

# Balancing Efficiency and Flexibility for DNN Acceleration via Temporal GPU-Systolic Array Integration

Cong Guo<sup>1</sup>, Yangjie Zhou<sup>1</sup>, Jingwen Leng<sup>1</sup>, Yuhao Zhu<sup>2</sup>, Zidong Du<sup>3</sup>,  
Quan Chen<sup>1</sup>, Chao Li<sup>1</sup>, Bin Yao<sup>1</sup>, Minyi Guo<sup>1</sup>

<sup>1</sup>Shanghai Jiao Tong University, <sup>2</sup>University of Rochester,

<sup>3</sup>Institute of Computing Technology, Chinese Academy of Sciences



UNIVERSITY of  
**ROCHESTER**



**中国科学院计算技术研究院**  
Institute Of Computing Technology Chinese Academy Of Sciences



# Biography

- Cong Guo
  - First-year Ph.D. student at Shanghai Jiao Tong University
  - Interested in computer architecture and high performance computing.
- Jingwen Leng
  - Associate professor
  - John Hopcroft Center for Computer Science  
Department of Computer Science and Engineering  
Shanghai Jiao Tong University



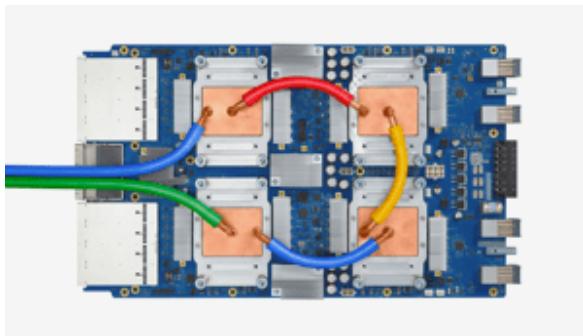
# Outline

- Introduction
- Simultaneous Multi-mode Architecture
- Evaluation



# Introduction

Efficiency



Google TPU [ISCA'16]

Flexibility



Nvidia GPU [Volta, 17]

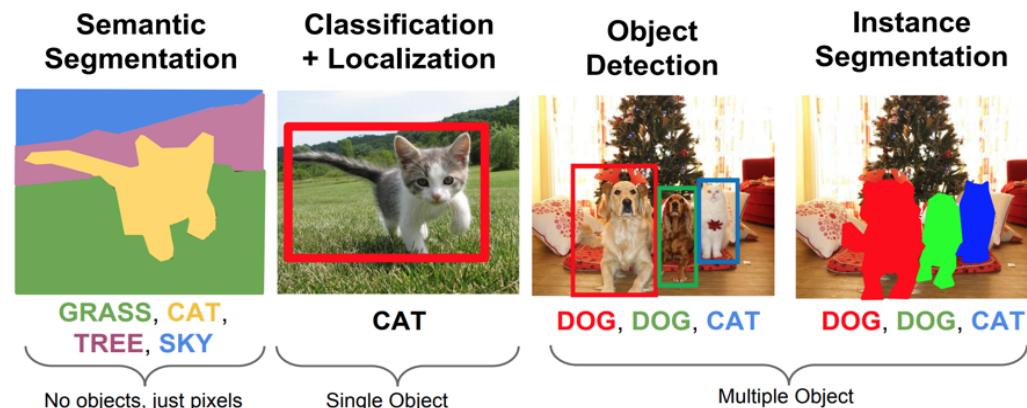
Specialized for  
GEMM (General Matrix Multiply)

General-purpose Core

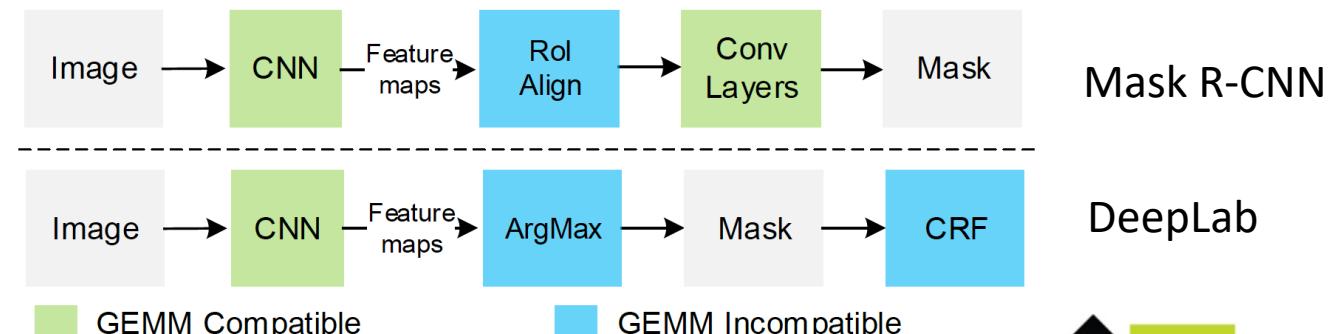


# Inefficiency for TPU

- TPU v2, 1 core, 22TFLOPS
- GPU V100, 15TFLOPS
- TPU profiling
  - Mask R-CNN
    - CNN/FC: **20%** faster than GPU
    - Total: **75%** slower than GPU
  - DeepLab
    - CNN/FC: **40%** faster than GPU
    - ArgMax: **2x** slowdown than GPU
    - CRF: **10x** worse than GPU

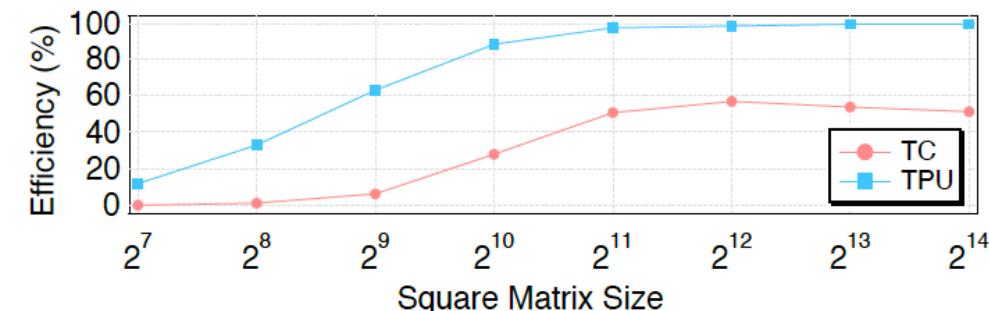
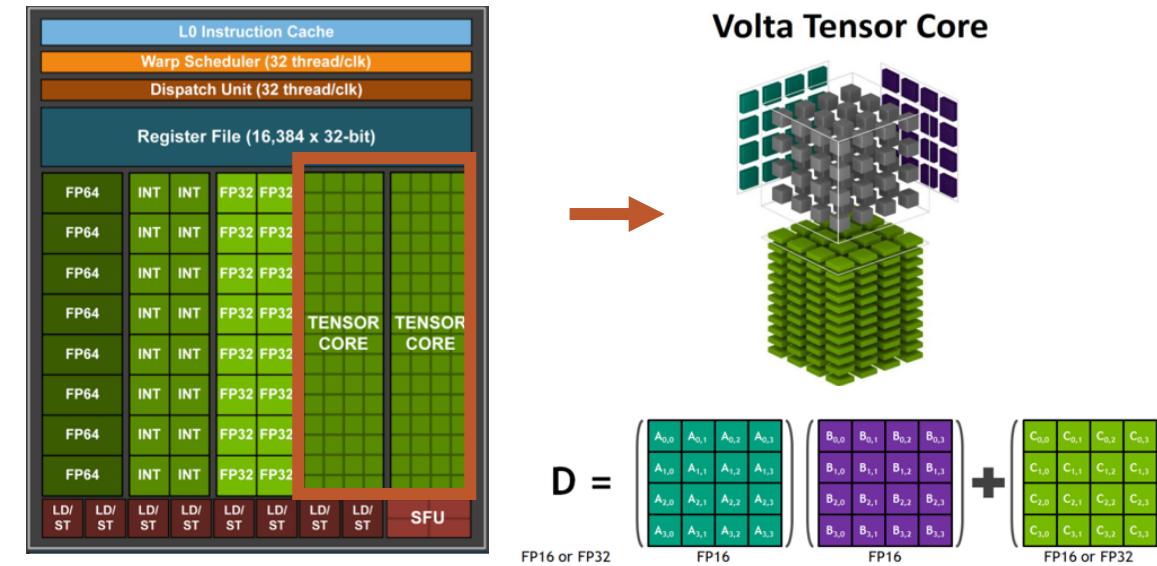


## Hybrid models:



# Inefficiency for GPU

- Spatial integration
- Explicit synchronization
- Fixed shape GEMM
- **Performance inefficiency**
- Tensor Core
  - Efficiency < 60%
- TPU
  - Efficiency > 99%



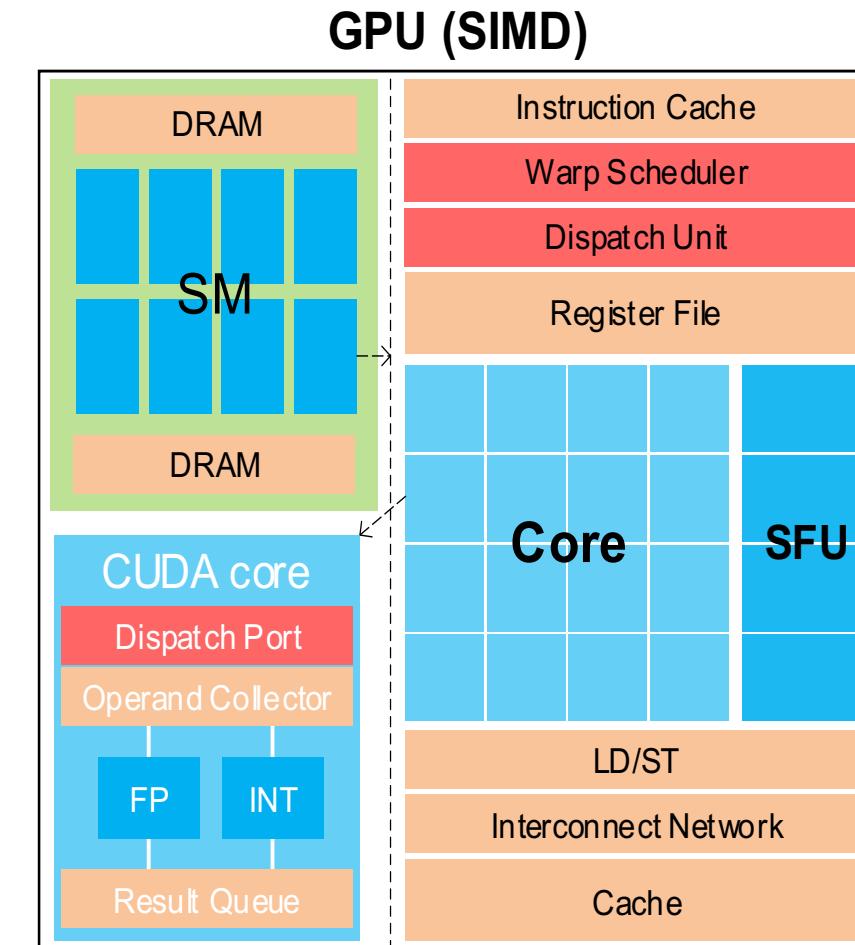
Efficiency: Achieved FLOPS divided by the peak FLOPS

# Outline

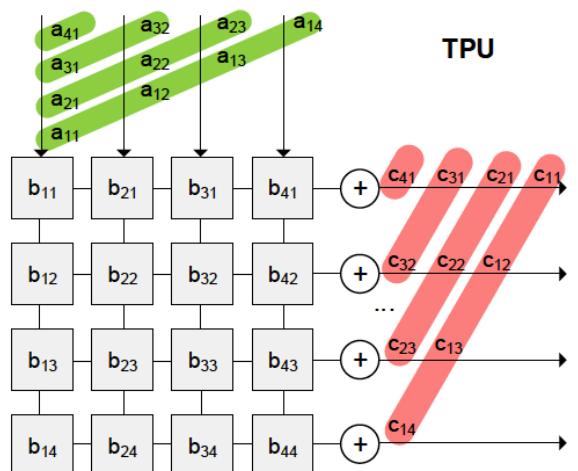
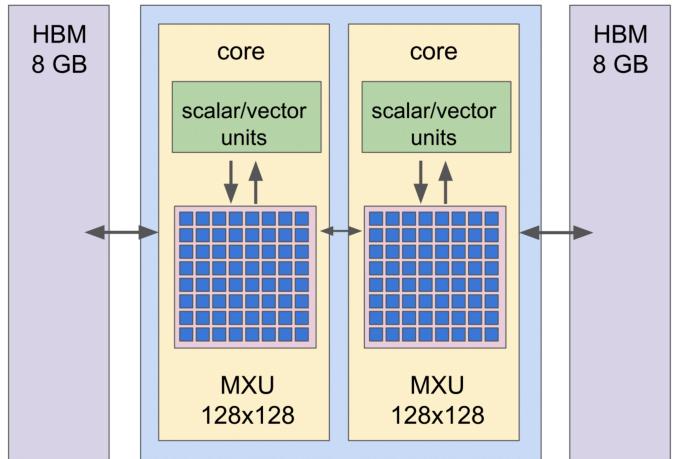
- Introduction
- Simultaneous Multi-mode Architecture
- Evaluation

# GPU

- GPU
- Parallelism
  - Single Instruction Multiple Data
  - Massive threads
  - Warp active mask
- Memory
  - Register file, vector access
  - Shared memory, scalar access
- Communication
  - PE array with shared memory
    - Warp Shuffle Instruction



# TPU



- TPU
- Papalism
  - 2D Systolic Array (MISD)
  - High concurrency
  - Active/Non-Active (one inst. two status)
- Memory
  - Weight, vector access (continuous)
  - Input/output, scalar access(discontinuous)
- Communication
  - Interconnected PE array

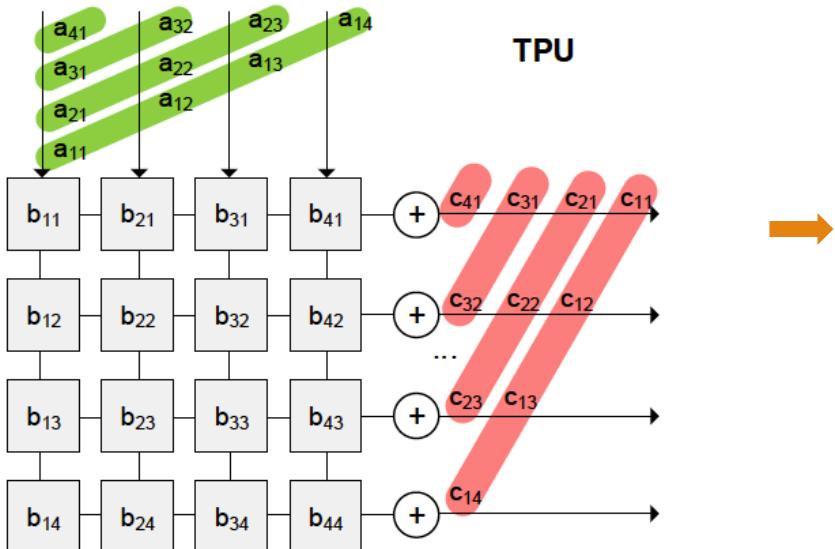


# Similarity

- GPU
- Parallelism
  - Single Instruction Multiple Data
  - Massive threads
  - Warp active mask
- Memory
  - Register file, vector access
  - Shared memory, scalar access
- Communication
  - PE array with shared memory
- TPU
- Parallelism
  - 2D Systolic Array (MISD)
  - High concurrency
  - Active/Non-Active (one inst. two status)
- Memory
  - Weight, vector access (continuous)
  - Input/output, scalar access(discontinuous)
- Communication
  - Interconnected PE array



# SMA Hardware Design



Similar to GPU

## Challenges:

- 1: Massive output scalar accesses
- 2: Inter-PE Communication
- 3: How to control the systolic array

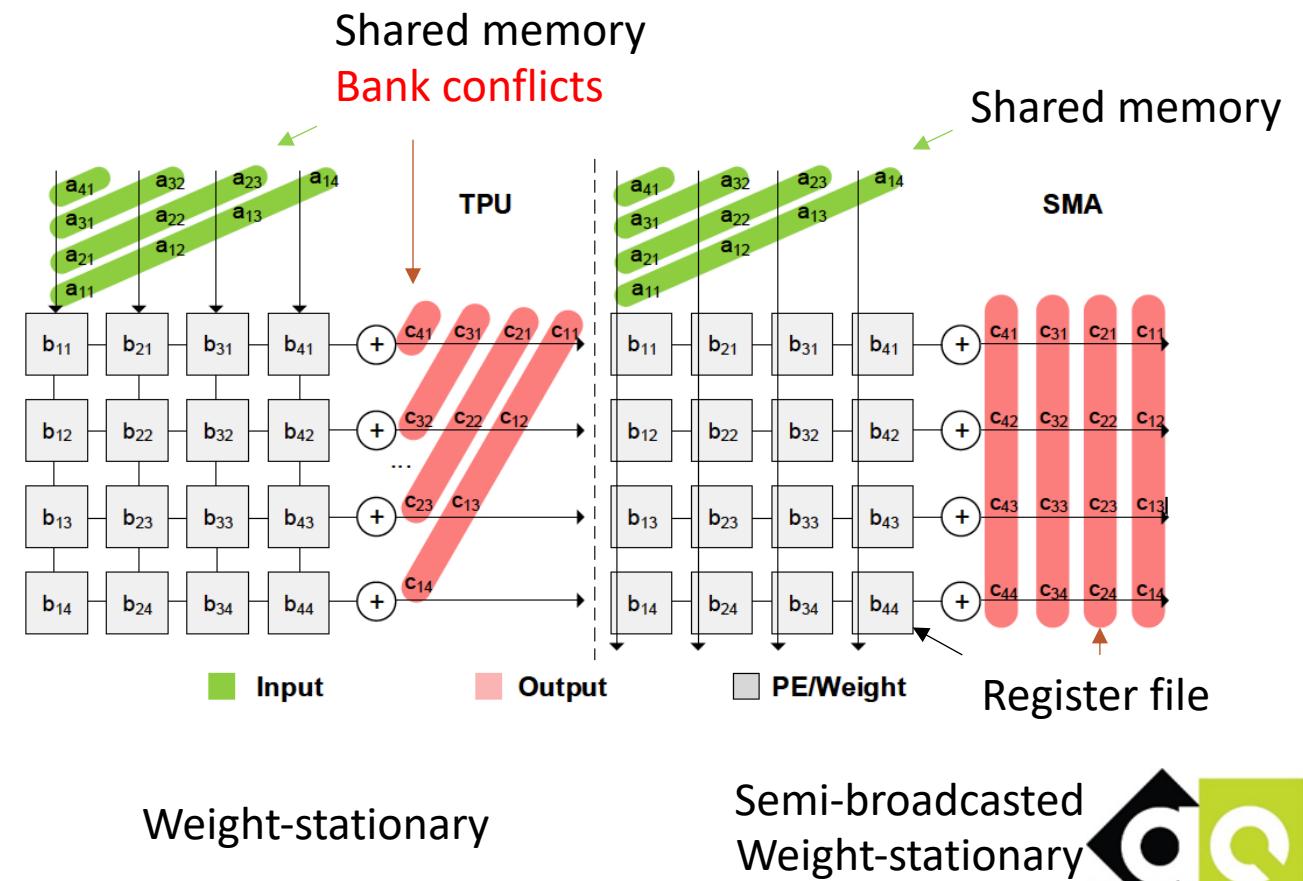


Simultaneous Multi-mode Architecture (SMA)



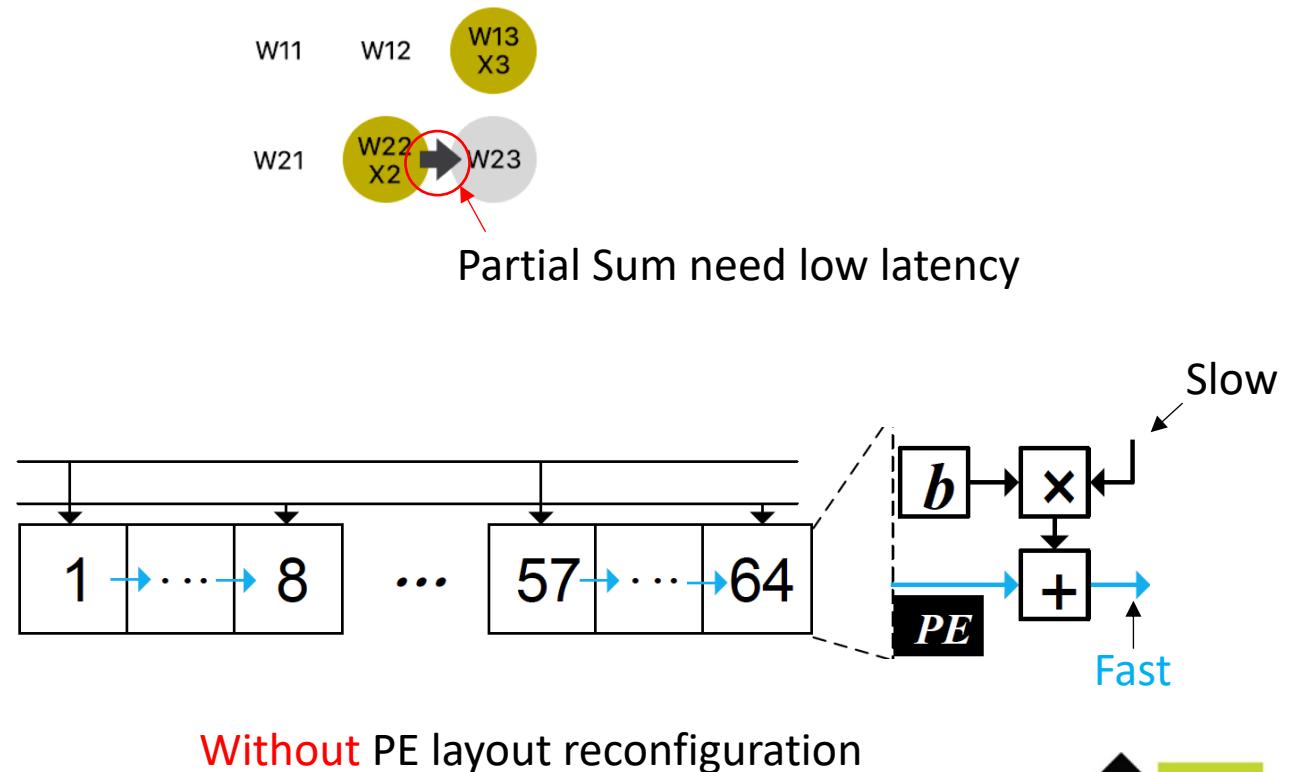
# Massive output scalar accesses

- Semi-broadcast
  - Weight
    - Preload
    - Register file
  - Output
    - Vector access
    - Register file
  - Input
    - Scalar access
    - Shared memory



# Partial Sum Communication

- One-way wires
  - Horizontal neighbor PEs
    - Fast
    - latency-sensitive
    - Negligible overhead
  - Vertical PEs
    - Slow
    - Broadcast, Prefetch
    - latency-insensitivity



# Instruction Control

- A new instruction:
  - LSMA (Load, Store and Multiply-accumulate)
- A systolic controller

Per SM

Input :  $8 \times 2 \times 4 Bytes$

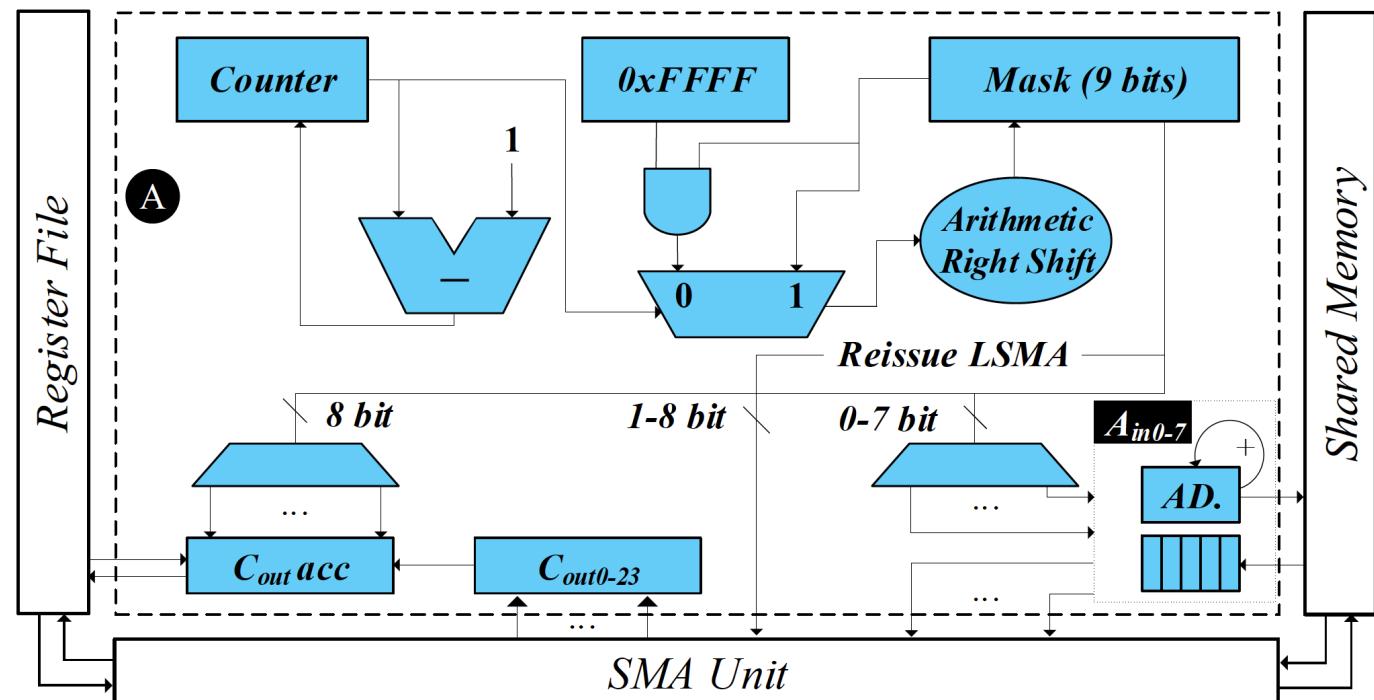
Output :  $24 \times 2 \times 4 Bytes$

Total : 256 Bytes

Area Overhead < 0.1%

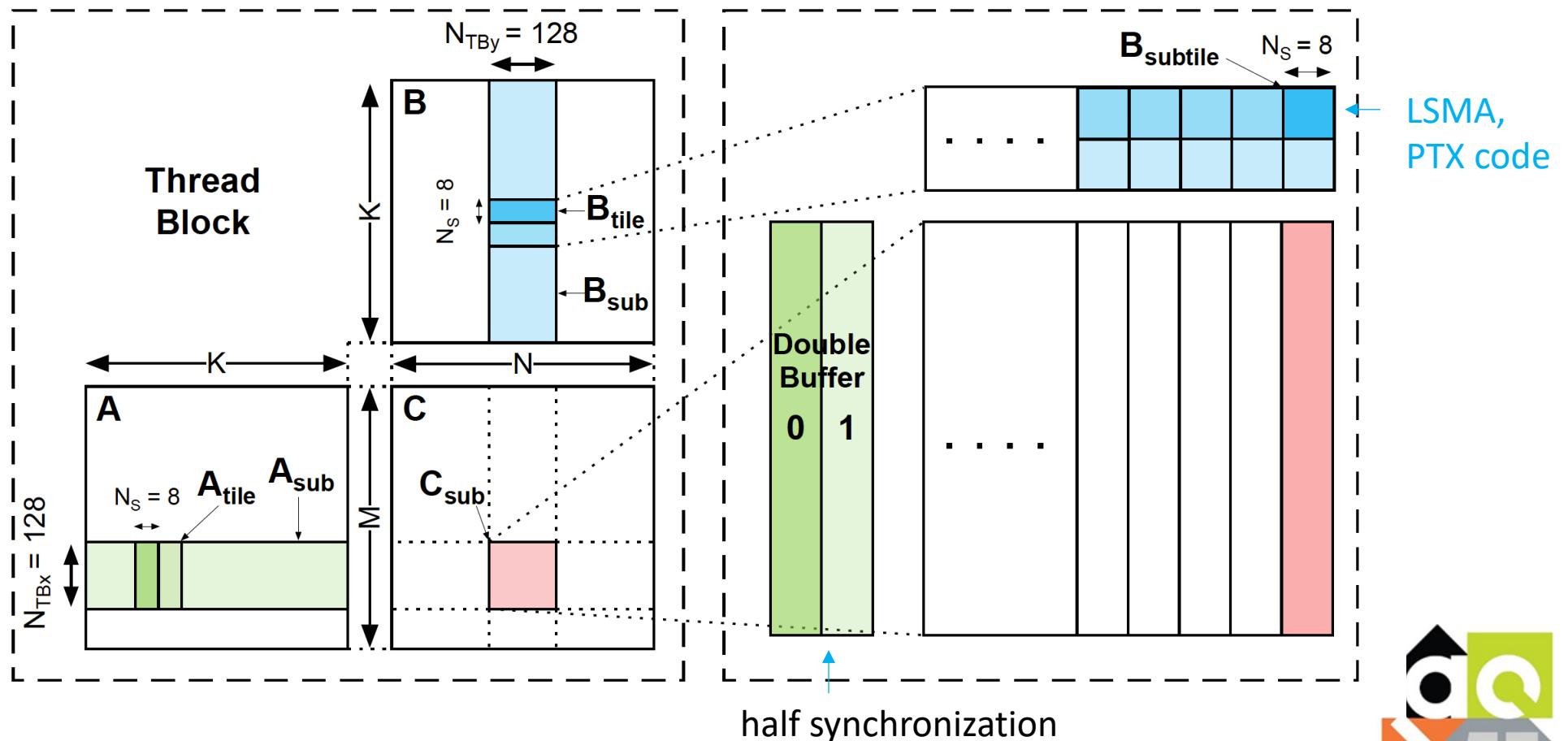
256KB Register file

128KB L1 Cache/Shared memory



$$\text{LSMA } B \Rightarrow C[\text{out}] \leftarrow A[\text{in}] \times B + C[\text{in}]$$

# Software Design: Tiling GEMM



# Outline

- Introduction
- Simultaneous Multi-mode Architecture
- Evaluation

# Evaluation

- Methodology
  - GPGPUsim-4.0
  - GEMM based on CUTLASS

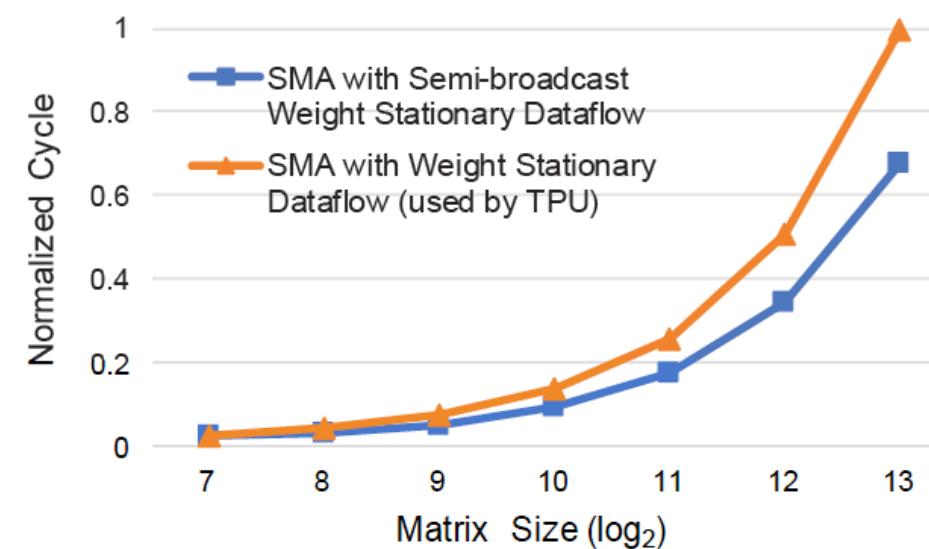
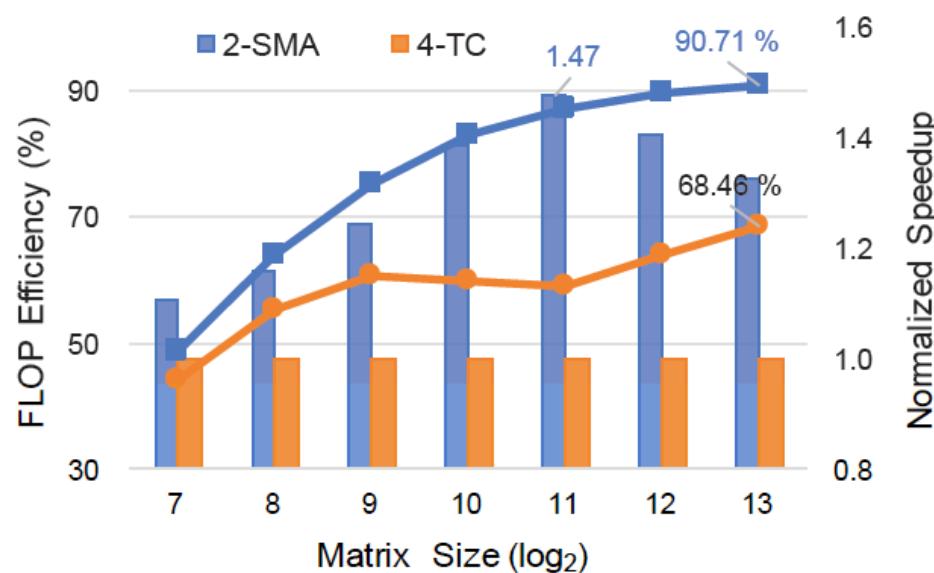
Platform	2-SMA	4-TC	3-SMA
Unit	2 SMA	4 TC	3 SMA
Precision	FP16	FP16	FP16
PE	256	256	384

	GPGPU	SMA
<b>Baseline</b>	Volta	Volta
<b>SMs</b>	80	80
<b>CUDA Core/SM</b>	64 FP32 units	3 $8 \times 8$ SMA unit
<b>Tensor Core/SM</b>	4 (256 FP16 units)	
<b>Shared Memory/SM</b>	32 banks Configurable up to 96KB	32 banks (8 for all SMA units) Configurable up to 96 KB
<b>Register File/SM</b>	256 KB	256 KB



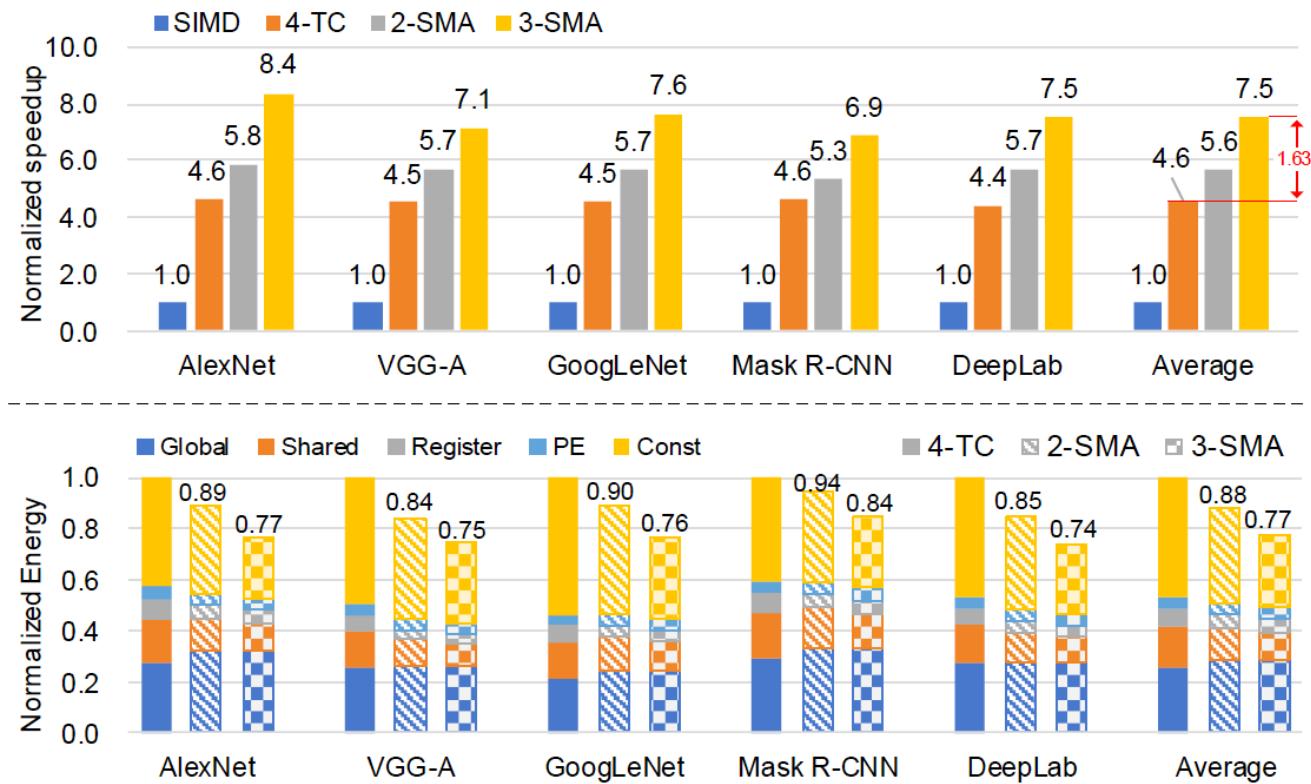
# Iso-FLOP

- Square GEMM
  - 2-SMA efficiency > 90%, 30% higher than 4-TC.
  - SMA (broadcast) 20% - 40% higher than non-broadcast.

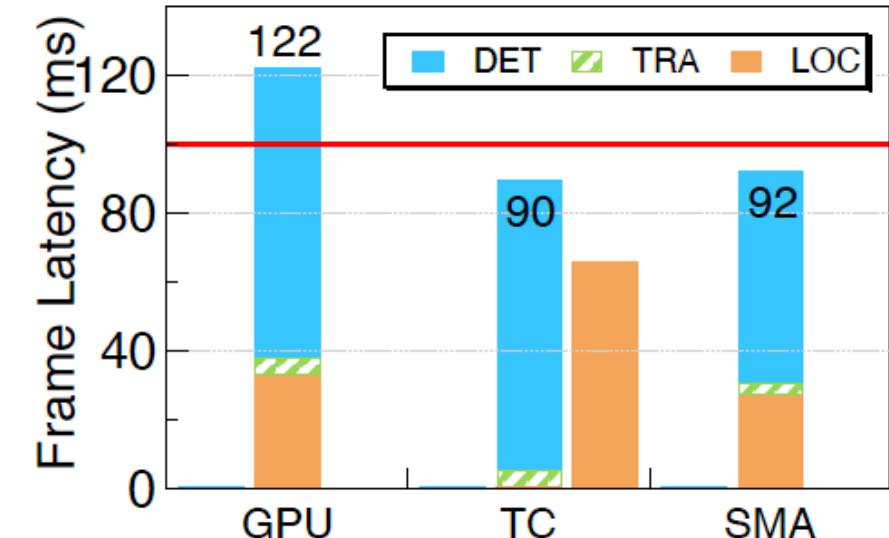
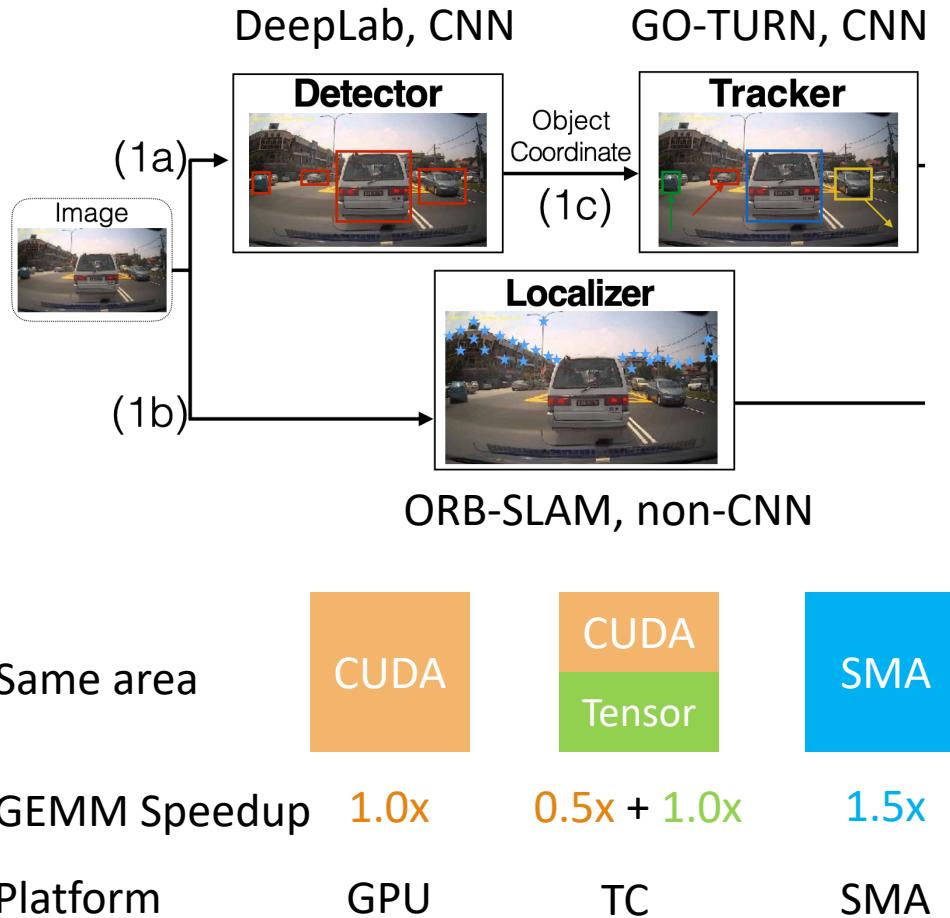


# ISO-Area

- 5 networks
  - 3-SMA **63% faster, 23% less energy** than 4-TC on average.

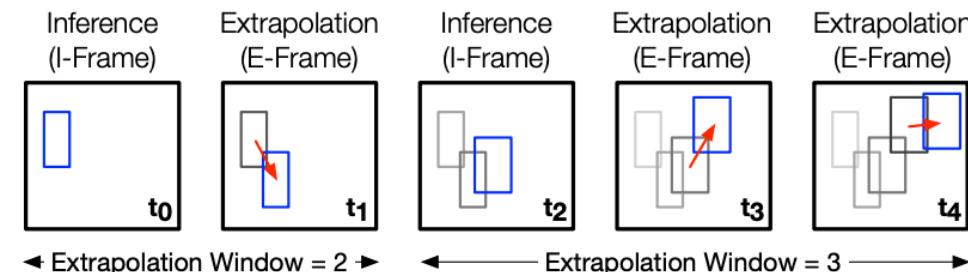
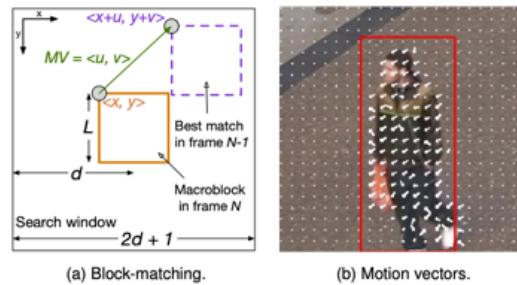


# End-to-end application (autopilot)

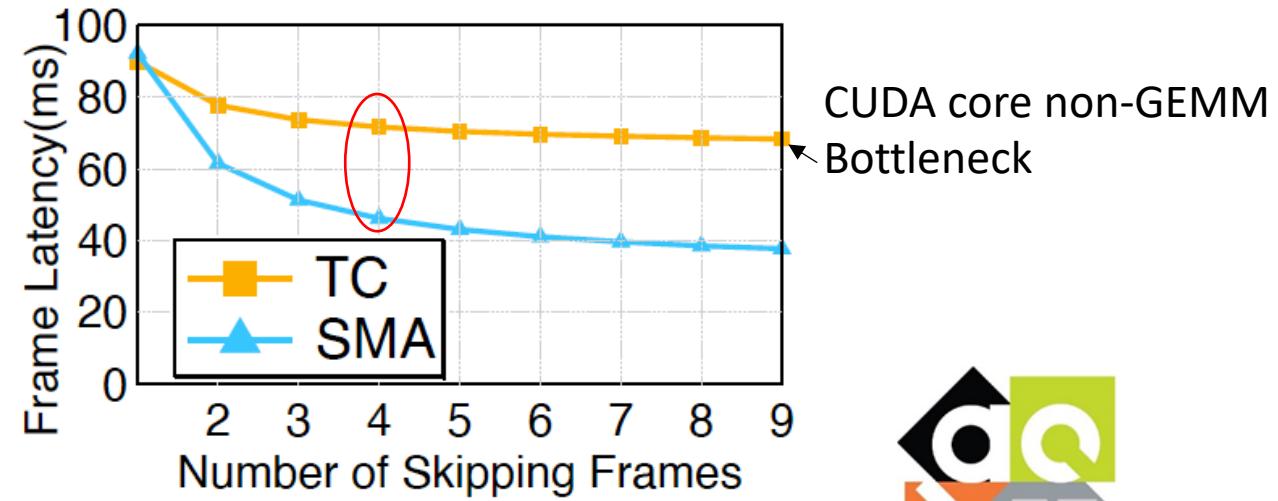
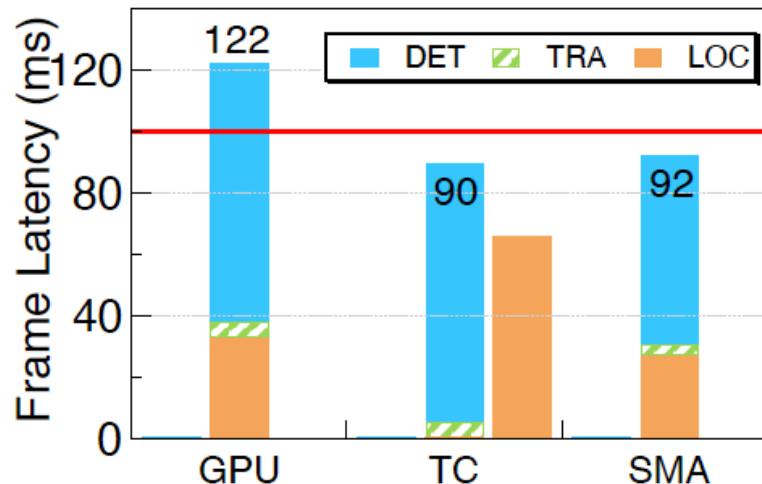


[Shih-Chieh Lin, etc., ASPLOS'18]

# End-to-end application (autopilot)



[Euphrates, ISCA'18]



$N = 4$ , SMA Latency Reduction 50%, More Flexibility

# Summary

- Hardware
  - Parallelism similarity
  - Memory and communication
  - Systolic controller
- Software
  - Tiling GEMM
- Evaluation
  - Efficiency
  - Flexibility



# Questions

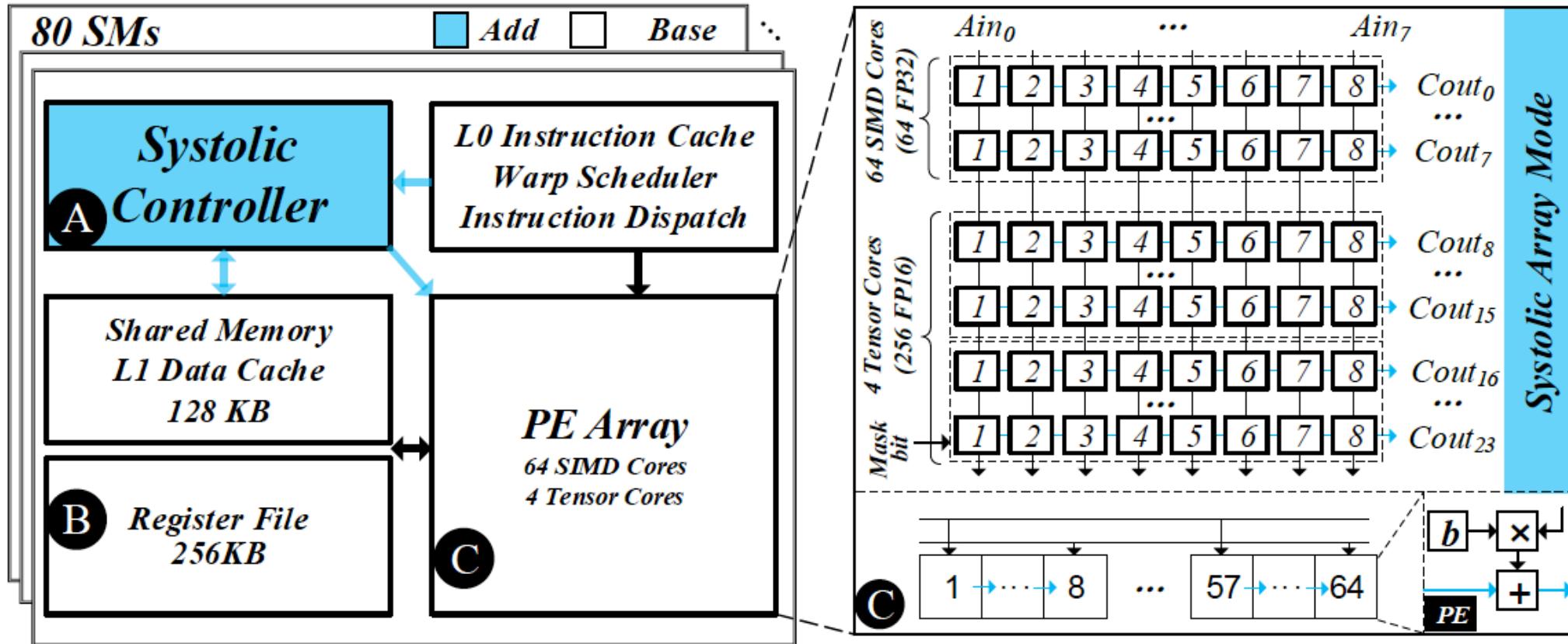
Thank you!



# Backup Slides



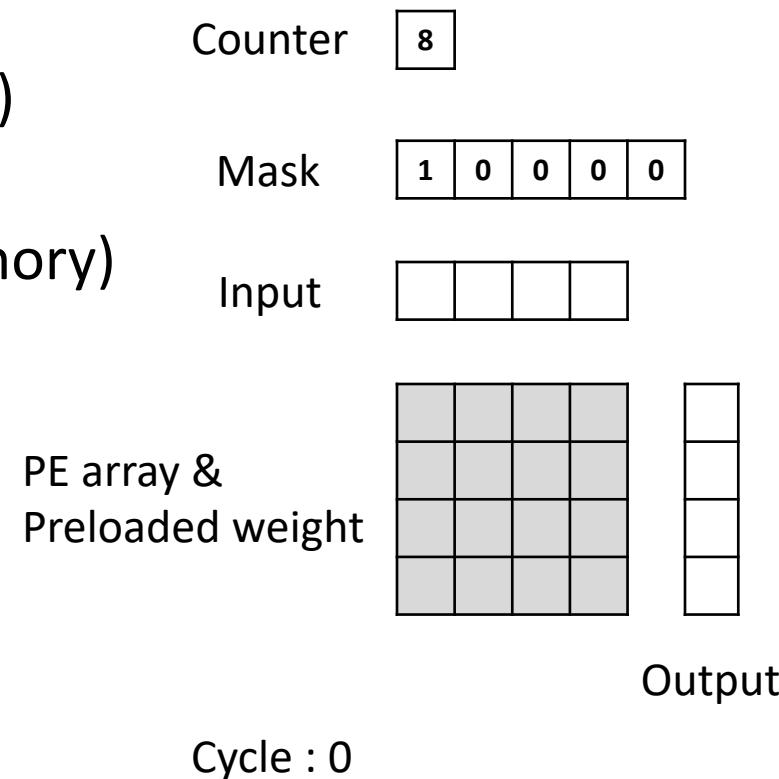
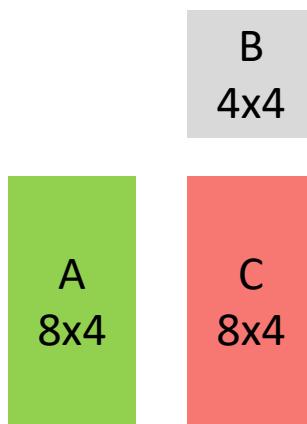
# SMA Hardware Design



Simultaneous Multi-mode Architecture (SMA)

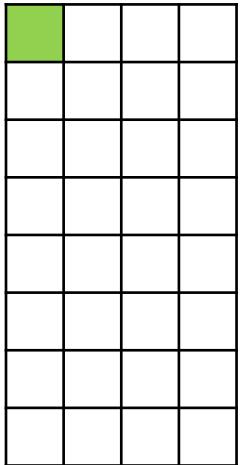
# LSMA execution

- 5-bit mask for **4x4** array
- Counter (input row number)
- Preload weight
- Prefetch input (Shared memory)
- Execute
- Store output (Register file)



# Cycle : 1

Discontinuous



Prefetch input

Counter    **8**

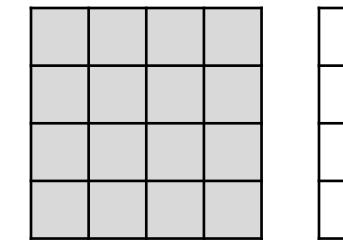
Mask    

1	0	0	0	0
---	---	---	---	---

Input    

Green	White	White	White	White
-------	-------	-------	-------	-------

PE array &  
Preloaded weight

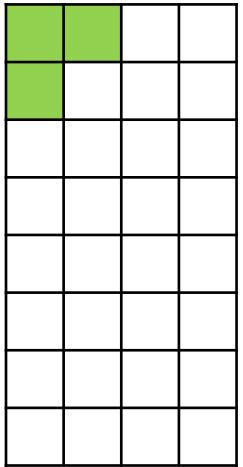


Output

Cycle : 1



# Cycle : 2



Execute

PE array &  
Preloaded weight

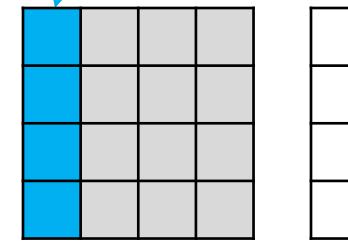
Cycle : 2

Counter 7

Mask 1 1 0 0 0

Input 

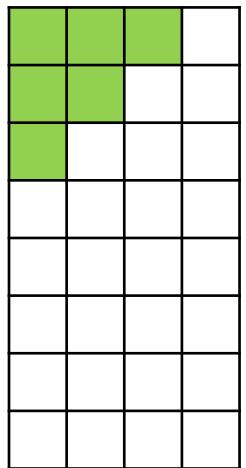
A horizontal vector of length 5, with the first two elements colored green and the remaining three elements white. An arrow points from this vector to the PE array.



Output



# Cycle : 3

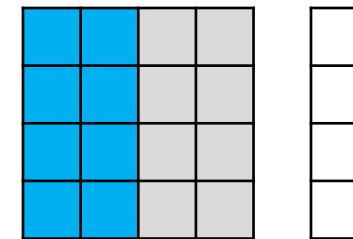


Counter 6

Mask 1 1 1 0 0

Input | Green | Green | Green | White |

PE array &  
Preloaded weight

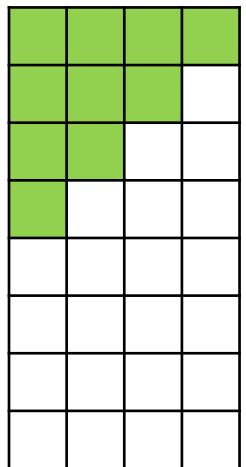


Output

Cycle : 3



# Cycle : 4

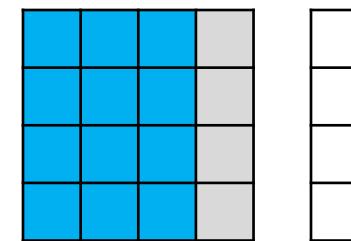


Counter 5

Mask 1 1 1 1 0

Input   |  |  |  |

PE array &  
Preloaded weight

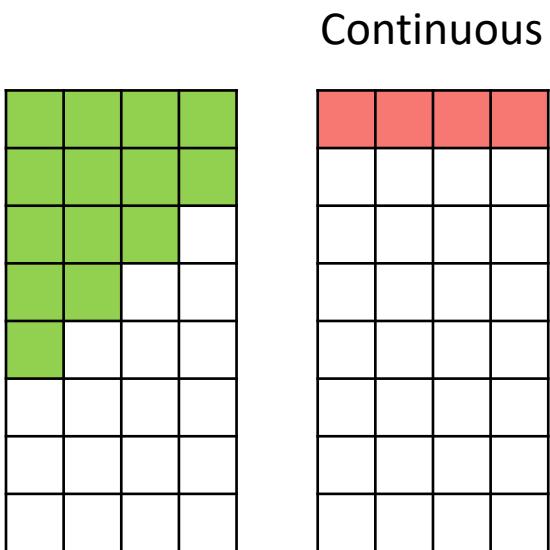


Output

Cycle : 4

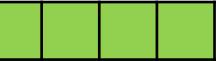


# Cycle : 5

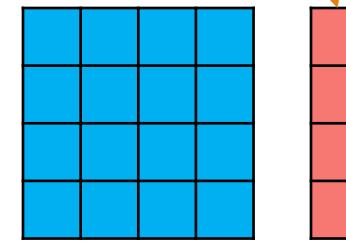


Counter 4

Mask 1 1 1 1 1

Input 

PE array &  
Preloaded weight



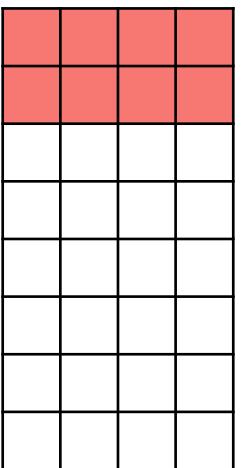
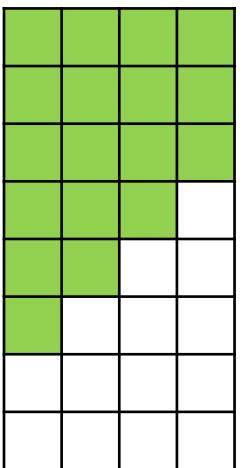
Store output

Output

Cycle : 5



# Cycle : 6



Counter

3

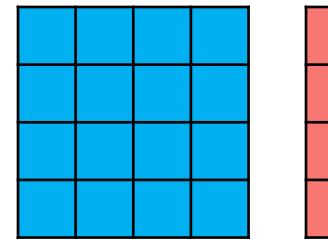
Mask

1 1 1 1 1

Input

A horizontal vector consisting of five green cells.

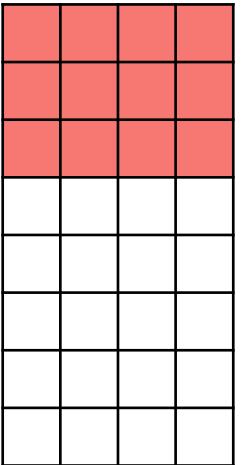
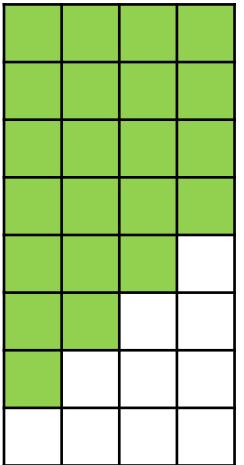
PE array &  
Preloaded weight



Output

Cycle : 6

# Cycle : 7

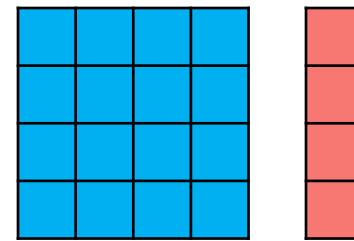


Counter 2

Mask 1 1 1 1 1

Input   |  |  |  |

PE array &  
Preloaded weight

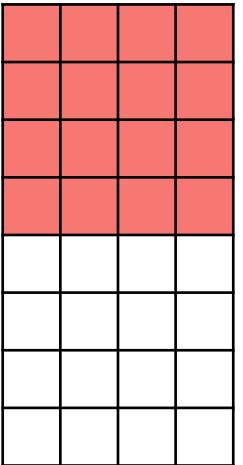
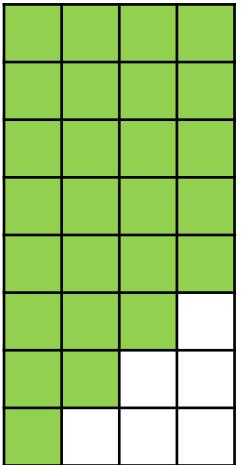


Output

Cycle : 7



# Cycle : 8

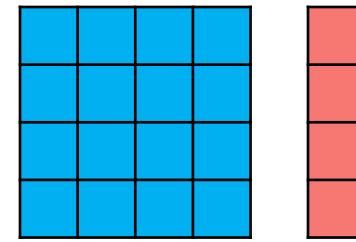


Counter 1

Mask 1 1 1 1 1

Input   |  |  |  |

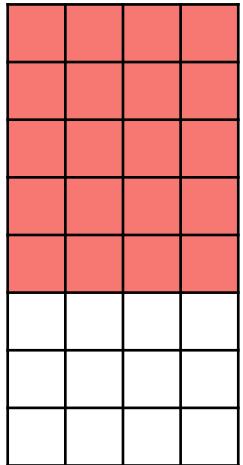
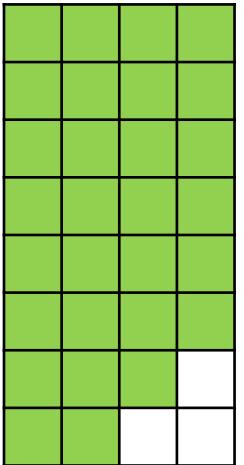
PE array &  
Preloaded weight



Output

Cycle : 8

# Cycle : 9



Counter

0

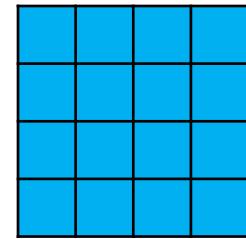
Mask

0 1 1 1 1

Input

A horizontal vector consisting of one white square followed by four green squares.

PE array &  
Preloaded weight

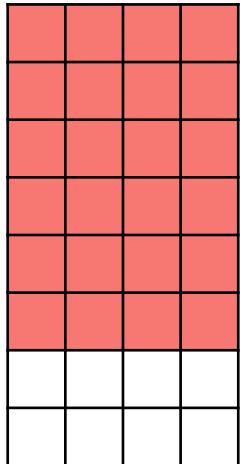
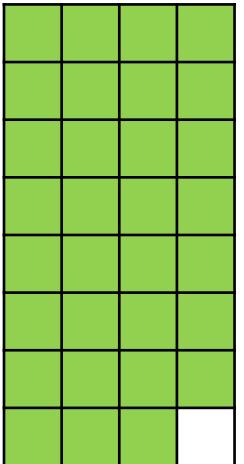


Output

Cycle : 9



# Cycle : 10

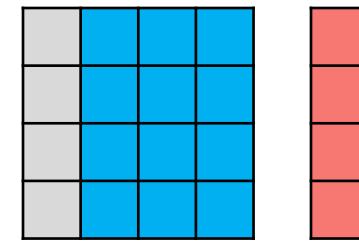


Counter 0

Mask 0 0 1 1 1

Input    |  |  |  |  |

PE array &  
Preloaded weight

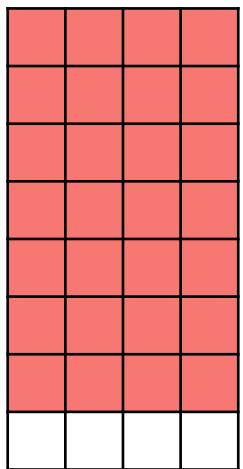
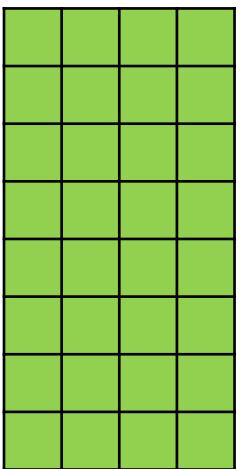


Output

Cycle : 10



# Cycle : 11

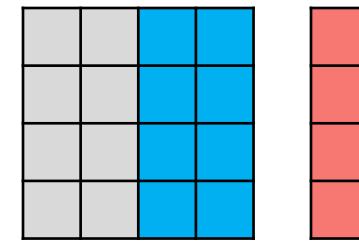


Counter 0

Mask 0 0 0 1 1

Input    |  |  |  |  
      |      |  
          |      |  
            |      |  
              |      |  
                |      |

PE array &  
Preloaded weight

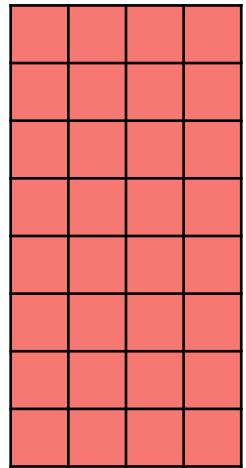


Output

Cycle : 11



# Cycle : 12

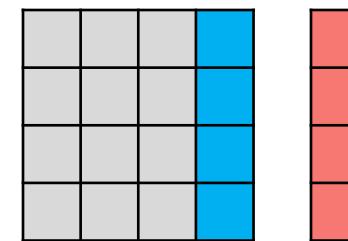


Counter 0

Mask 0 0 0 0 1

Input   |  |  |  |

PE array &  
Preloaded weight



Output

Cycle : 12

# Software Design: Tiling GEMM

