

# Lecture 11: Al Accelerator Introduction and CNN Accelerators

#### **Notes**

- First round of team meeting on Dec 1 and Dec 2.
- Extra-credit quiz today.



## Recap

- Federated Learning
- Machine Learning Compiler
- Machine Learning System



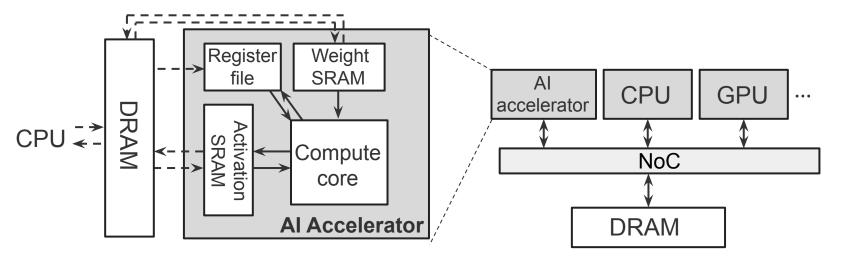
## **Topics**

- Hardware accelerator: Overview
- Convolutional operation conversion
- Systolic array
- Convolutional Neural Network System
- Popular accelerator design



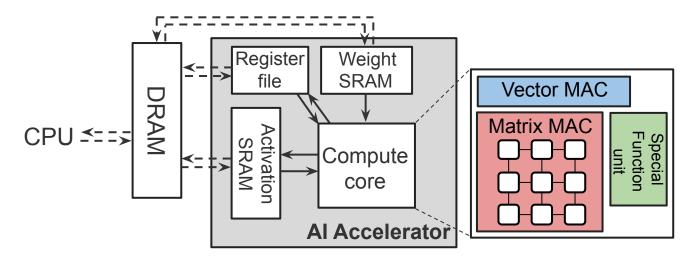
#### **AI Accelerator**

• The Al accelerator can execute part of the machine code that is related to the Al workload.





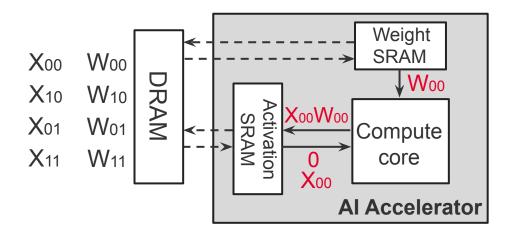
#### **AI Accelerator**



- 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.



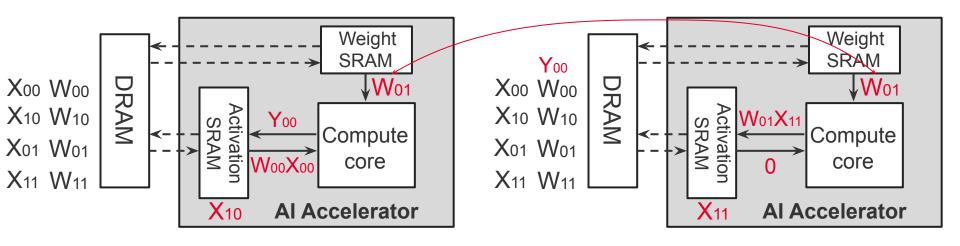
#### **AI Accelerator**



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



## **Memory Access Reduction**



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

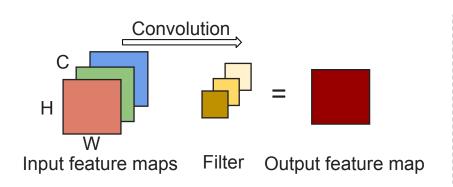


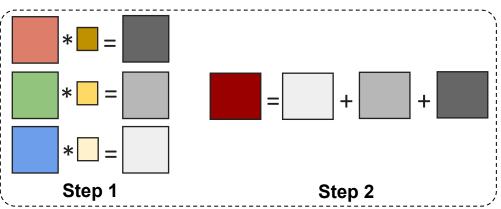
## **Topics**

- Hardware accelerator: Overview
- Convolutional operation conversion
- Systolic array
- Convolutional Neural Network System
- Popular accelerator design



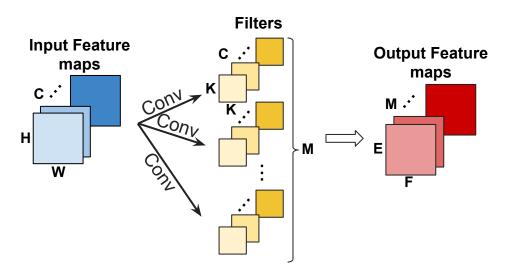
## **Convolutional Layers**





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





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

C: number of input channels

H,W: size of the input feature maps

M: number of weight filters

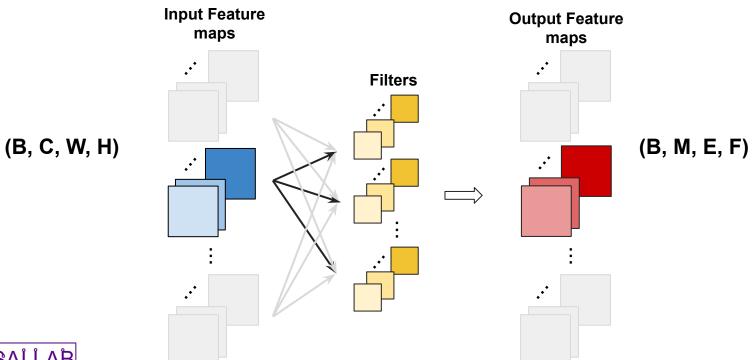
K: weight kernel size

E,F: size of the output feature maps



**Input Feature Output Feature** maps maps **Filters** (B, C, W, H) (B, M, E, F)





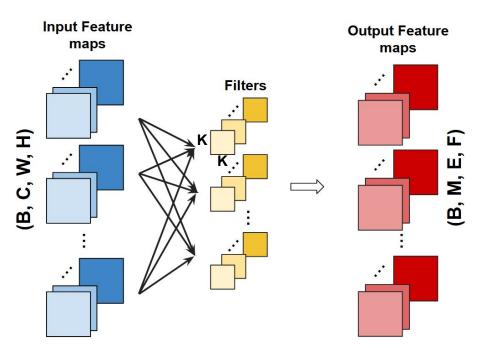


(B, C, W, H)

**Input Feature Output Feature** maps maps **Filters** (B, M, E, F)



#### **Computational Cost: Standard Convolution**



- Number of MACs: B×M×K×K×C×E×F
- Storage cost: 32×(M×C×K×K+B×C×H×W+B×M×E×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:
  - o 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 k<sub>1</sub> = 1 to K

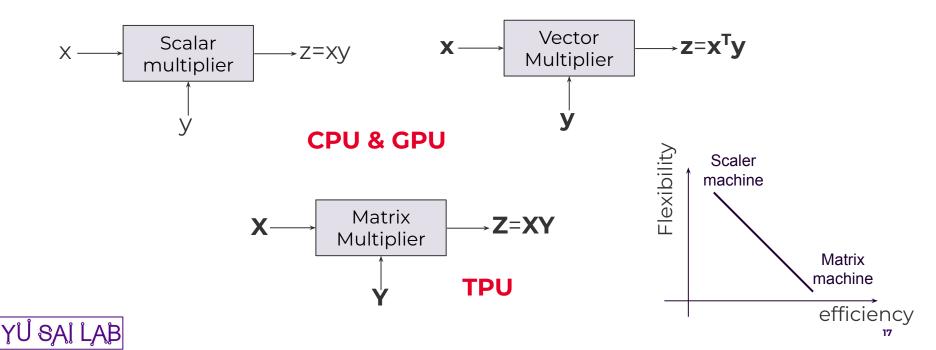
for k<sub>2</sub> = 1 to K

out[b][m][e][f] += in[b][c][e+k<sub>1</sub>-(K+1)/2][f+k<sub>2</sub>-(K+1)/2] * filter[m][c][k<sub>1</sub>][k<sub>2</sub>];
```

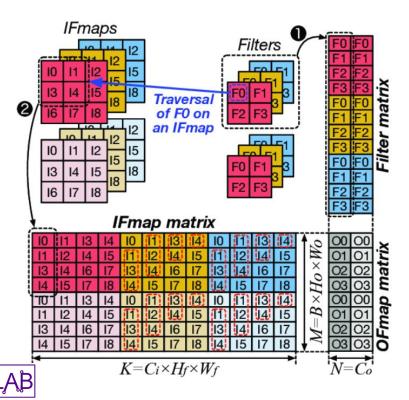
- This simple loop nest can be transformed in numerous ways to capture different reuse patterns of the activations and weights and to map the computation to a hardware accelerator implementation.
- A CNN's dataflow defines how the loops are ordered, partitioned, and parallelized
- We can use the scaler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



# **Computational Dataflow for CNN**

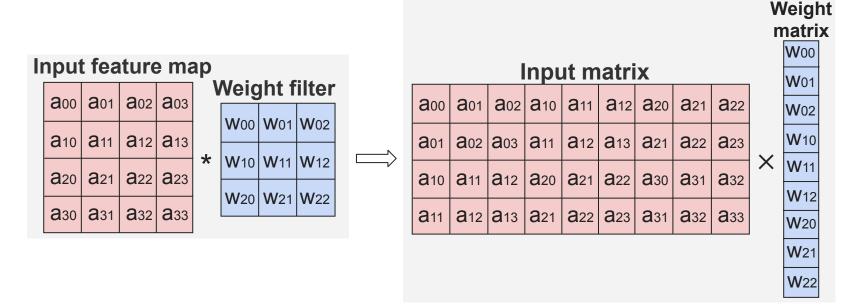


#### **How to Convert to Matrix Multiplication?**



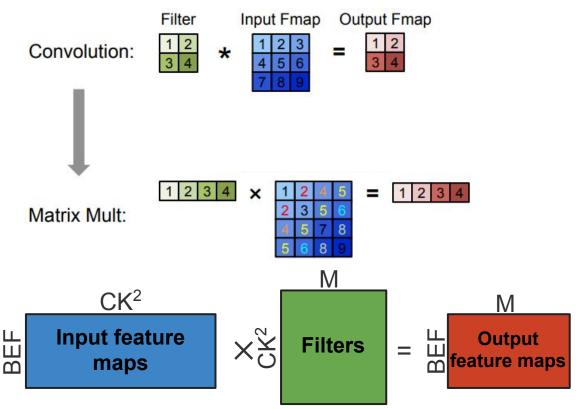
 A standard Convolutional operation can be converted to 2D matrix multiplication using Im2Col operations.

#### **How to Convert to Matrix Multiplication?**





#### **How to Convert to Matrix Multiplication?**



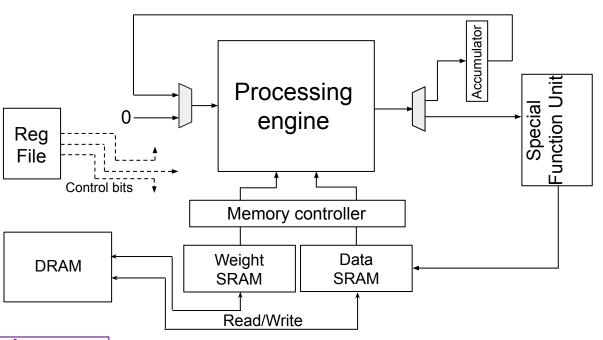


## **Topics**

- Hardware accelerator: Overview
- Convolutional operation conversion
- Systolic array
- Convolutional Neural Network System
- Popular accelerator design



## **Hardware Architectures for DNN Processing**

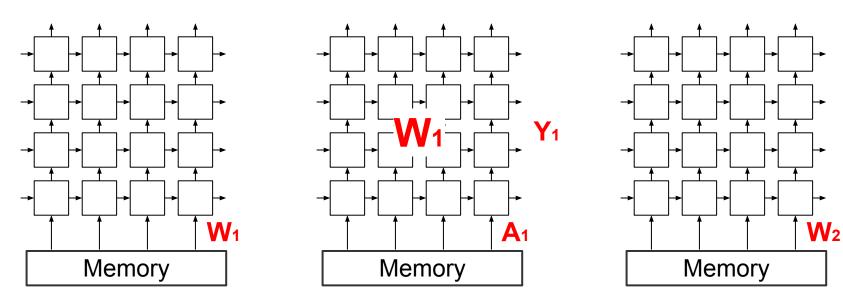


Major building blocks:

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



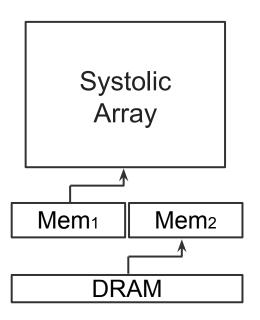
# **Computing Paradigms**

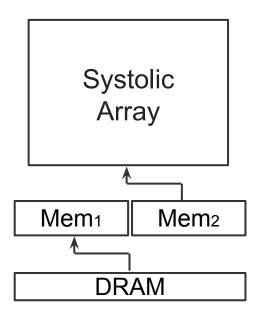


Spatial architecture can achieve great reuse of the extracted content, leading to a reduced memory access cost.



# **Double Buffering**



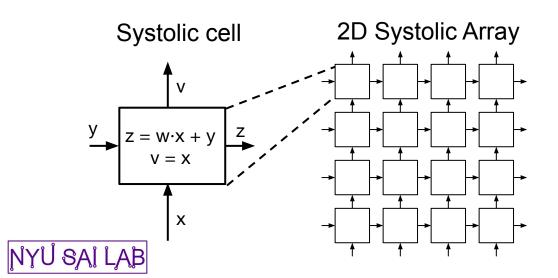


- Double buffering in hardware design is a technique used to improve the efficiency and performance of data processing, especially in systems that require smooth and continuous data transfer.
- The idea is to overlap the data production and consumption processes to avoid delays.



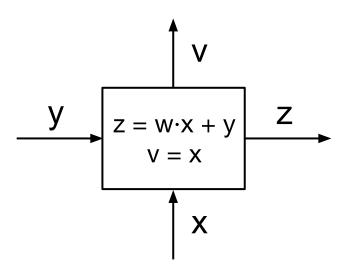
## Systolic Array (Weight Stationary Version)

- 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



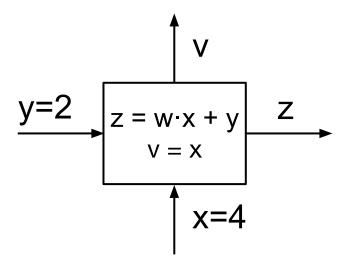


TPU (Google)

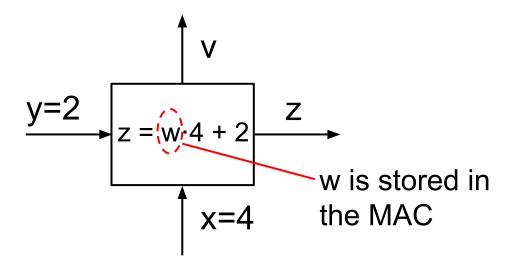


- Takes data (x and y) as input
- w stays in the systolic cell
- Performs a multiply-accumulate operation

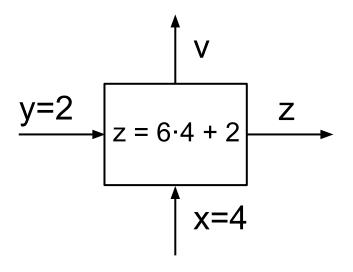




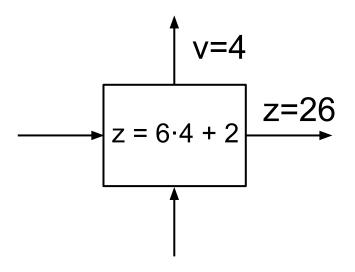






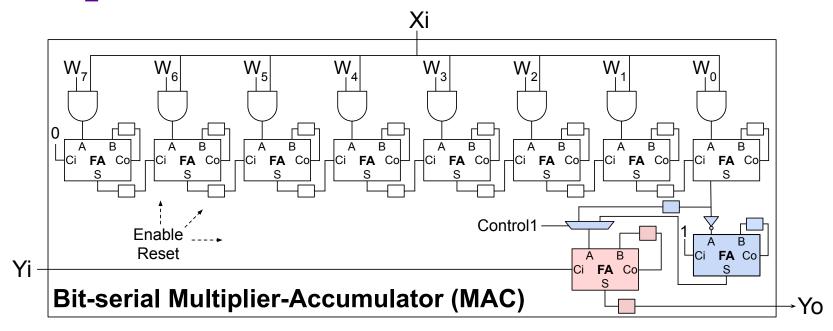






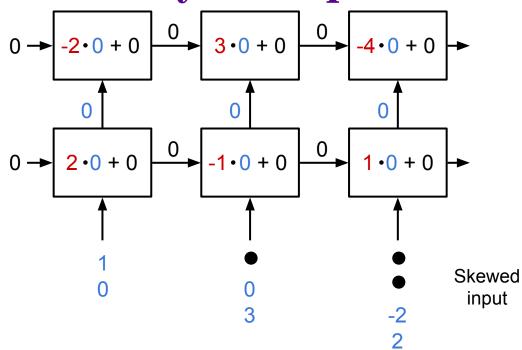


## Multiplier Accumulator



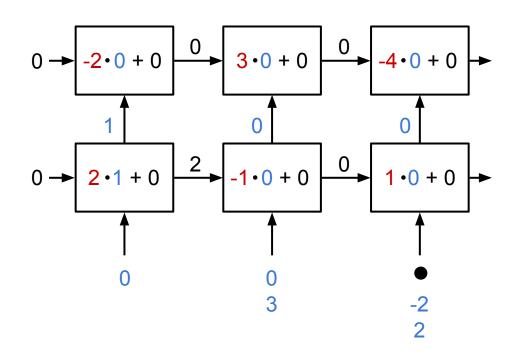


Weight Data Result Matrix Matrix 
$$\begin{bmatrix} 2 & -1 & 1 \\ -2 & 3 & -4 \end{bmatrix} \times \begin{bmatrix} 1 & 0 \\ 0 & 3 \\ -2 & 2 \end{bmatrix} = \begin{bmatrix} 0 & -1 \\ 6 & 1 \end{bmatrix}$$



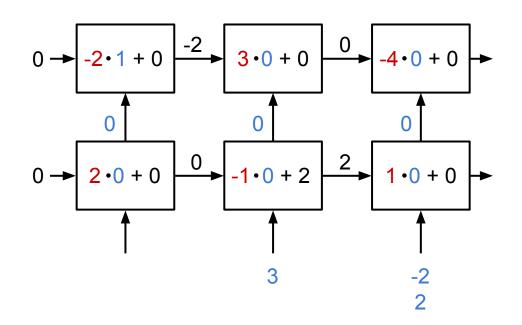


Weight Data Result Matrix Matrix 
$$\begin{bmatrix} 2 & -1 & 1 \\ -2 & 3 & -4 \end{bmatrix} \times \begin{bmatrix} 1 & 0 \\ 0 & 3 \\ -2 & 2 \end{bmatrix} = \begin{bmatrix} 0 & -1 \\ 6 & 1 \end{bmatrix}$$





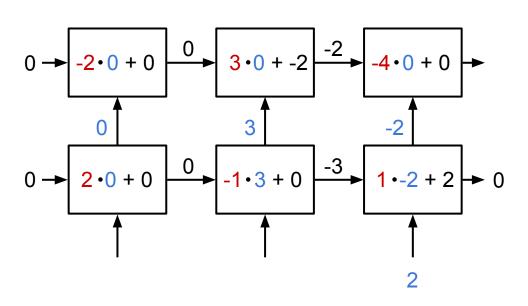
Weight Data Result Matrix Matrix 
$$\begin{bmatrix} 2 & -1 & 1 \\ -2 & 3 & -4 \end{bmatrix} \times \begin{bmatrix} 1 & 0 \\ 0 & 3 \\ -2 & 2 \end{bmatrix} = \begin{bmatrix} 0 & -1 \\ 6 & 1 \end{bmatrix}$$



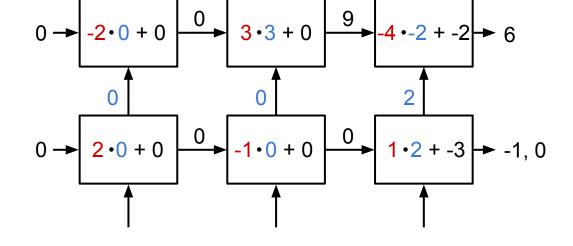


Weight Data Result Matrix Matrix 
$$\begin{bmatrix} 2 & -1 & 1 \\ -2 & 3 & -4 \end{bmatrix} \times \begin{bmatrix} 1 & 0 \\ 0 & 3 \\ -2 & 2 \end{bmatrix} = \begin{bmatrix} 0 & -1 \\ 6 & 1 \end{bmatrix}$$





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

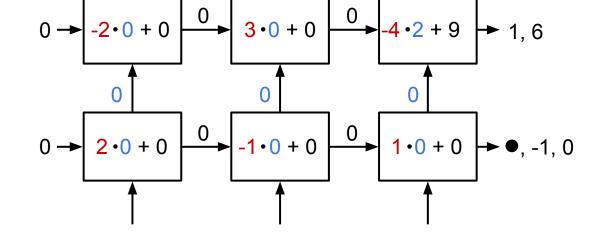




## Visualizing Systolic Array Multiplication

Weight Data Result Matrix 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



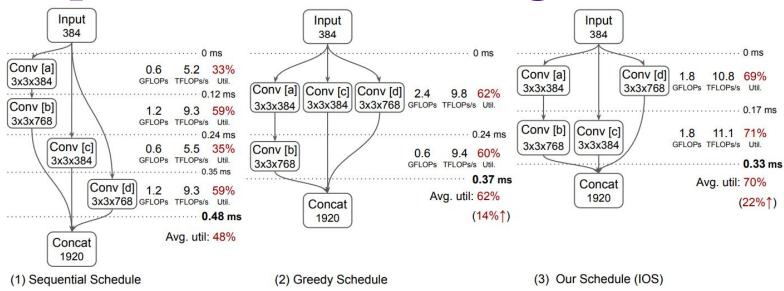


## **Topics**

- Hardware accelerator: Overview
- Convolutional operation conversion
- Systolic array
- Convolutional Neural Network System
- Popular accelerator design



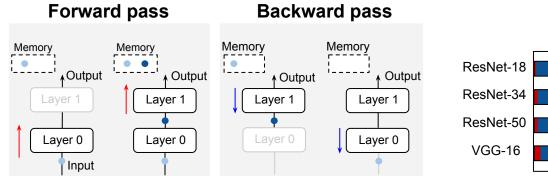
### **Computational Scheduling**

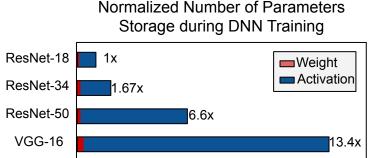


- The branchy CNN can be scheduling and computed in a much higher efficiency.
- Two convolutional operations can be combined to achieve less memory cost.



Ding, Yaoyao, et al. "los: Inter-operator scheduler for cnn acceleration." *Proceedings of Machine Learning and Systems* 3 (2021): 167-180.

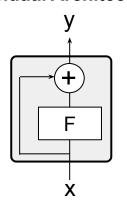




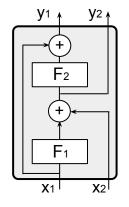
- The memory footprint grows proportional with the layer depth.
- On top of this, small edge devices typically have limited on-chip storage, leading to frequent and costly accesses to off-chip memories.



#### **Residual Architecture**



#### **Reversible Architecture**



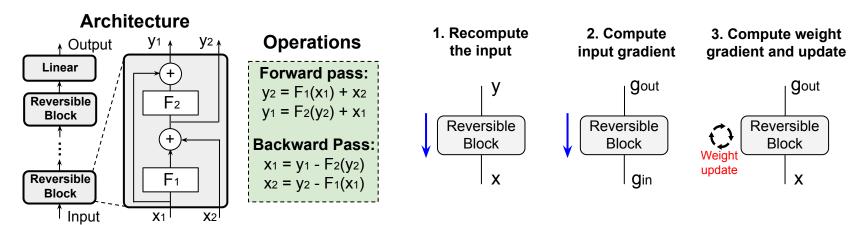
#### Forward pass:

$$y_2 = F_1(x_1) + x_2$$
  
 $y_1 = F_2(y_2) + x_1$ 

#### **Backward Pass:**

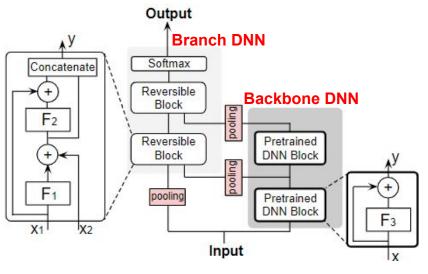
$$x_1 = y_1 - F_2(y_2)$$
  
 $x_2 = y_2 - F_1(x_1)$ 

 A reversible residual network (RevNet) is a variant of the canonical residual neural network (ResNet).



- The reversible architecture enables the backward pass computations to be performed without the need to store the input activations.
- Given the output y, the input activations are first recomputed. Afterwards, the input and weight gradients are computed with standard backward pass operations.



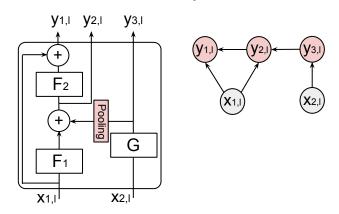


**Duplex DNN** 

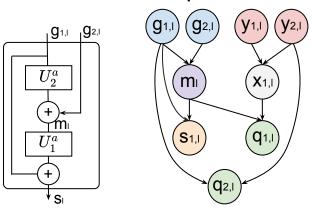
- This approach in turn imposes higher compute demands.
- We propose to judiciously train a subset of the model parameters to minimize training.
- The backbone DNN is frozen during the backward pass of the DNN.
- The normalization layers are removed from the branch DNN to facilitate the training process.



#### Forward pass



#### **Backward pass**



- We carefully schedule the computations of DuDNN during both forward and backward passes to achieve optimal system performance.
- We propose the optimal compute pattern to minimize memory usage and lifetime.

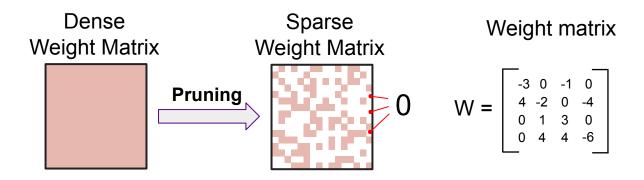


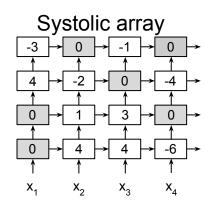
### **Topics**

- Hardware accelerator: Overview
- Convolutional operation conversion
- Systolic array
- Convolutional Neural Network System
- Popular accelerator design



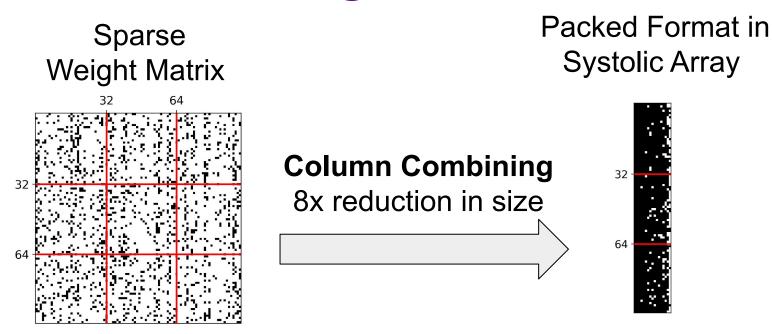
- Unimportant (i.e., small weights) are set to 0
- However, it is hard to leverage these zero weights to reduce the hardware cost





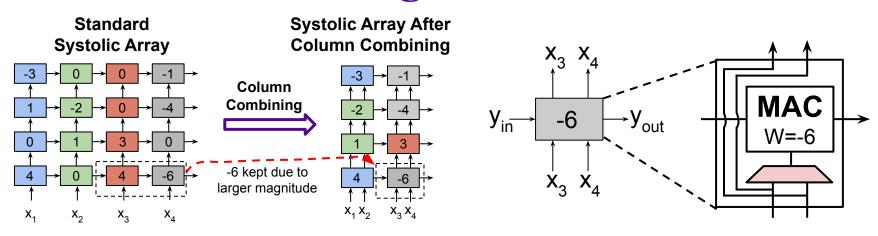


Kung, H. T., Bradley McDanel, and Sai Qian Zhang. "Packing sparse convolutional neural networks for efficient systolic array implementations: Column combining under joint optimization." *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*. 2019.



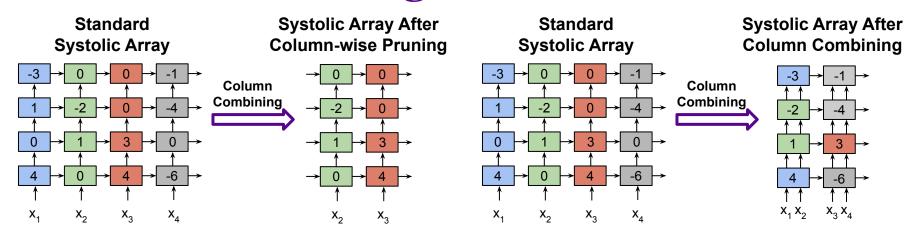


Kung, H. T., Bradley McDanel, and Sai Qian Zhang. "Packing sparse convolutional neural networks for efficient systolic array implementations: Column combining under joint optimization." *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*. 2019.



Column combining can greatly increase the utilization of the systolic array.

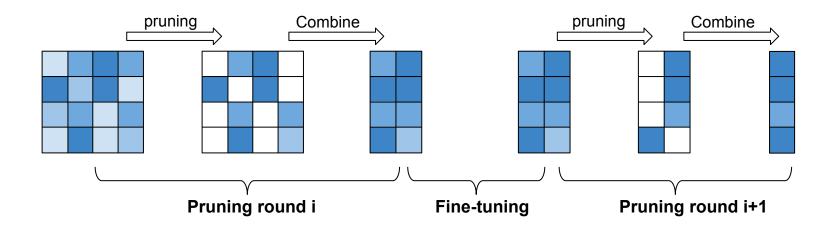




- Compared with column-wise Pruning (filterwise pruning), Column Combining allows for a much more flexible pruning pattern
- We can also apply Column Combining pruning on the input of each layer.



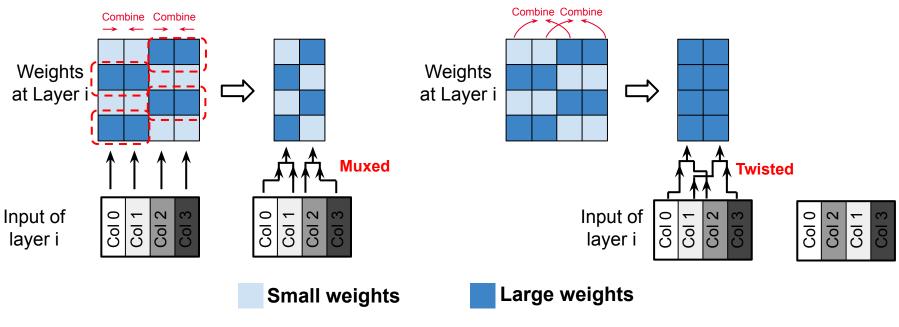
# **Column-Combining Pruning**





Kung, H. T., Bradley McDanel, and Sai Qian Zhang. "Packing sparse convolutional neural networks for efficient systolic array implementations: Column combining under joint optimization." *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*. 2019.

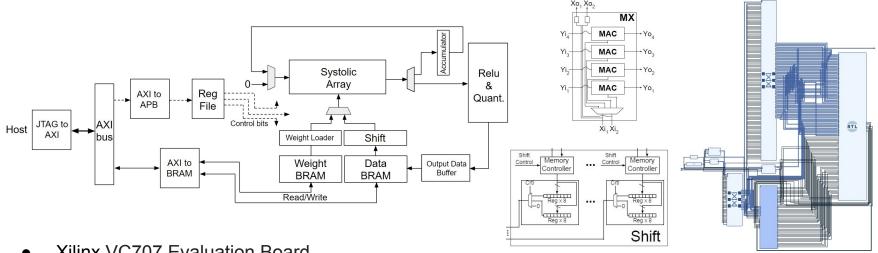
#### Column-Combining Pruning: Row Permutation





Kung, H. T., Bradley McDanel, and Sai Qian Zhang. "Packing sparse convolutional neural networks for efficient systolic array implementations: Column combining under joint optimization." *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*. 2019.

#### **Hardware Implementation**

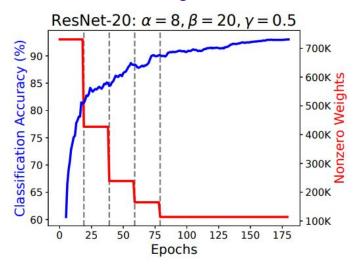


- Xilinx VC707 Evaluation Board
- Total hardware available: Lookup Table (303600), Flip-Flops (607200), BRAM (1030, each 36Kb)
- More than 15K lines of verilog code
- 128 by 64 systolic array



Kung, H. T., Bradley McDanel, and Sai Qian Zhang. "Packing sparse convolutional neural networks for efficient systolic array implementations: Column combining under joint optimization." Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems. 2019.

#### **Accuracy Evaluation Results**



ResNet-20 on CIFAR-10

|                               | Accuracy |
|-------------------------------|----------|
| Original DNN                  | 72.08%   |
| Structured Filterwise Pruning | 69.0%    |
| Column Combining              | 71.81%   |

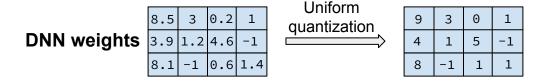
VGG-19 on ImageNet (87.5% sparsity)



Kung, H. T., Bradley McDanel, and Sai Qian Zhang. "Packing sparse convolutional neural networks for efficient systolic array implementations: Column combining under joint optimization." *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*. 2019.

## **Prior Work: Efficient DNN Data Types**

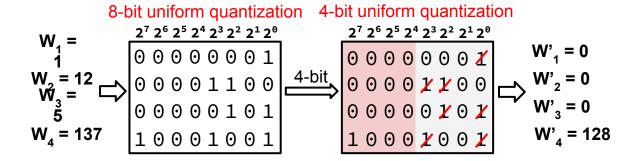
• Filter weights and activations can be quantized with low precision to accelerate the inference and reduce the model size.



Low-precision quantization leads to large accuracy loss.



#### **Problem on Low Precision Uniform Quantization**

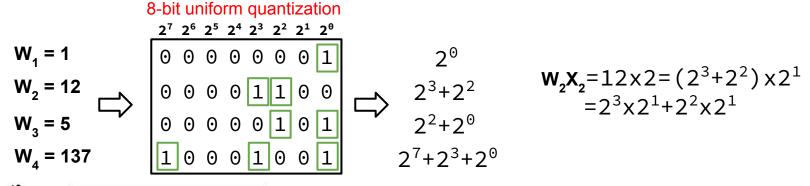


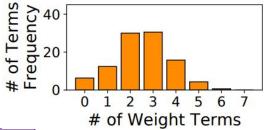
- Low-precision quantization leads to significant quantization error.
- Both weights and input activation are highly biased in values.



#### Representing Values in Power-of-two Terms

An integer value can be represented as a summation of power-of-two term(s).

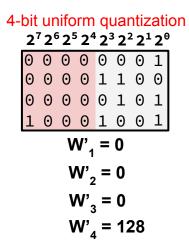


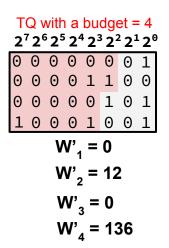


 Most quantized weight and data values can be represented with 2 or 3 power-of-two terms.



#### **Term Quantization**





- We can control the term-level computations by setting a group term budget.
- For a group of values, we rank and remove the small terms based on this budget.



#### Hybrid Encoding for Shortened Expressions (HESE)

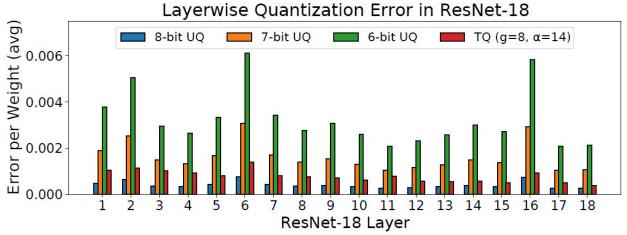
- To minimize the number of power-of-two terms in the binary input, we propose Hybrid Encoding for Shortened Expression (HESE)
- HESE offers:
  - Signed power-of-two expression with minimum length
  - Much less term-pair multiplications

| Binary expression | $31 = 2^4 + 2^3 + 2^2 + 2^1 + 2^0$ | 5 terms |
|-------------------|------------------------------------|---------|
|                   | $27 = 2^4 + 2^3 + 2^1 + 2^0$       | 4 terms |

| HESE | $31 = 2^5 - 2^0$       | 2 terms |
|------|------------------------|---------|
|      | $27 = 2^5 - 2^2 - 2^0$ | 3 terms |



# **Quantization Error Analysis**

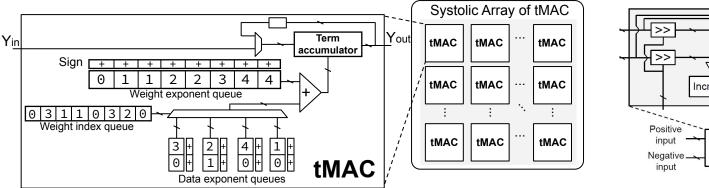


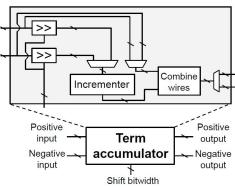
- We represent a group of 8 weights with 14 terms, each weight only requires 1.75 terms on average
- Term Quantization (TQ) introduces a small amount of quantization error over 8-bit uniform quantization (UQ)
- TQ achieves a much lower quantization error than 7-bit and 6-bit uniform quantization
- TQ with 1.75 term per value achieves a similar quantization error as 5-bit UQ



### Multiplier-accumulator Design

- We propose the term MAC (tMAC) for the efficient implementation of TQ.
- A tMAC processes all term-pair multiplications across a group of weight and data values.

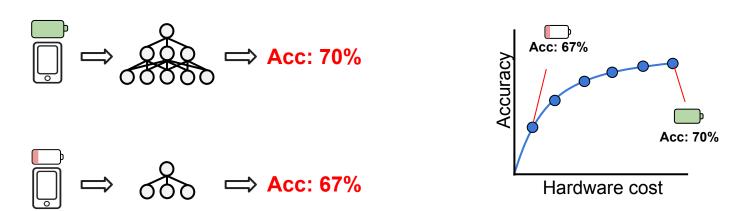




- Each term is represented by their corresponding exponent (2-3 bits).
- The term accumulation can be implemented using half adders.



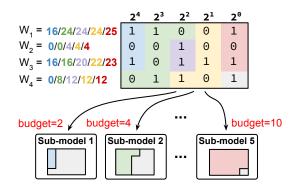
#### Multi-resolution DNN Inference with TQ



 DNN is expected to run at different resolution to achieve a good trade-off between hardware cost and accuracy.



#### Multi-resolution DNN Inference with TQ

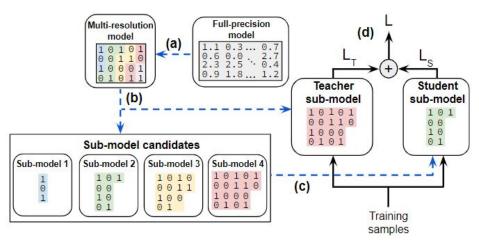


**Multi-resolution DNN model** 

 A meta multi-resolution DNN model which can work under different term budget needs to be trained.



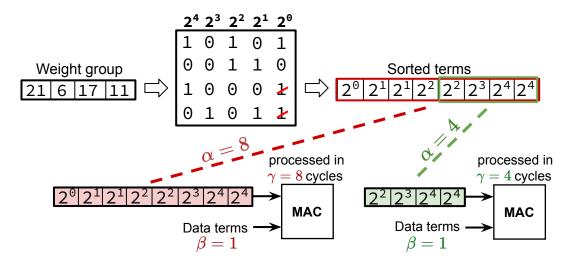
### **Multi-resolution DNN Training**



- We develop a multi-resolution training scheme to jointly train multiple DNN models under different term budgets.
- Instead of jointly train all the sub-models together, we apply the knowledge distillation framework to jointly train two sub-models per iteration.



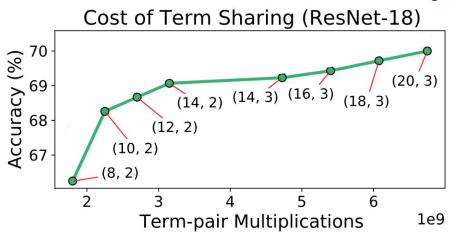
## **Multi-resolution DNN System**

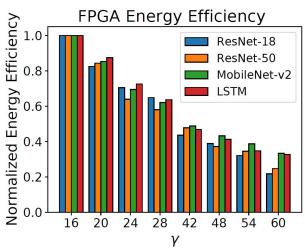


The energy consumption will scale with the term budgets of the weight and data.



### **Evaluation: Accuracy Performance**

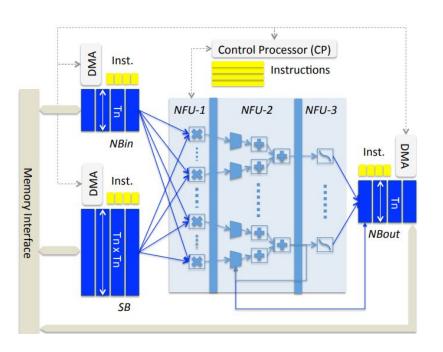




- The multi-resolution DNN incurs 0.4%-3.8% degradation compared to the original floating-point DNN (70.2%).
- The energy efficiency grows (3.25x on average), as term budget reduces from 60 to 16.



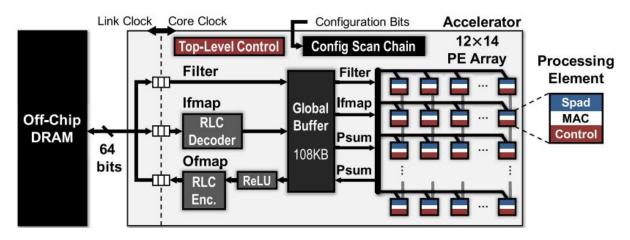
#### Diannao



- 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



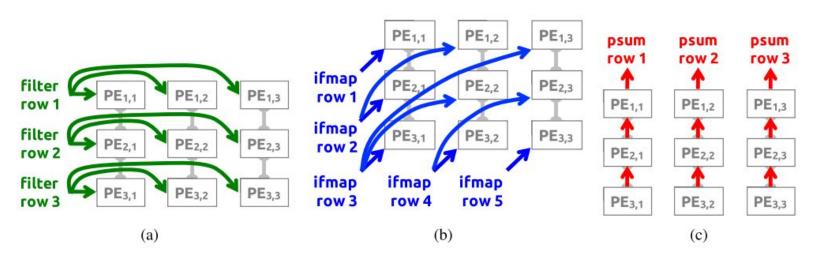
## **Eyeriss**



- 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 × 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, ...

Run Level Run Level Run Level Term

Output (64b): 2 12 4 53 2 22 0 5b 16b 5b 16b 5b 16b 1b

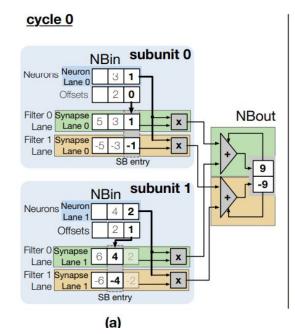
RLC is used for compressing the input activation.

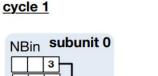


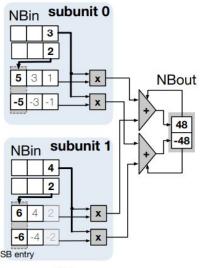
#### Cnvlutin

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

Weight: [1, 3, 5]







(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