

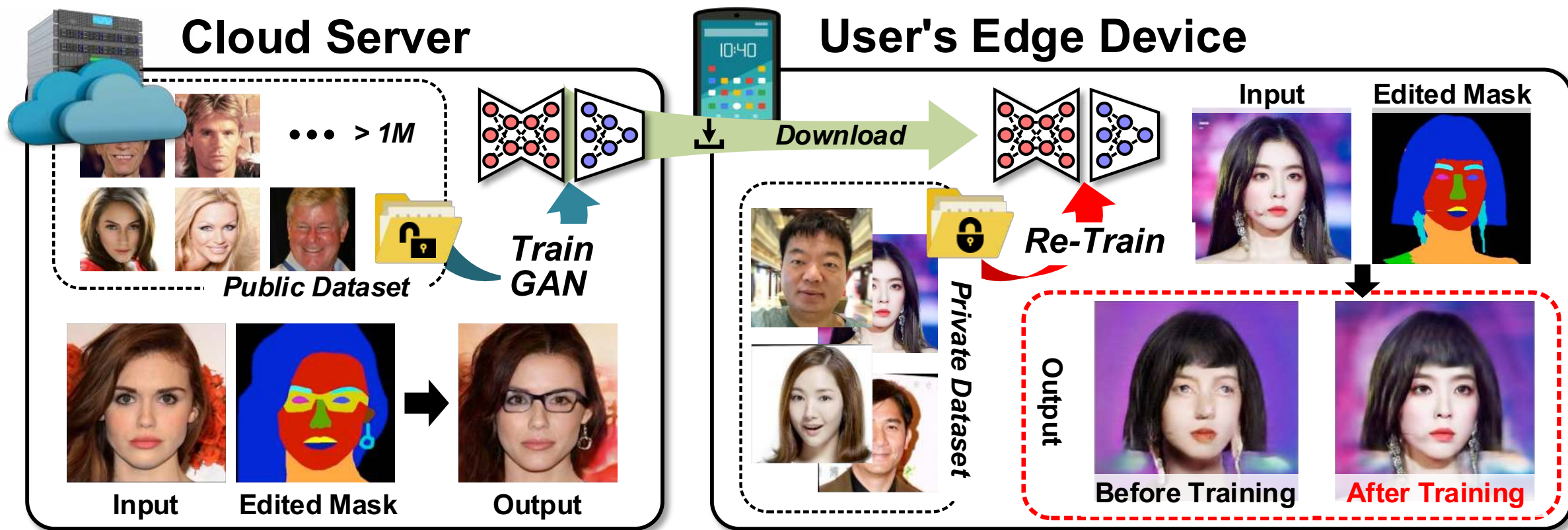
# **GANPU: A Versatile Many-Core Processor for Training GAN on Mobile Devices with Speculative Dual-Sparsity Exploitation**

**Sanghoon Kang**, Donghyeon Han, Juhyoung Lee,  
Dongseok Im, Sangyeob Kim, Soyeon Kim, Junha Ryu, and Hoi-Jun Yoo

**Semiconductor System Lab.  
School of EE, KAIST**

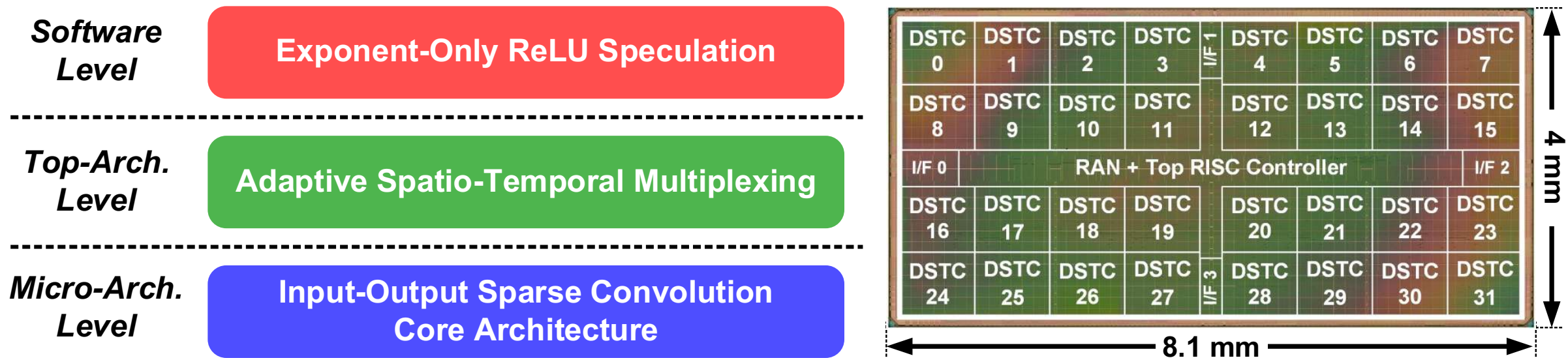
# GAN with On-Device Training

- GAN Applications on Mobile Need **On-Device Training**
  - To restore performance from distortion due to **data discrepancy**



# GANPU: A GAN Training Processor

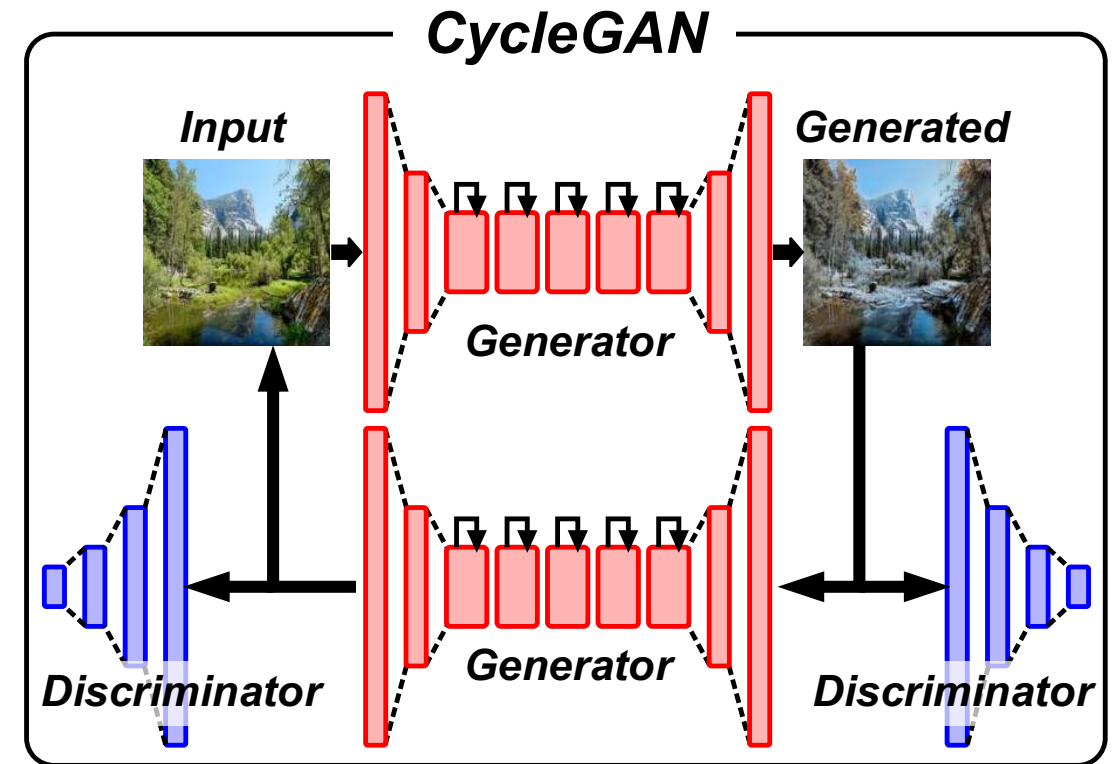
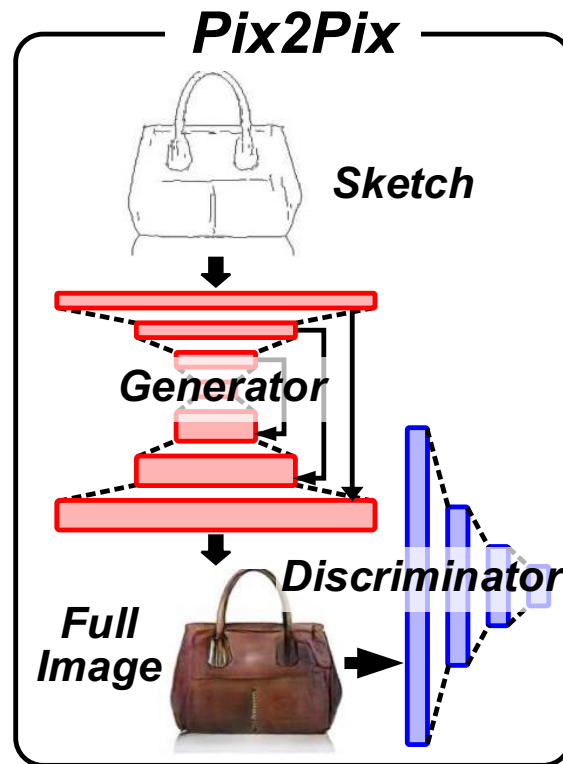
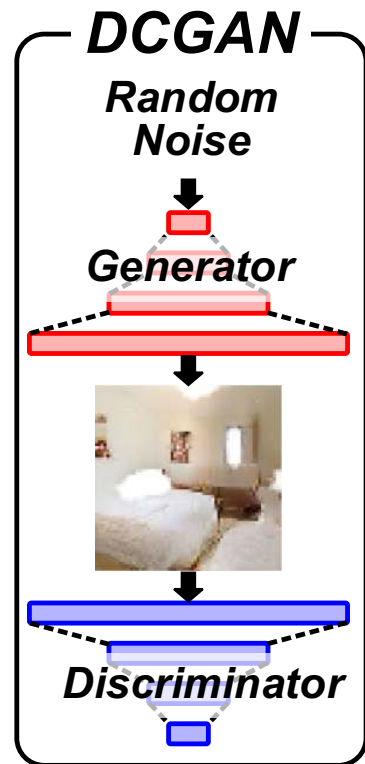
- A 135 TFLOPS/W Multi-DNN Accelerator for Mobile Devices



- Heterogeneity of multi-DNNs → Adaptive core & B/W allocation
- Large # of MAC ops in GANs → HW architecture & SW algorithm for dual sparsity exploitation

# Multi-DNN Architecture of GAN

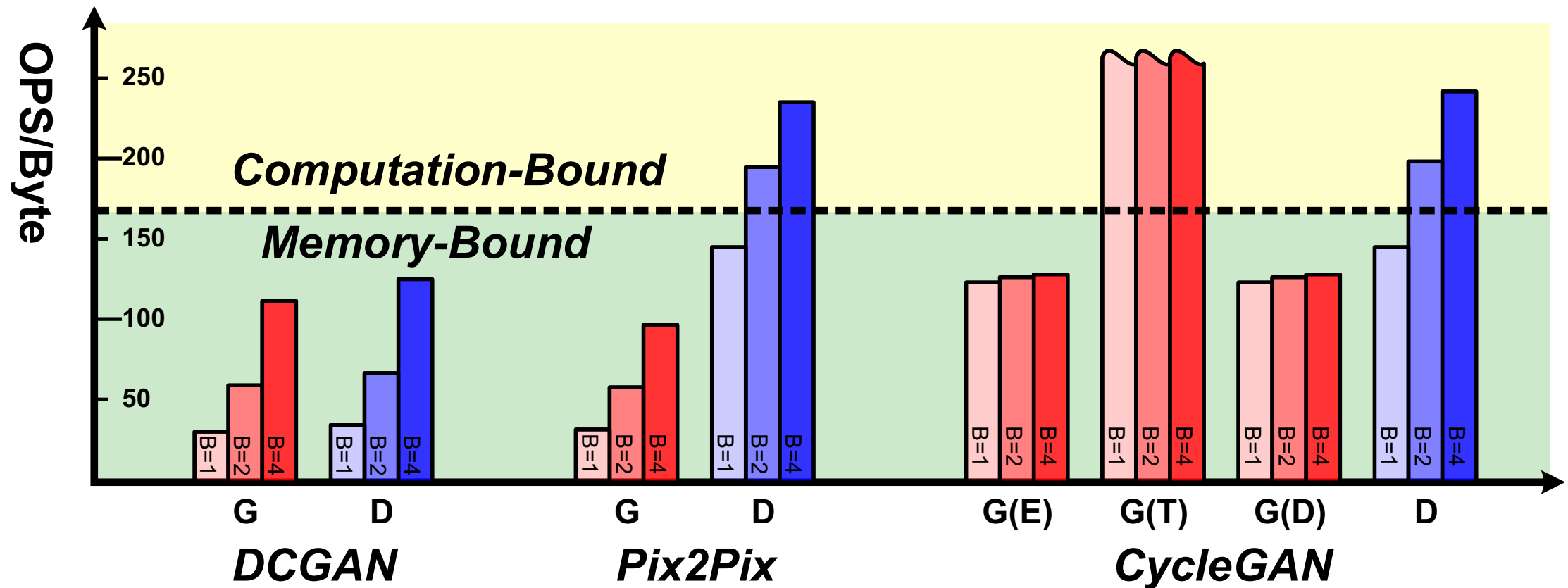
- Various Multi-DNN Architectures



**Complex Architectures w/ Many DNNs**

# Multi-DNN Architecture of GAN

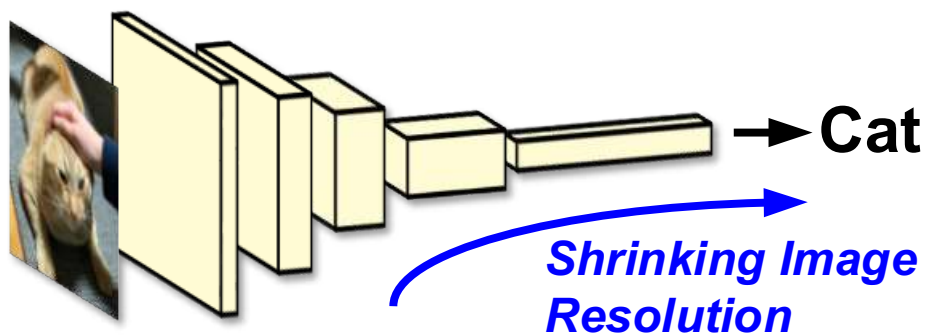
- **Co-Existing Memory Bottleneck & Computation Bottleneck**
  - **Underutilization** of resources (Core & B/W)



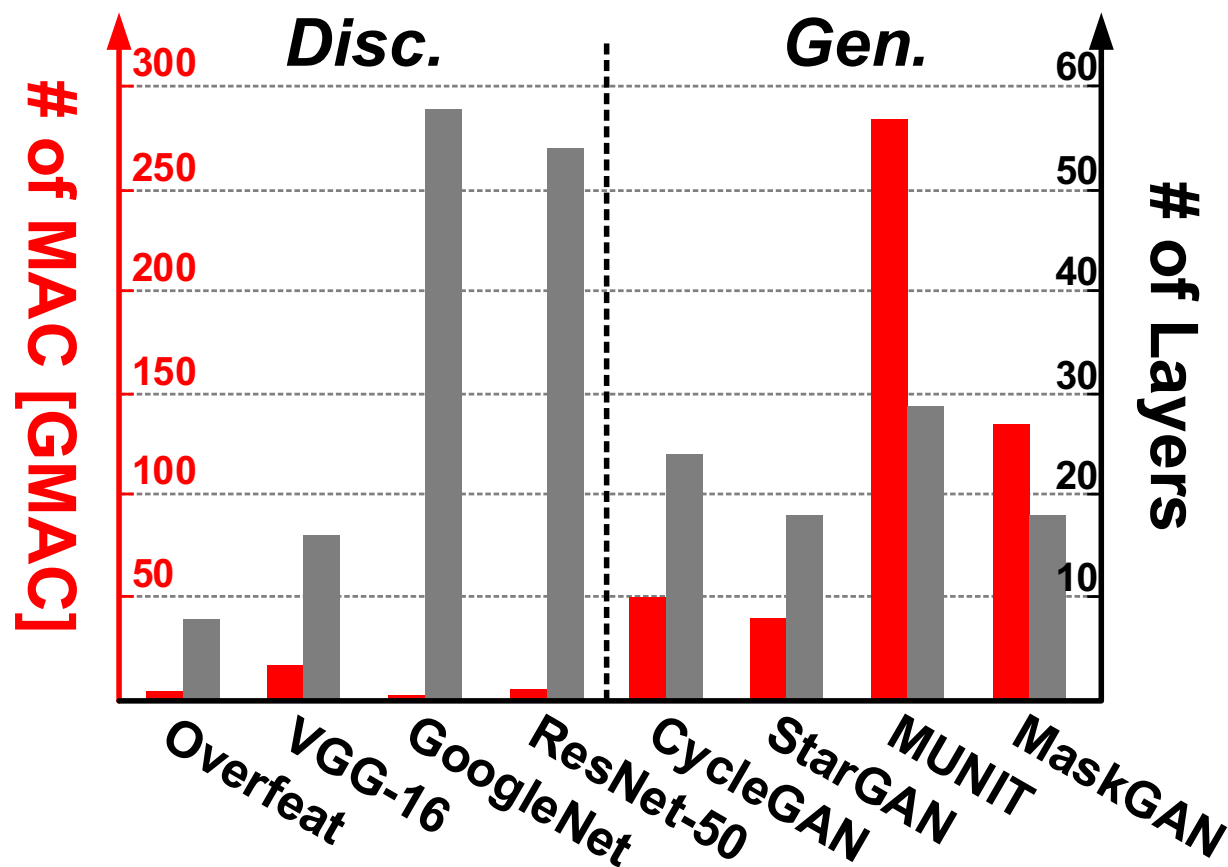
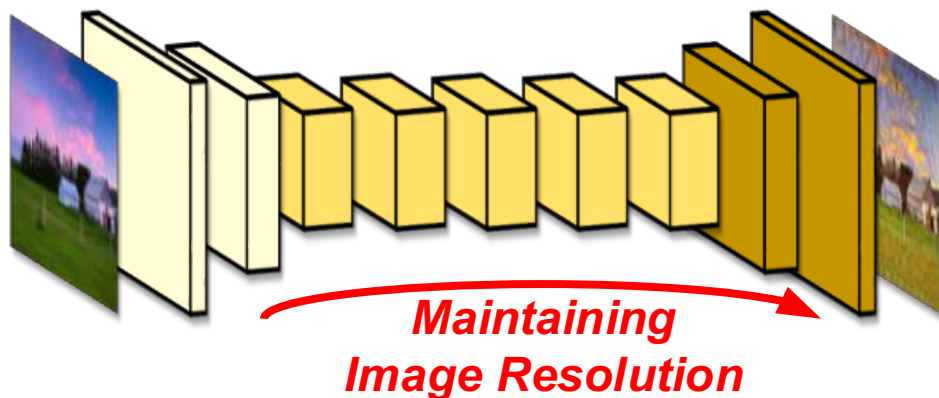
# High Computational Requirement

- More Computation from Larger Intermediate Activations

## Discriminative Models

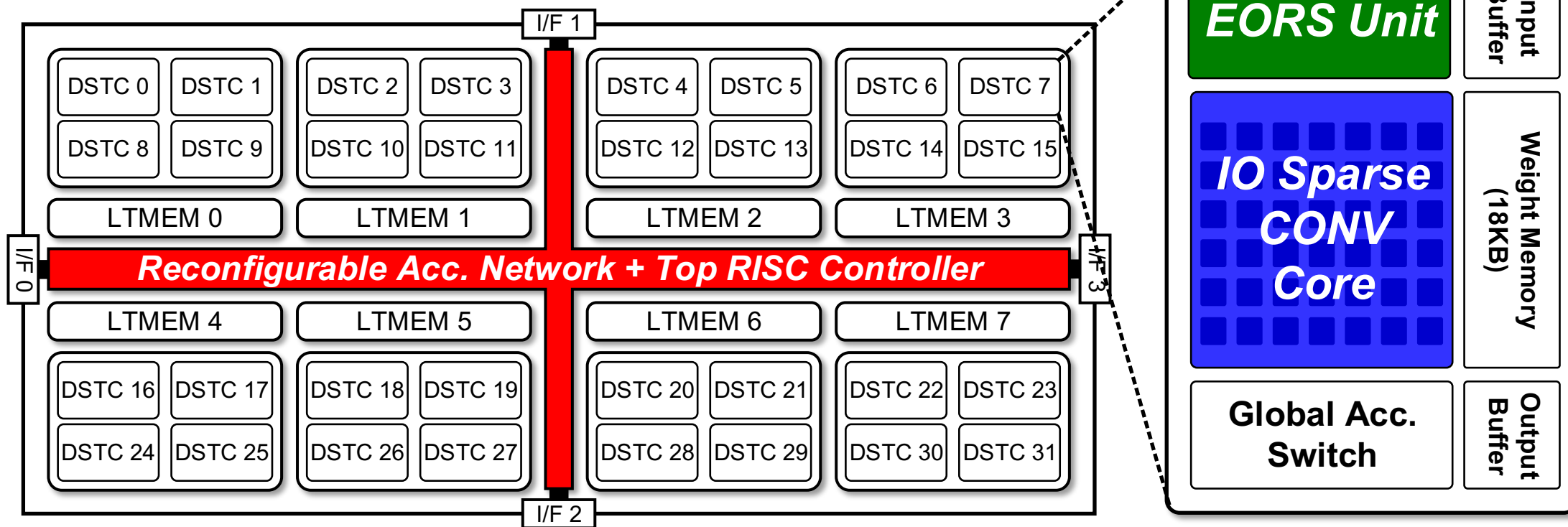


## Generative Models



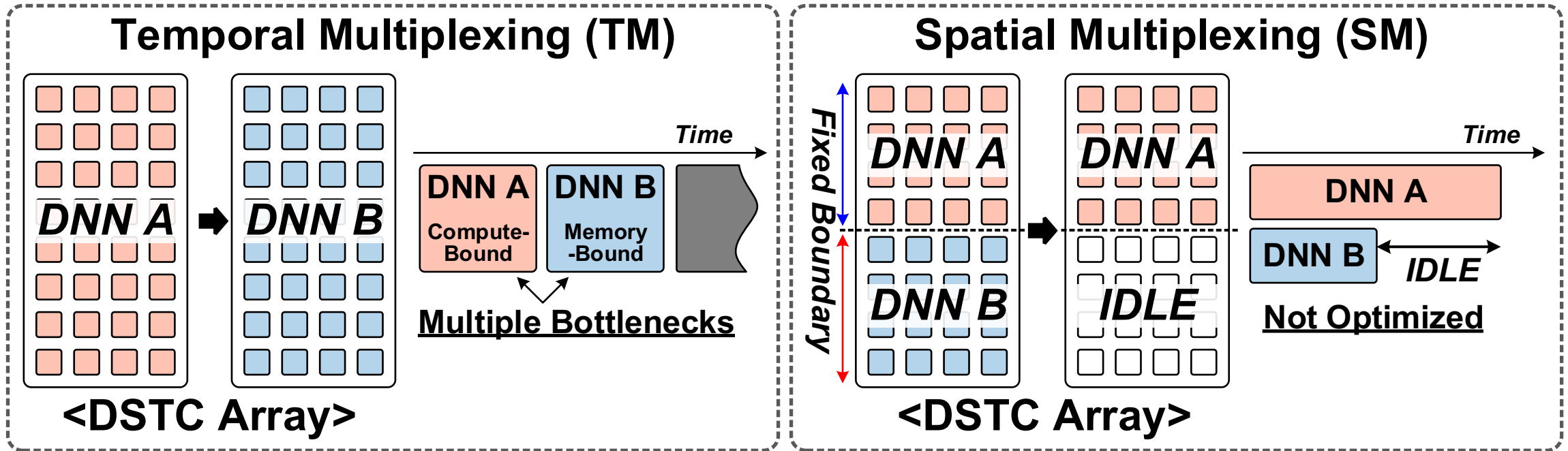
# Overall Architecture of GANPU

1. *Adaptive Spatio-Temporal Multiplexing (ASTM)*
2. *Input-Output Sparse Convolution Core Architecture (IOSC)*
3. *Exponent-Only ReLU Speculation (EORS)*



# Multi-DNN Workload Allocation

## Conventional Workload Multiplexing Methods

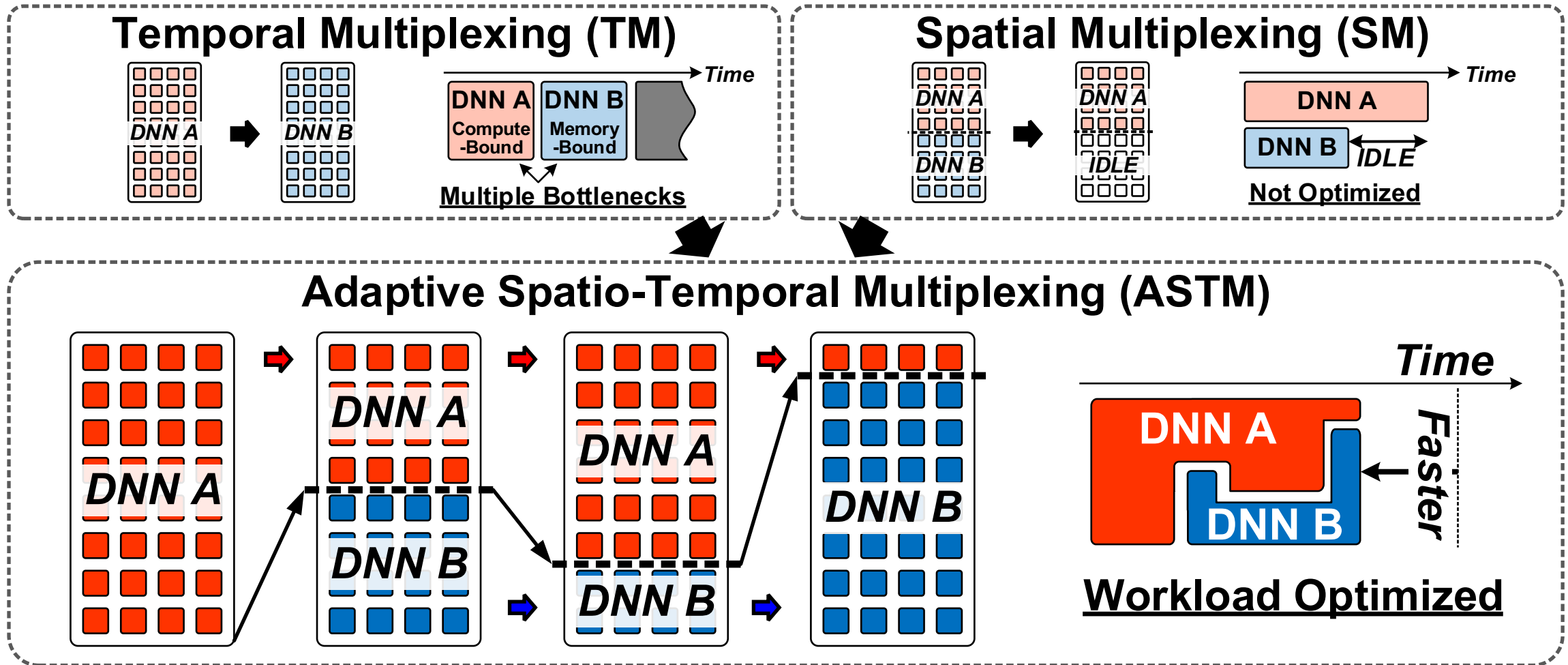


- Varying computation & B/W workloads → Multiple bottlenecks
- Fixed DSTC allocation → No adaptation for optimal processing



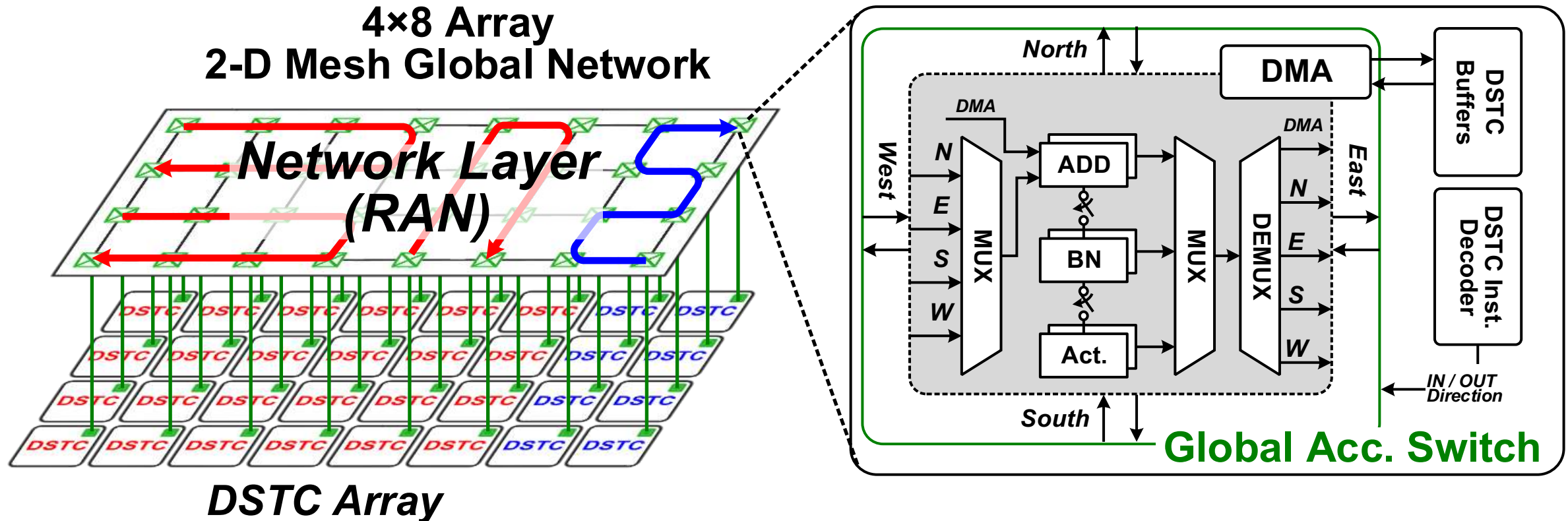
# Multi-DNN Workload Allocation

- Adaptive Spatio-Temporal Multiplexing (ASTM)



# Reconfigurable Accumulation Network

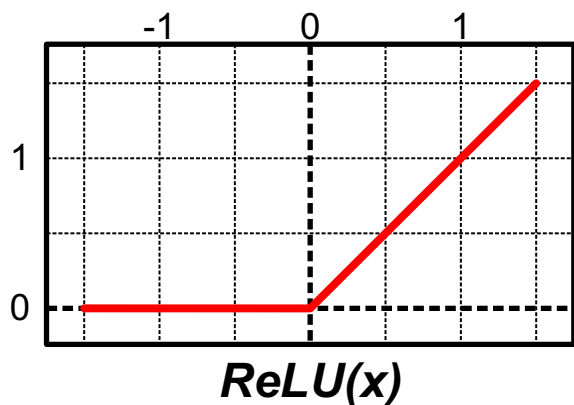
- Supporting Versatile Accumulation Patterns
  - Stream partial sums across DSTCs through **global accumulation switch**
  - 14% throughput increase** in Pix2Pix



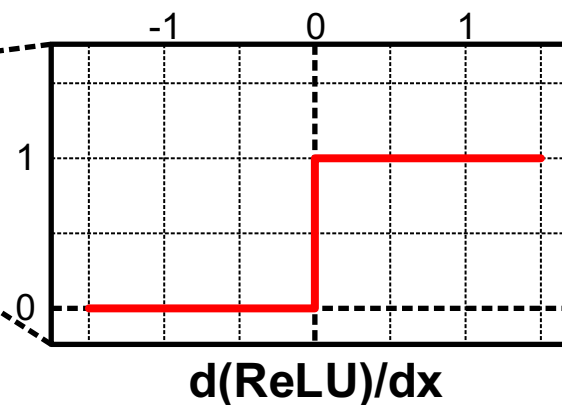
# Sparse Activation in DNN Training

- Activatio Sparsity @ FF & BP Stages Due to **ReLU**

Feed Forward (FF)



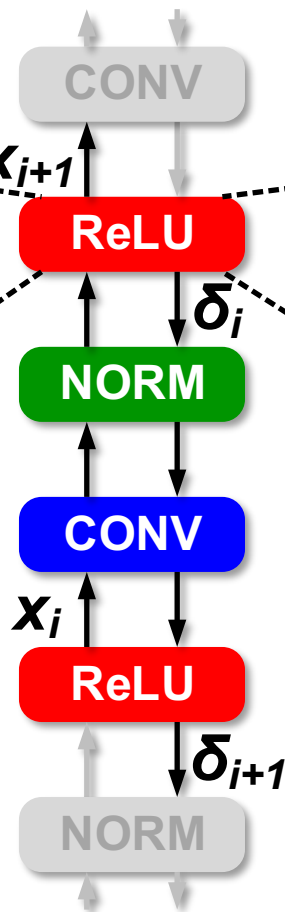
Backward Propagation (BP)



$$x_{i+1} = ReLU(NORM(CONV(x_i)))$$

OA: Zeros                      IA: Zeros

☹️ OA Sparsity Unknown



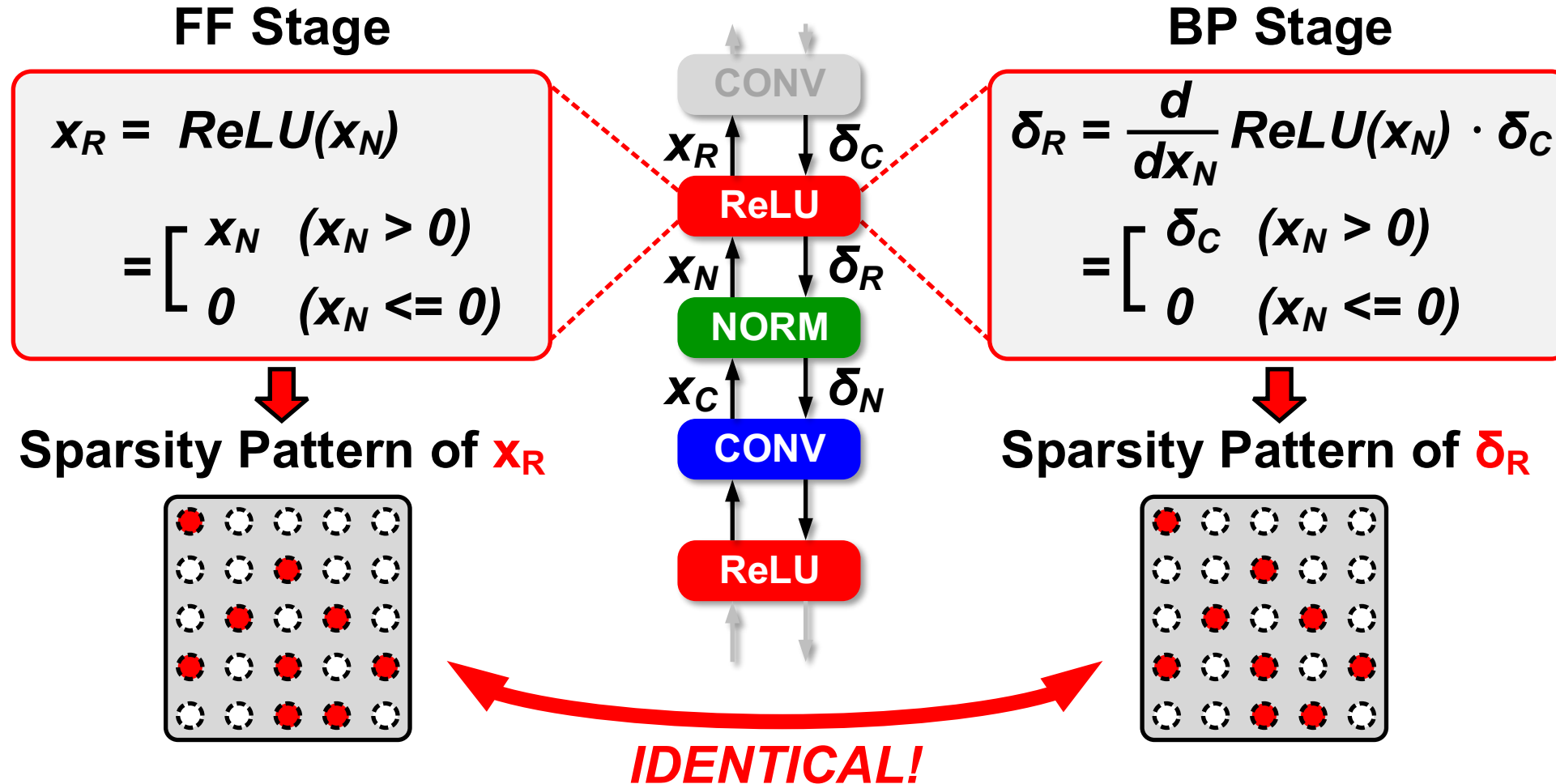
$$\delta_{i+1} = ReLU_{BP}(CONV_{BP}(NORM_{BP}(\delta_i)))$$

OA: Zeros                      IA: No Zeros

☺️ OA Sparsity Known

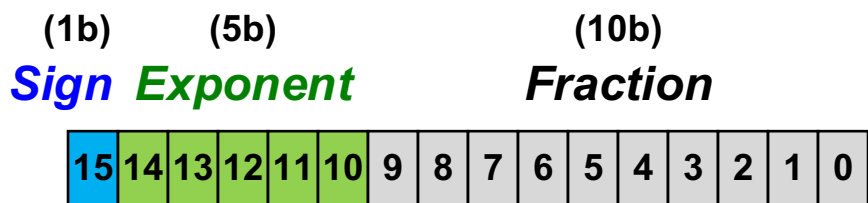
# OA Sparsity @ BP Stage

- **Identical Pattern** → OA Sparsity Pattern Pre-determined



# OA Sparsity @ FF Stage

## Exponent-Only ReLU Speculation (EORS)

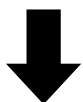


< FP16 (Half Precision) Format >

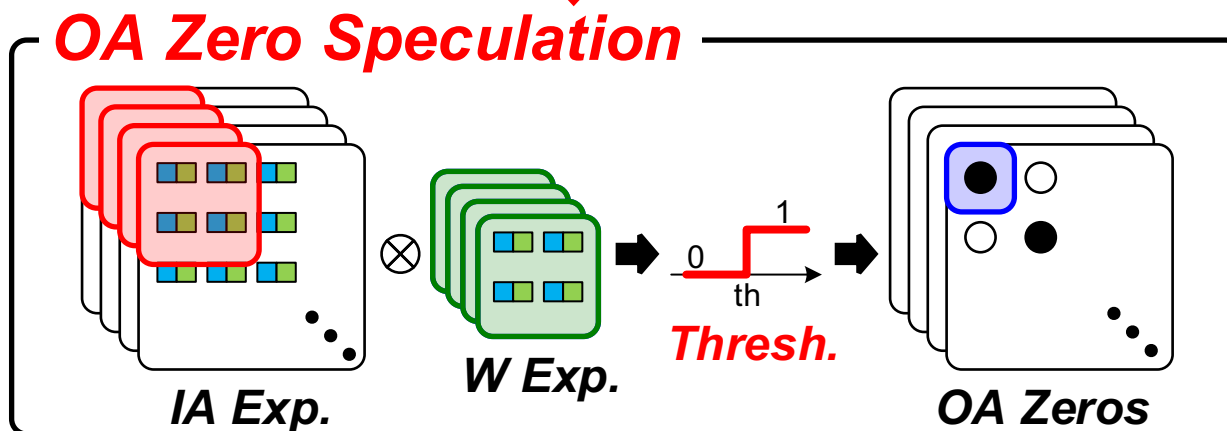
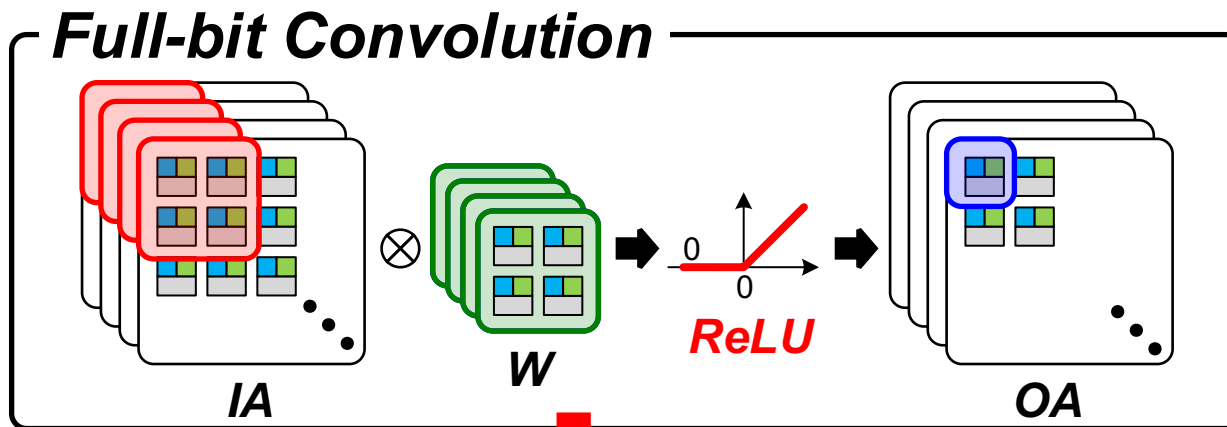
$$\text{Data} = (-1)^{\text{Sign}} \times 2^{\text{Exponent}-15} \times 1.\text{Fraction}_2$$

(less important)

**Scale of IA & W**

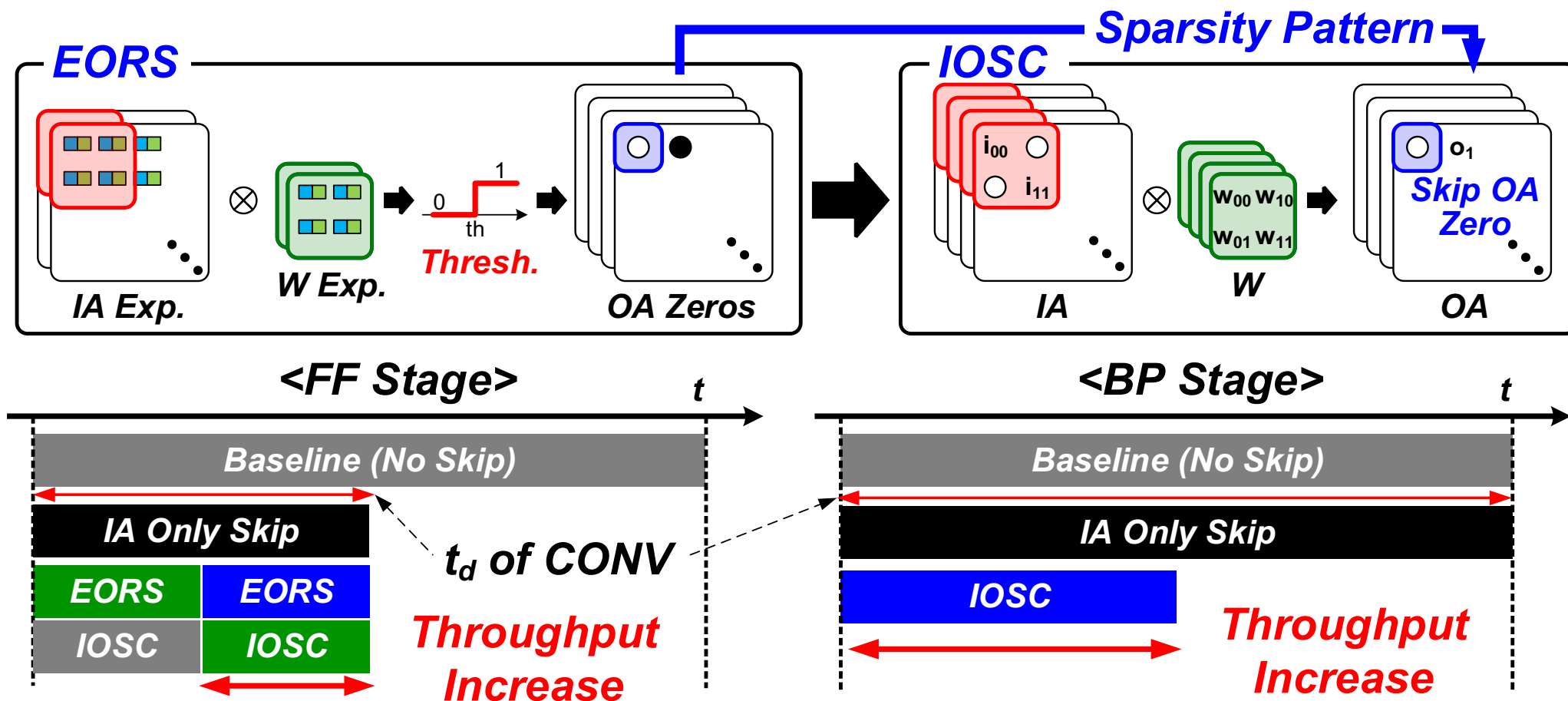


**Simple to Estimate  
OA Sparsity Pattern**



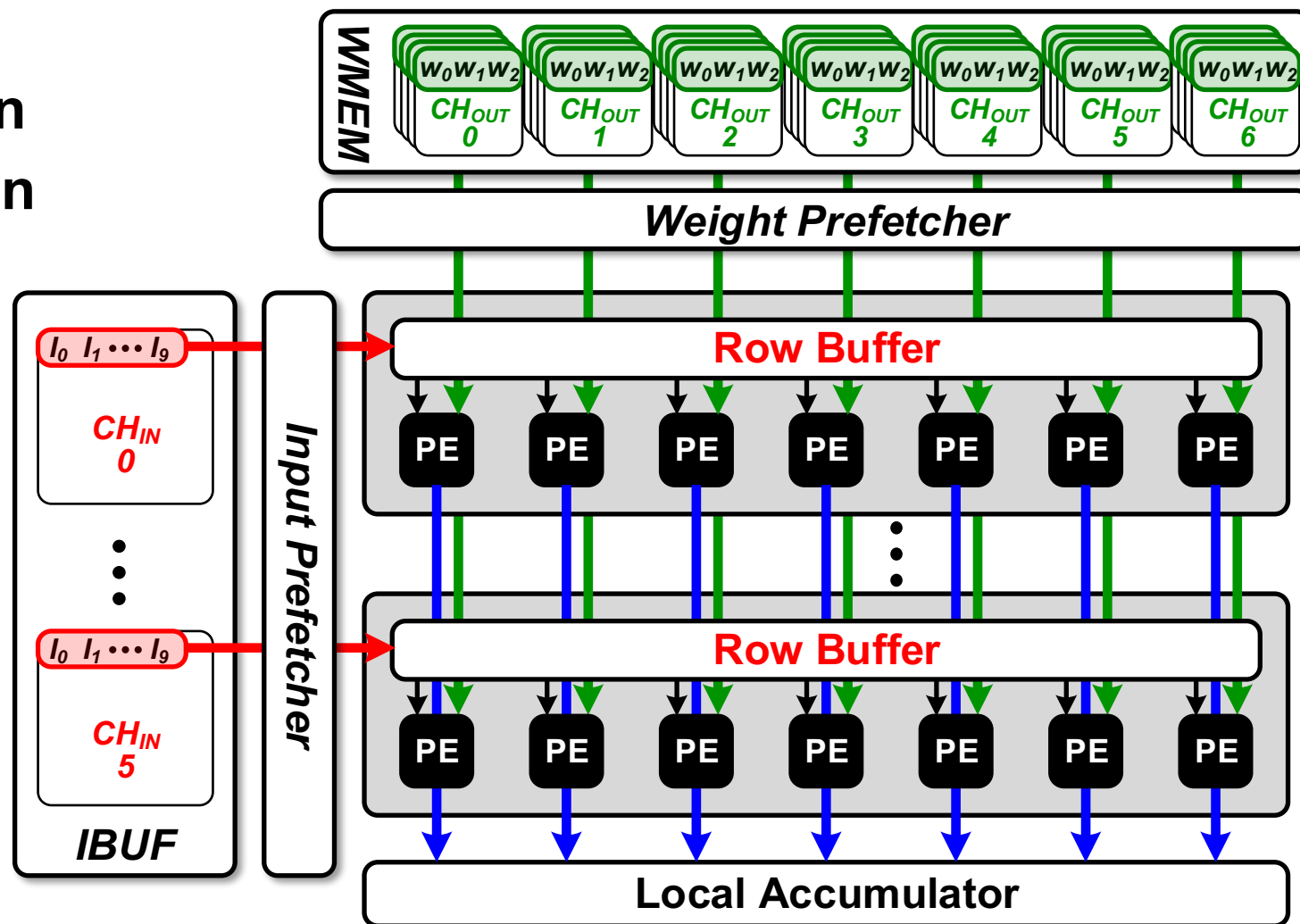
# Dual Sparsity Exploitation with EORS

- Higher Throughput @ Both FF & BP Through IOSC + EORS



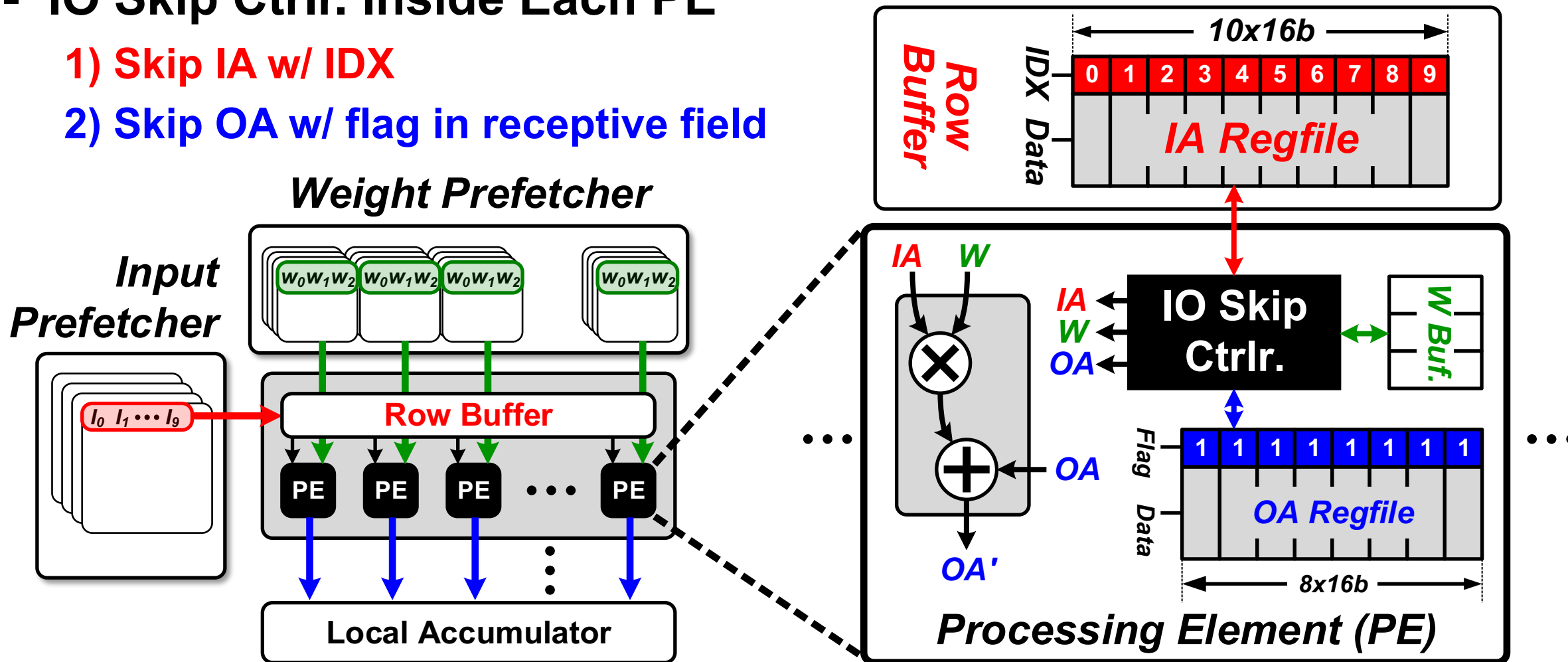
# IOSC Architecture

- **6×7 PE Array in IOSC**
  - **CH<sub>IN</sub>**: Row-wise division
  - **CH<sub>OUT</sub>**: Col-wise division
  - Every PE allocated w/ different CH!
- **Hierarchical Dataflow**
  - Row stationary (Inside row buffer)
  - Input stationary (Inside PE)



# IOSC PE Architecture

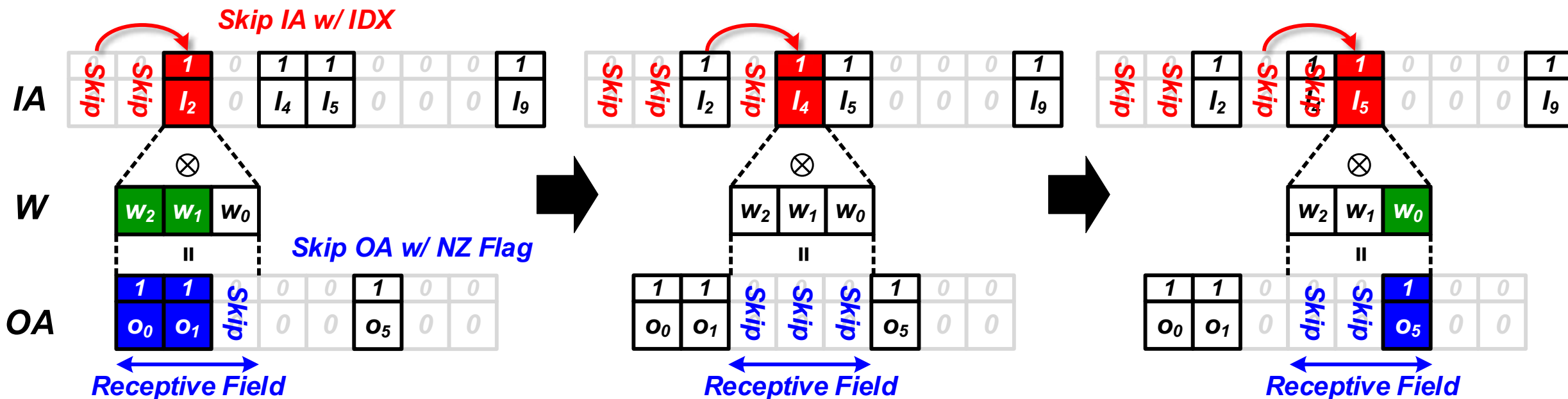
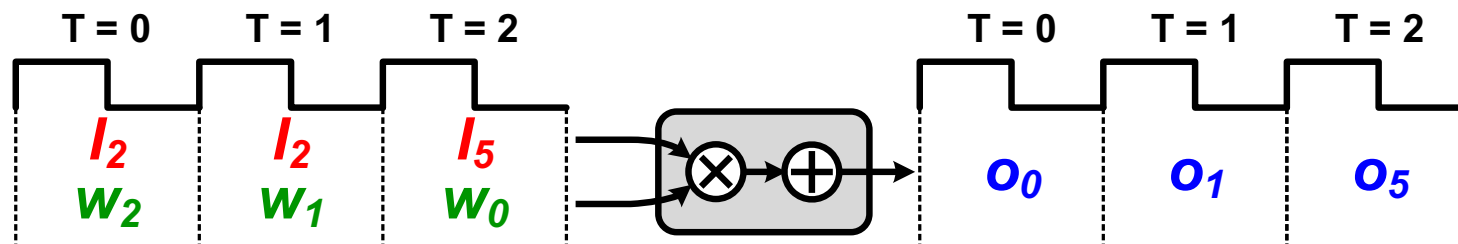
- IO Skip Ctrlr. Inside Each PE
  - 1) Skip IA w/ IDX
  - 2) Skip OA w/ flag in receptive field





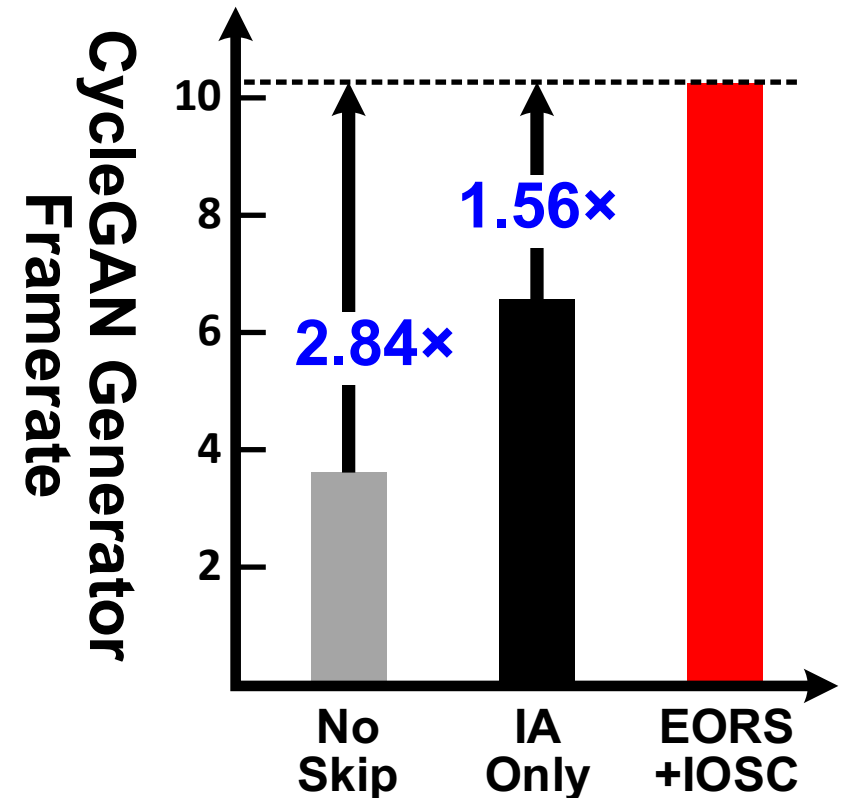
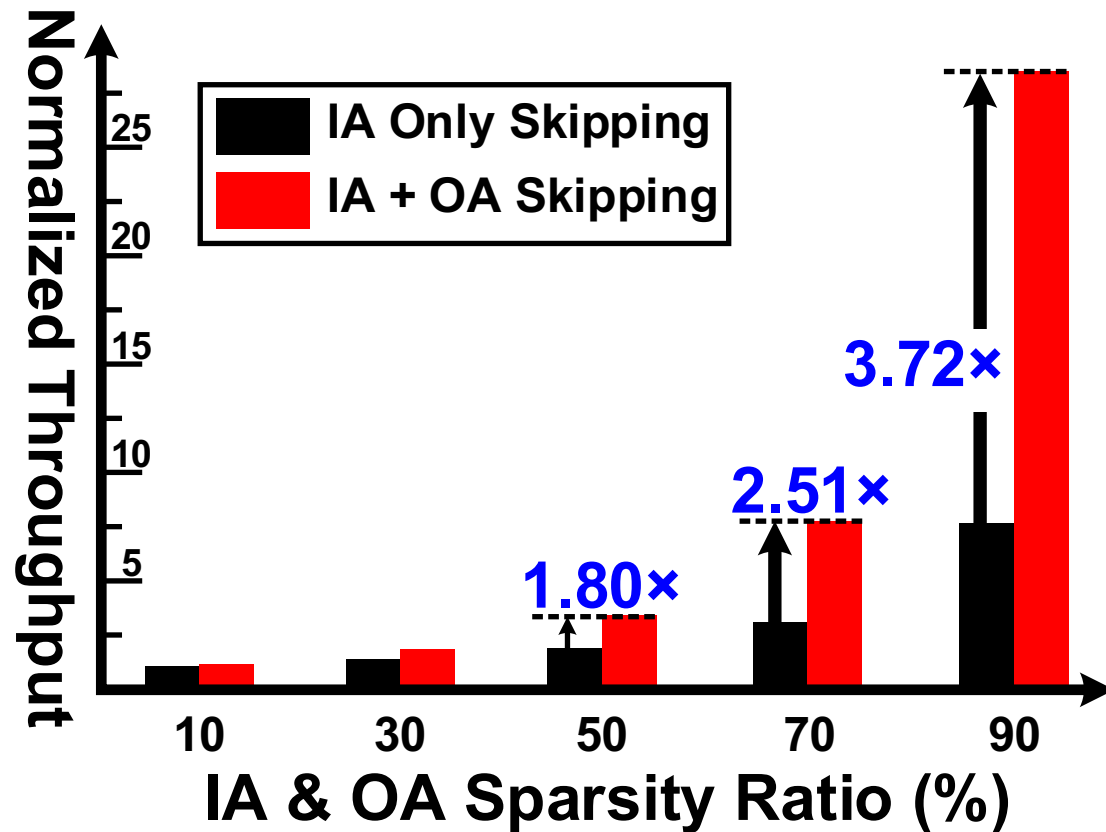
# IOSC PE Architecture

- Zero-Skipping Inside PE

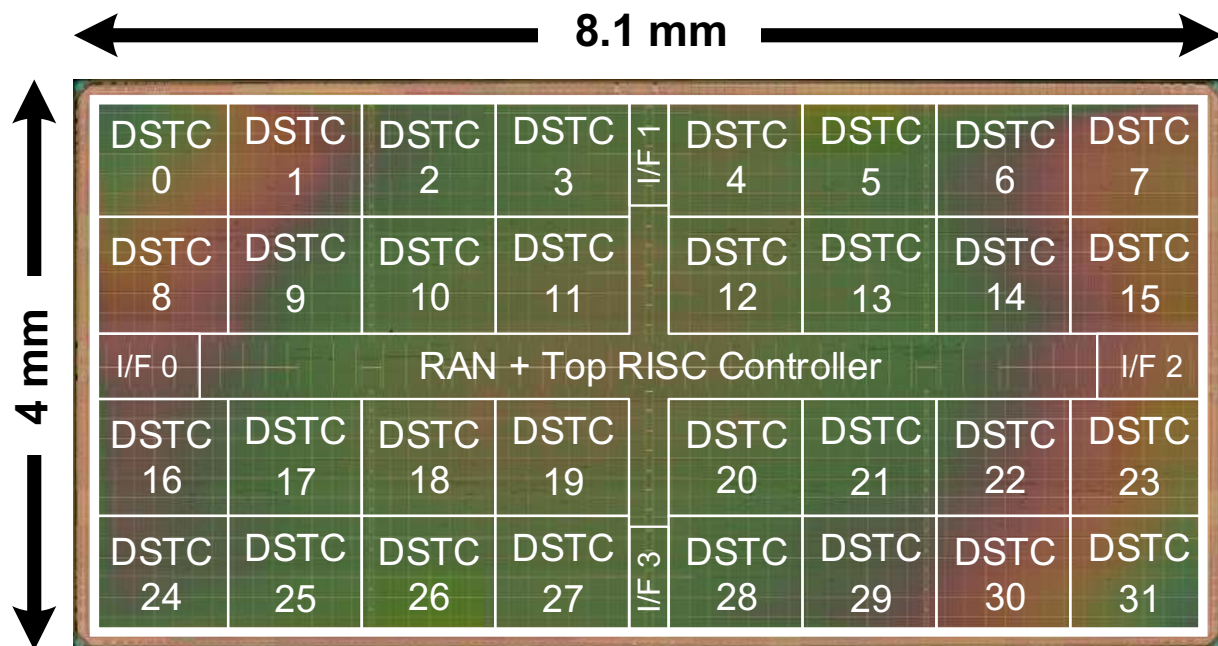


# IOSC Performance

- Throughput Increase with Dual Sparsity Exploitation
  - 2.84x** increase in CycleGAN inference framerate



# Chip Photography and Summary



	Specifications	
Technology	Samsung 65nm 1P8M CMOS	
Die Area	4mm × 8.1mm (32.4mm <sup>2</sup> )	
SRAM	676 KB	
Supply	0.70V ~ 1.1V	
Frequency	~ 200MHz	
Power Consumption	58mW @(25MHz, 0.70V)	
	647mW @(200MHz, 1.1V)	
	FP16	FP8
Peak Performance	0.54 - <b>14.02*</b> TFLOPS	1.08 - <b>24.13*</b> TFLOPS
Power Efficiency [TFLOPS/W]	0.83 - 37.80* @(200MHz, 1.1V)	1.66 - 68.12* @(200MHz, 1.1V)
	1.81 - <b>75.68*</b> @(50MHz, 0.75V)	3.62 - <b>135.10*</b> @(50MHz, 0.75V)

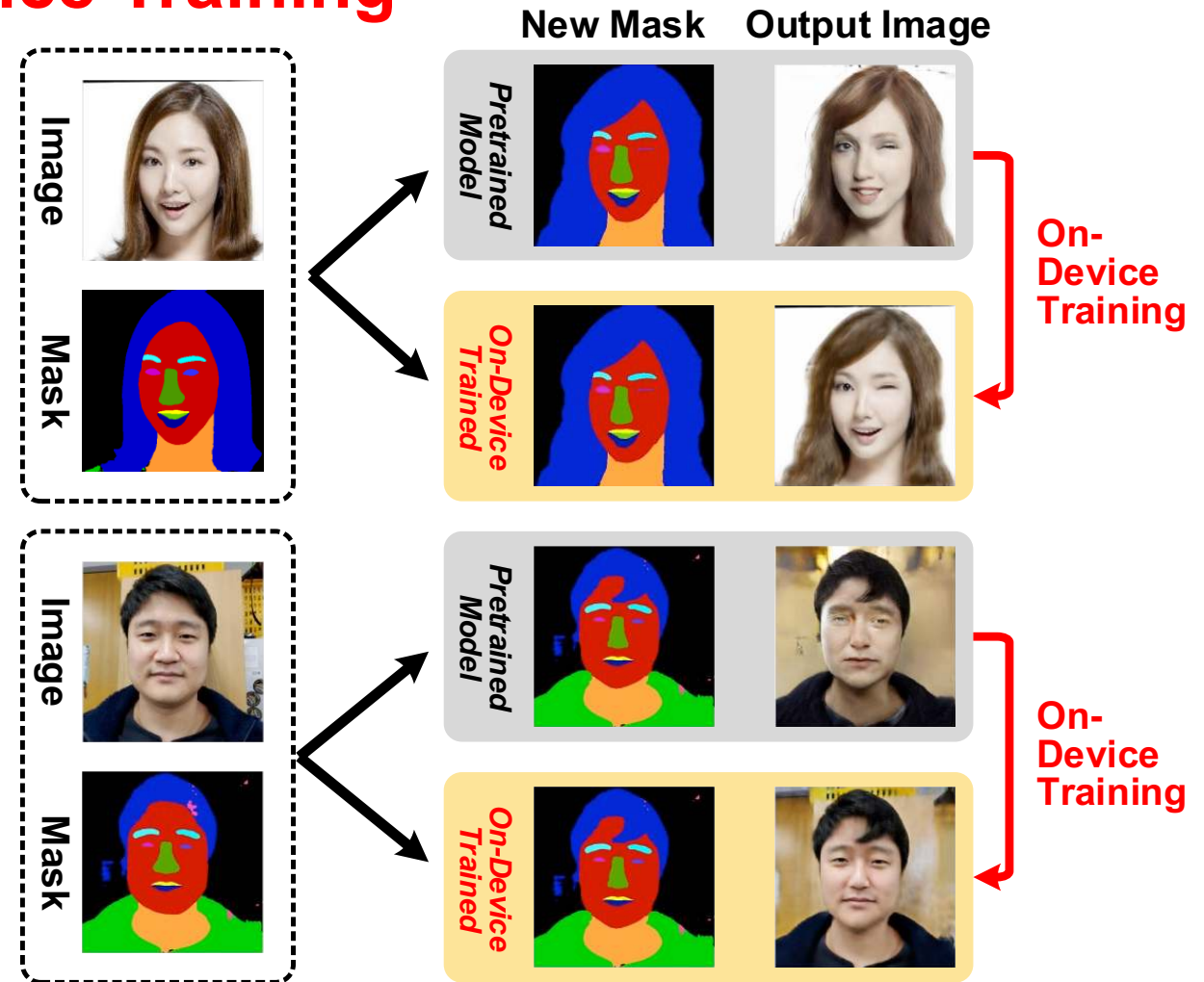
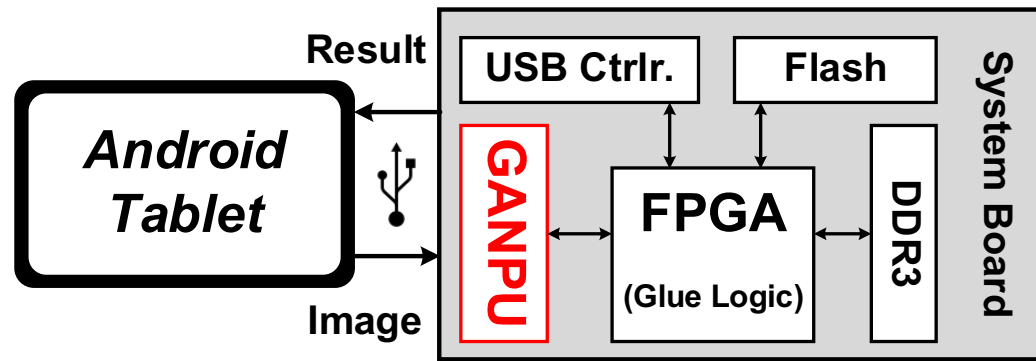
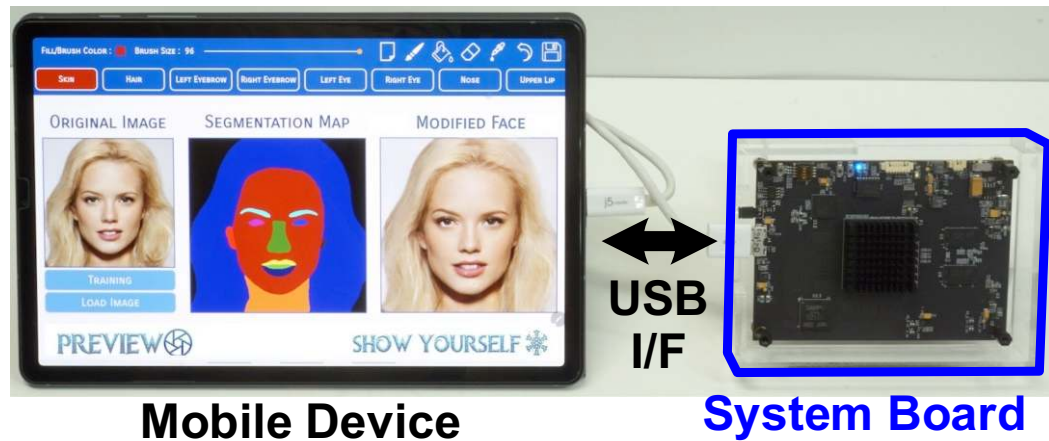
## ▪ CycleGAN Training

- **2.04×** throughput increase
- **76.25%** energy-efficiency increase

\* Logical Performance @ 90% IA&OA Sparsity

# Demonstration System

## Face Manipulation w/ On-Device Training



# Conclusion

- **GANPU: An Energy-Efficient Processor for Training GANs on Mobile Devices**
- **For Energy-Efficient GAN Training**
  - *Adaptive Spatio-Temporal Workload Multiplexing*
  - *IA-OA Sparsity Exploitation Convolution Core Architecture*
  - *Exponent-Only ReLU Speculation*

**A 135 TFLOPS/W Multi-DNN Training Processor  
for Mobile GAN Applications**

# Thank You!

- **Questions? Feel Free to Contact Me!**

- E-mail: [sanghoon\\_kang@kaist.ac.kr](mailto:sanghoon_kang@kaist.ac.kr)

- LinkedIn: <https://www.linkedin.com/in/sanghoonkang94/>

- Zoom Meeting:

- <https://zoom.us/j/97717782673?pwd=Z1BSMGQ2bXVZcmZNVIFzS3ZBL3B1QT09> (Password: HC\_GANPU)

- **Acknowledgement**

- This work was supported by the Samsung Research Funding & Incubation Center for Future Technology under Grant SRFC-TB1703-09.