix}$$

Weights in red are preloaded into the systolic array

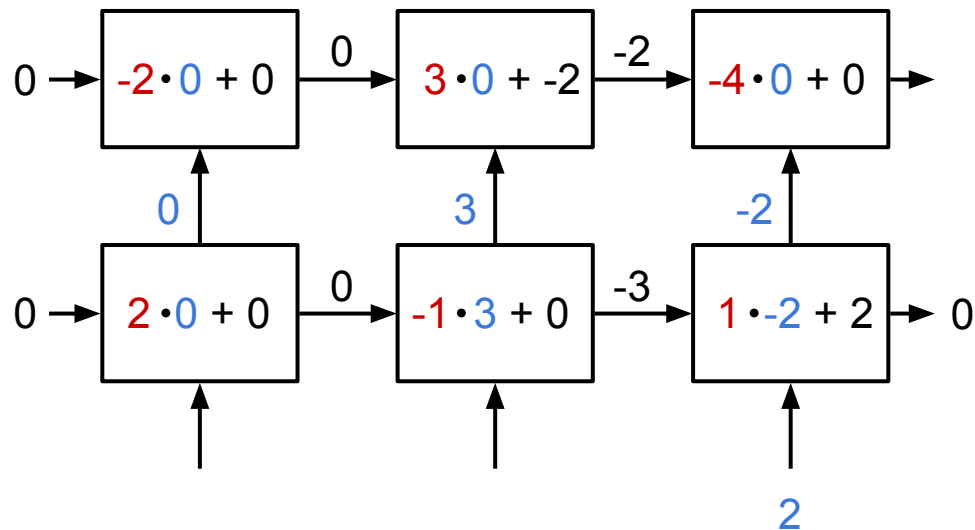


# Visualizing Systolic Array Multiplication

Weight Matrix      Data Matrix      Result Matrix

$$\begin{bmatrix} 2 & -1 & 1 \\ -2 & 3 & -4 \end{bmatrix} \times \begin{bmatrix} 1 & 0 \\ 0 & 3 \\ -2 & 2 \end{bmatrix} = \begin{bmatrix} 0 & -1 \\ 6 & 1 \end{bmatrix}$$

Weights in **red** are preloaded into the systolic array

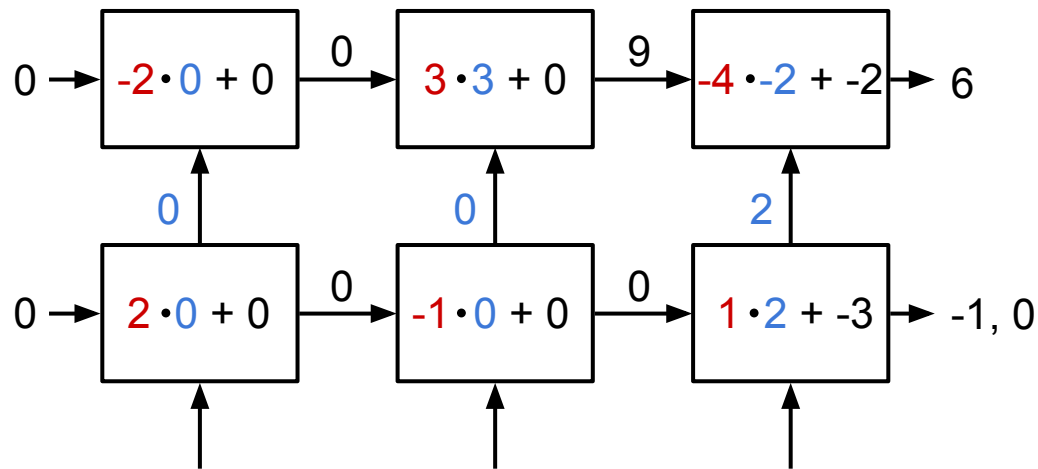


# Visualizing Systolic Array Multiplication

Weight Matrix      Data Matrix      Result Matrix

$$\begin{bmatrix} 2 & -1 & 1 \\ -2 & 3 & -4 \end{bmatrix} \times \begin{bmatrix} 1 & 0 \\ 0 & 3 \\ -2 & 2 \end{bmatrix} = \begin{bmatrix} 0 & -1 \\ 6 & 1 \end{bmatrix}$$

Weights in red are preloaded into the systolic array

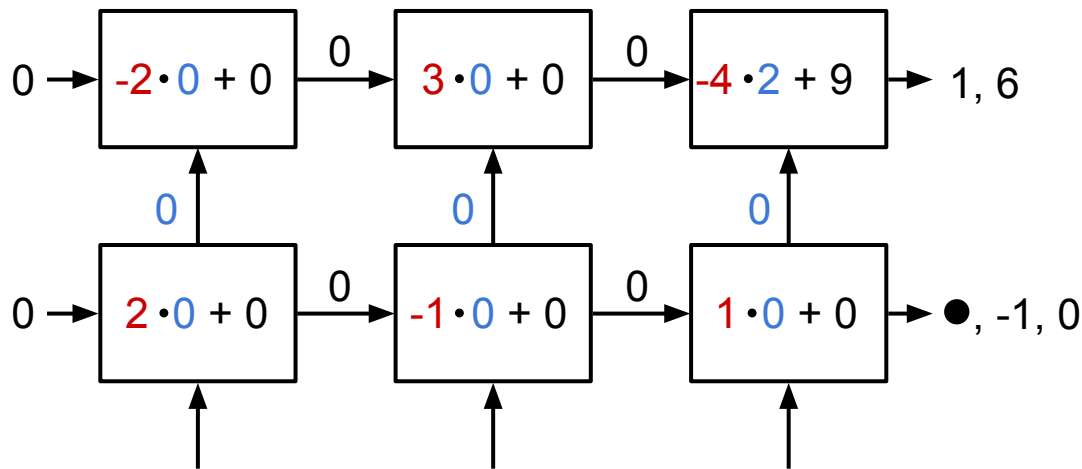


# Visualizing Systolic Array Multiplication

Weight Matrix      Data Matrix      Result Matrix

$$\begin{bmatrix} 2 & -1 & 1 \\ -2 & 3 & -4 \end{bmatrix} \times \begin{bmatrix} 1 & 0 \\ 0 & 3 \\ -2 & 2 \end{bmatrix} = \begin{bmatrix} 0 & -1 \\ 6 & 1 \end{bmatrix}$$

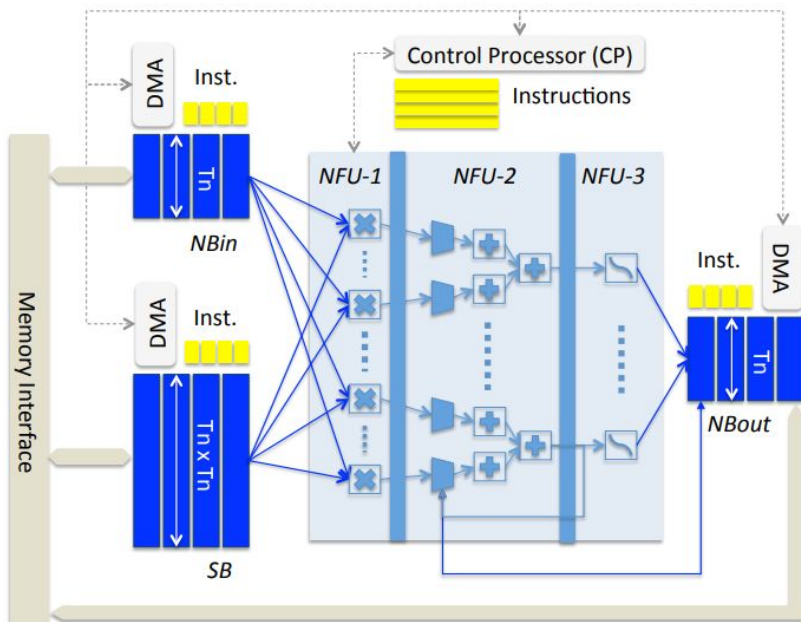
Weights in **red** are preloaded into the systolic array



# Topics

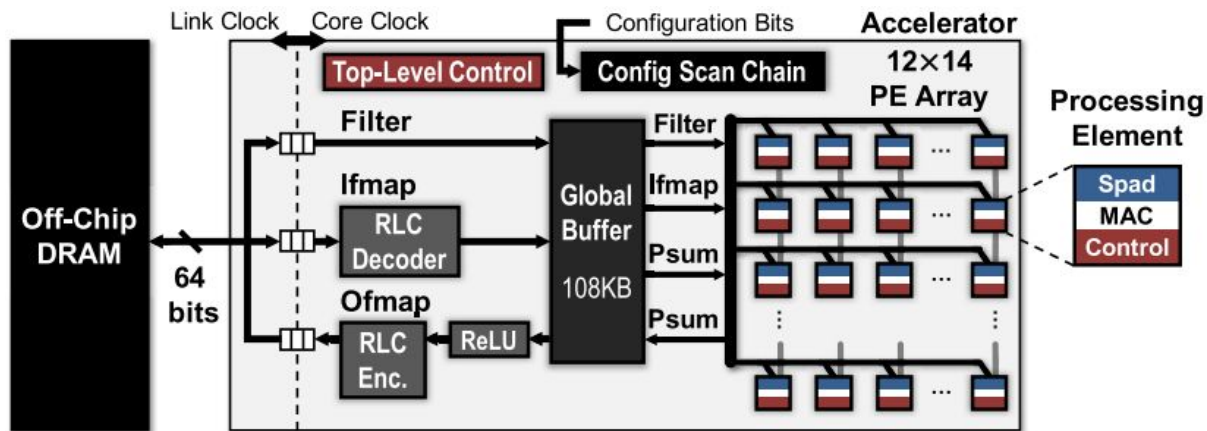
- Convolutional operation conversion
- Hardware architecture of CNN accelerator
- Systolic array
- Popular accelerator design
  - Eyeriss
  - Diannao
  - Cnvlutin
  - EIE

# Diannao



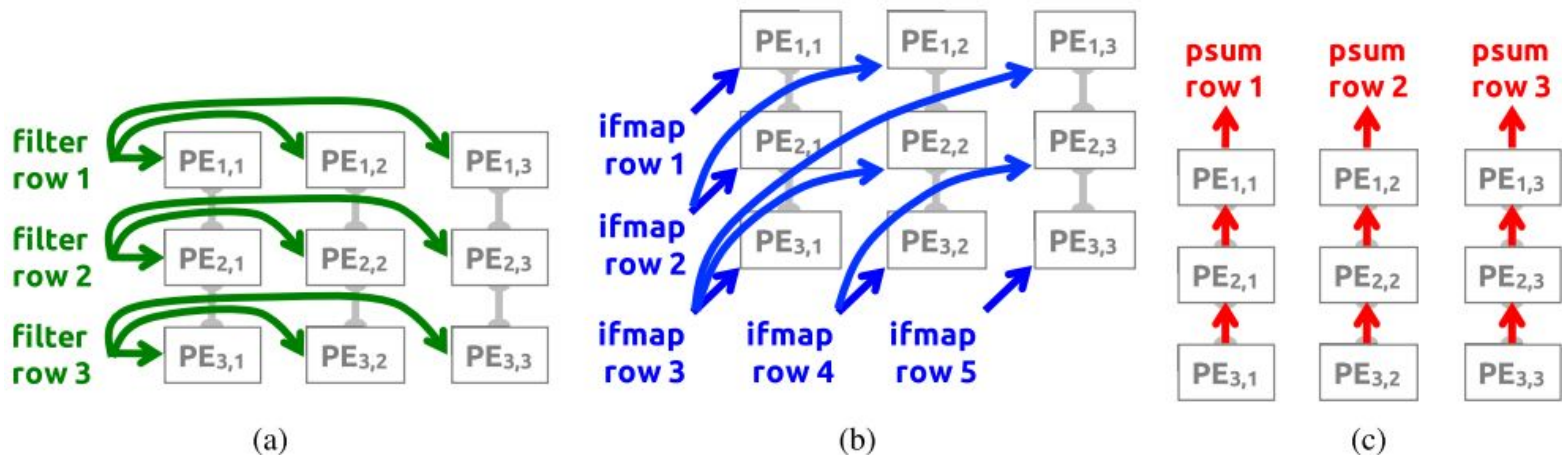
- The first popular end-to-end DNN (CNN) accelerator.
- Diannao is synthesized with 65nm using Synopsys tools, achieving a throughput of 482 GOP/s.
- NFU consists of three stages:
  - Multiplier units
  - Adder tree
  - Nonlinear unit

# Eyeriss



- Eyeriss optimizes for the energy efficiency of the entire system, including the accelerator chip and off-chip DRAM, for various CNN shapes by reconfiguring the architecture.
- The core clock domain consists of a spatial array of 168 PEs organized as a  $12 \times 14$  rectangle, a 108-kB GLB, an RLC CODEC, and an ReLU module.

# Data Reuse for Memory Access Reduction

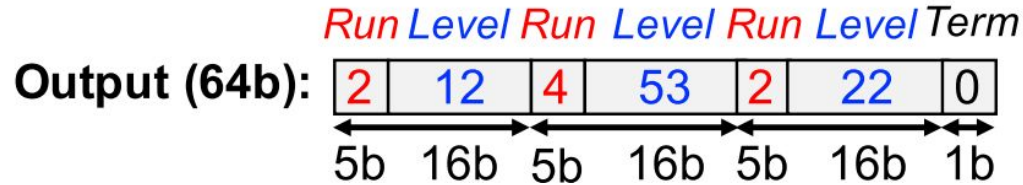


- Reuse and accumulation of data within a PE set reduce accesses to the GLB and DRAM, saving data movement energy cost.



# Rerun-length encoding

Input: 0, 0, 12, 0, 0, 0, 0, 53, 0, 0, 22, ...

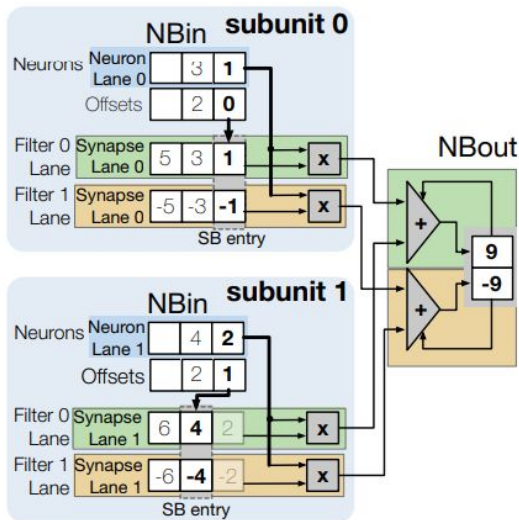


- RLC is used for compressing the input activation.

# Cnvlutin

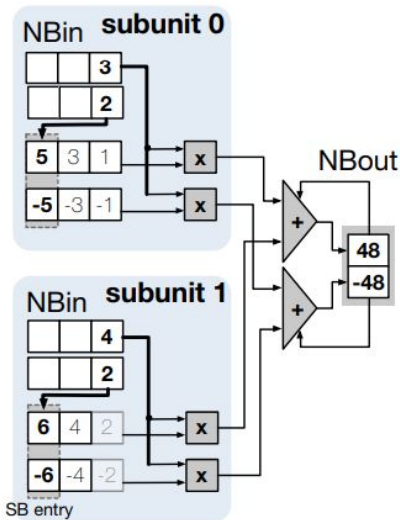
Input:  $[1, 0, 3] \rightarrow [1, 3]$  (input)  $[0, 2]$  (offset)  
 Weight:  $[1, 3, 5]$

cycle 0



(a)

cycle 1



(b)

- A large fraction of the computations performed by CNNs are intrinsically ineffectual as they involve a multiplication where one of the inputs is zero.
- Cnvlutin is a value-based approach to hardware acceleration that eliminates most of these ineffectual operations, improving performance and energy over a state-of-the-art accelerator with no accuracy loss.

# Presentation

- [Kangaroo: Lossless Self-Speculative Decoding via Double Early Exiting](#) (Roshan)
- [SCNN: An Accelerator for Compressed-sparse Convolutional Neural Networks](#)  
(Lavanya, Murali)

