# An Efficient General-Purpose Optical Accelerator for Neural Networks

## ASP-DAC 2025

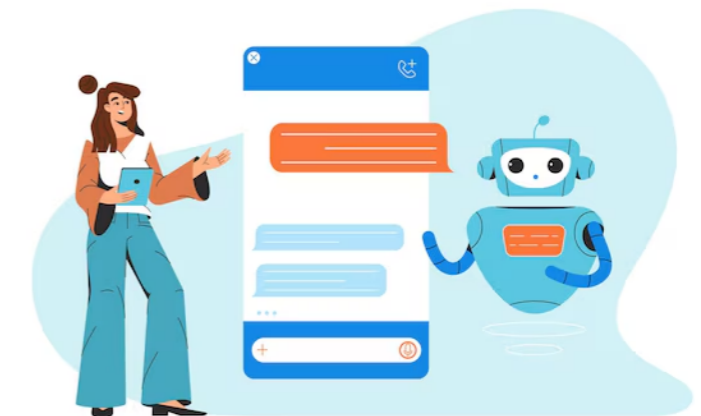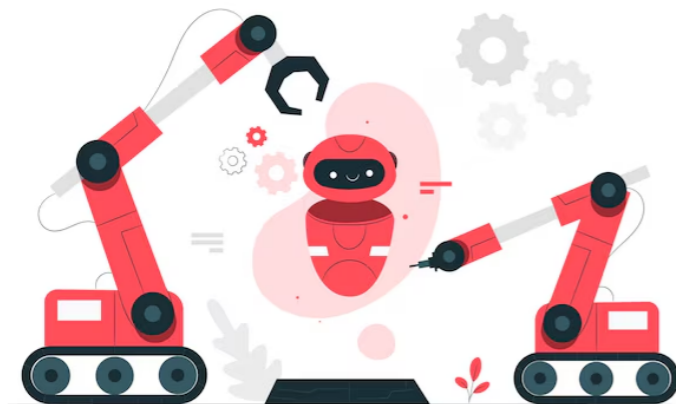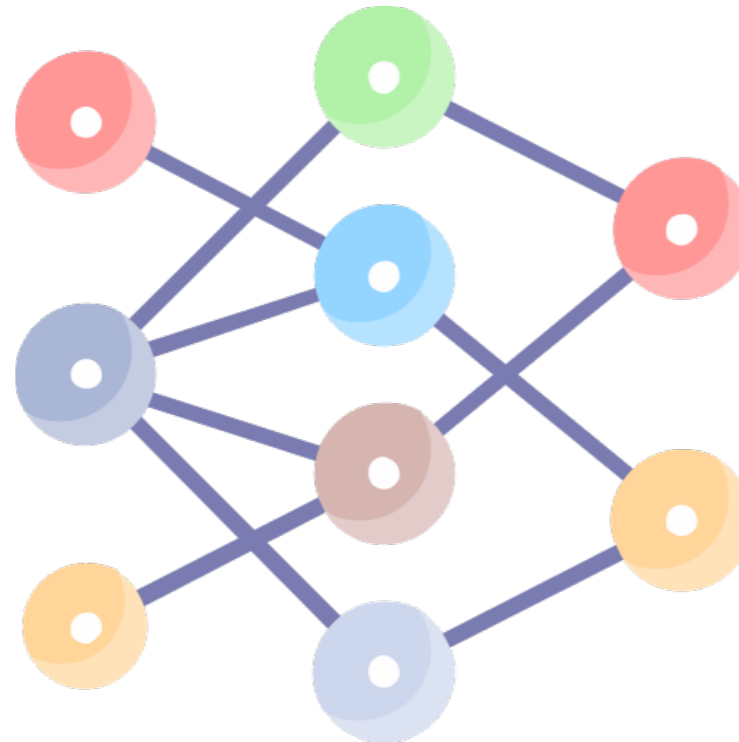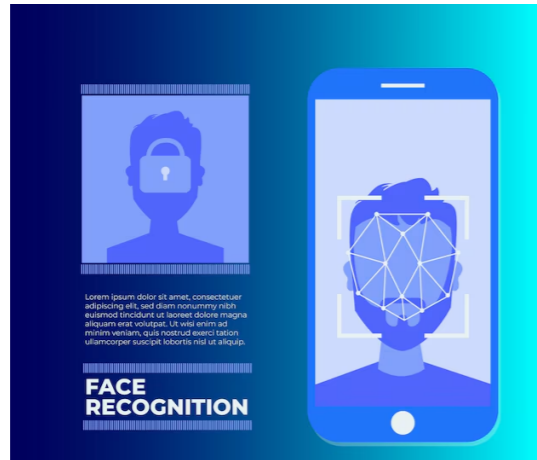**Sijie Fei** (Technical University of Munich)

Amro Eldebiky (Technical University of Munich)

Grace Li Zhang (Technical University of Darmstadt)

Bing Li (University of Siegen)

Ulf Schlichtmann (Technical University of Munich)

# Why **O**ptical **N**eural **N**etworks (**ONN**)
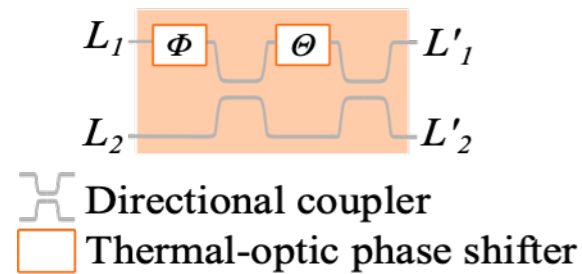
- Strict time constraints
- Millions/Trillions Parameters

  …

higher speed

**ONN**

# Basics of ONN

**Architecture dependent on matrix dimension**

Weight Matrices with cascading MZIs

- **Mach-Zehnder Interferometers** (MZI)

$L_1$ — $\Phi$ — $\Theta$ — $L'_1$

$L_2$ — $L'_2$

Directional coupler

Thermal-optic phase shifter

$$\begin{bmatrix} L'^c_1 \\ L'^c_2 \end{bmatrix} = ie^{\frac{i\theta}{2}} \begin{bmatrix} e^{i\phi}\sin\frac{\theta}{2} & \cos\frac{\theta}{2} \\ e^{i\phi}\cos\frac{\theta}{2} & -\sin\frac{\theta}{2} \end{bmatrix} \begin{bmatrix} L^c_1 \\ L^c_2 \end{bmatrix} = \mathbf{T} \begin{bmatrix} L^c_1 \\ L^c_2 \end{bmatrix}$$
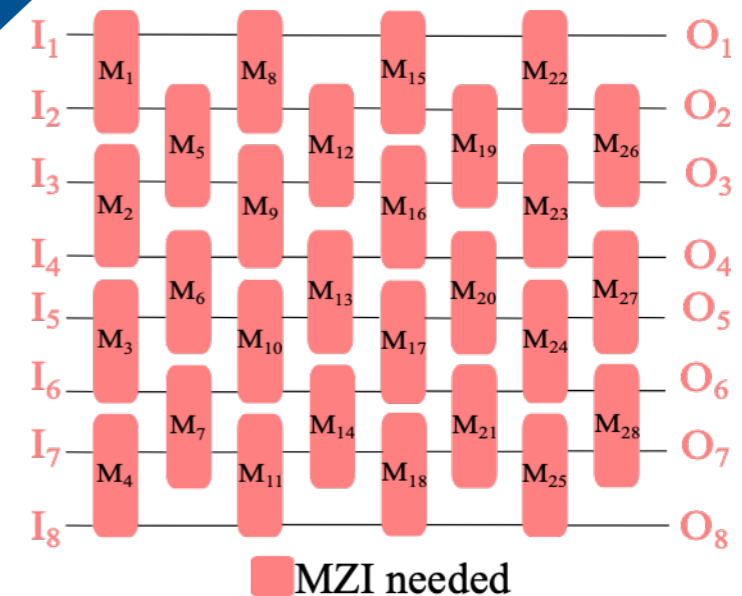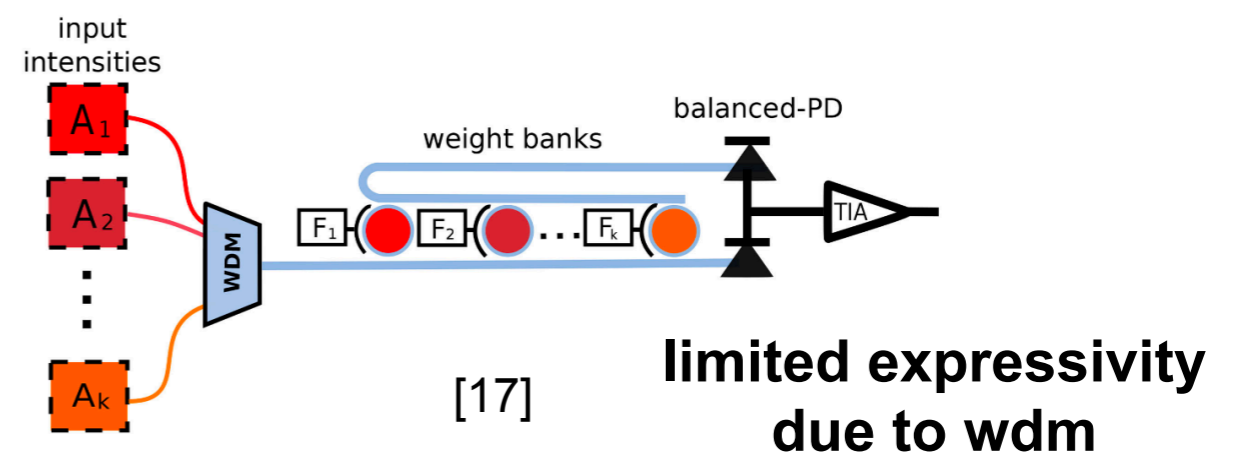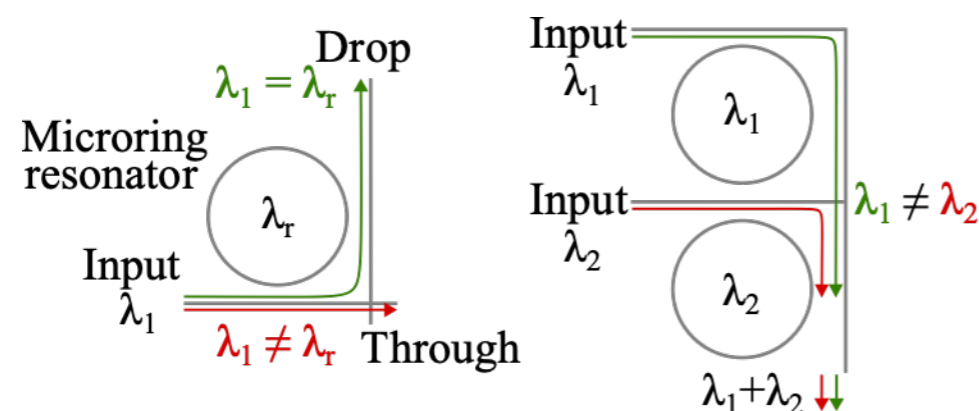
Unitary

SVD

[6]

$I_1, I_2, I_3, I_4, I_5, I_6, I_7, I_8$ — $O_1, O_2, O_3, O_4, O_5, O_6, O_7, O_8$

$M_1, M_8, M_{15}, M_{22}, M_5, M_{12}, M_{19}, M_{26}, M_2, M_9, M_{16}, M_{23}, M_6, M_{13}, M_{20}, M_{27}, M_3, M_{10}, M_{17}, M_{24}, M_7, M_{14}, M_{21}, M_{28}, M_4, M_{11}, M_{18}, M_{25}$

MZI needed

MxN Matrix ~ (M^2+N^2)/2 MZIs

- **Microring Resonator** (MRR)

Drop

$\lambda_1 = \lambda_r$

Microring resonator $\lambda_r$

Input $\lambda_1$

$\lambda_1 \neq \lambda_r$ Through

Input $\lambda_1$

Input $\lambda_2$

$\lambda_1$

$\lambda_2$

$\lambda_1 \neq \lambda_2$

$\lambda_1 + \lambda_2$ ↓↓

input intensities

$A_1$

$A_2$

$A_k$

WDM

weight banks

$F_1$ ... $F_k$
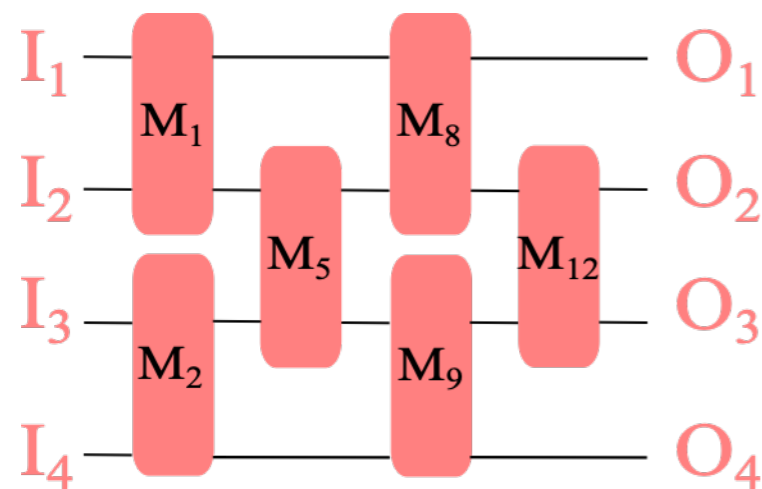
balanced-PD

TIA

[17]

**limited expressivity due to wdm**

[6] YichenShen, NicholasCHarris, Skirlo, etal. Deep learning with coherent nanophotonic circuits. Nature photonics, 11(7):441–446, 2017.
[17] Viraj Bangari, Bicky A Marquez, Heidi Miller, et al. Digital electronics and analog photonics for convolutional neural networks (DEAP-CNNs). IEEE Journal of Selected Topics in Quantum Electronics, 26(1):1–13, 2019.

# Limitation of **G**eneral-purpose **O**ptical **A**ccelerators(**GOA**)

• **GOA:** the same optical accelerator that can be reused for different neural networks
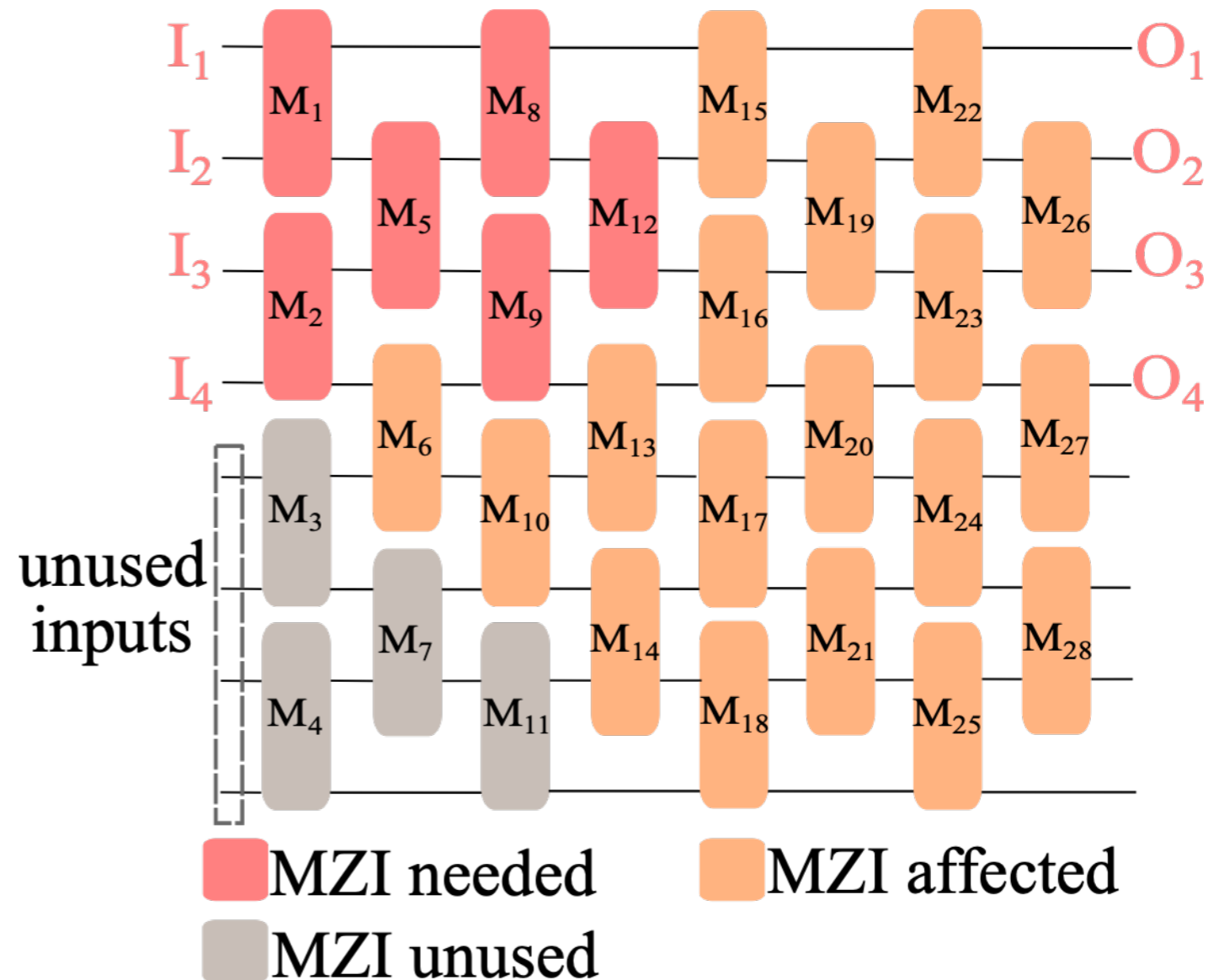


a 4x4 matrix

Mismatch of the matrix dimensions and GOA dimensions

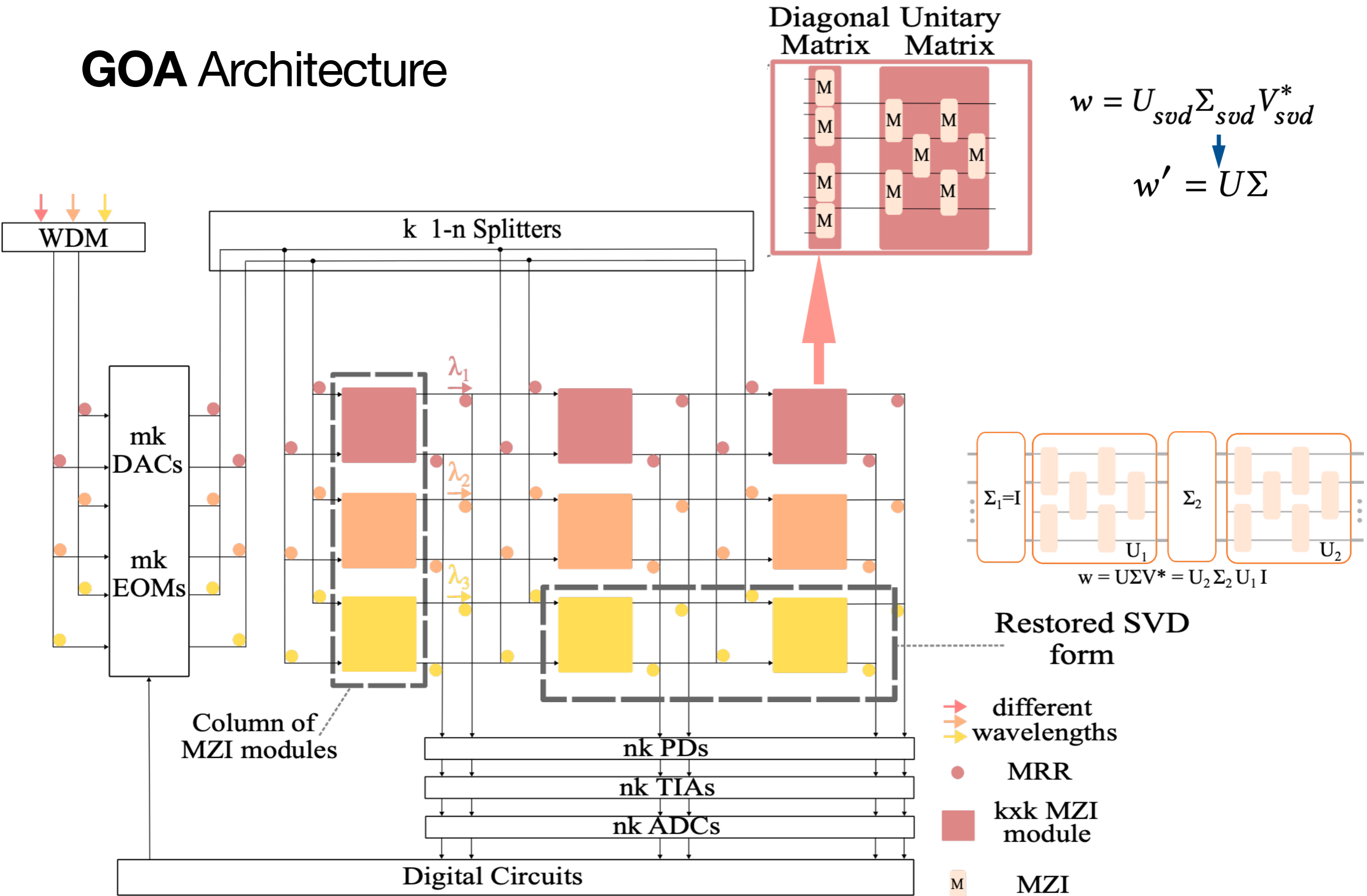**Low utilization efficiency**
**High mapping effort**
**(Mapping: tuning PSs on GOA one time)**

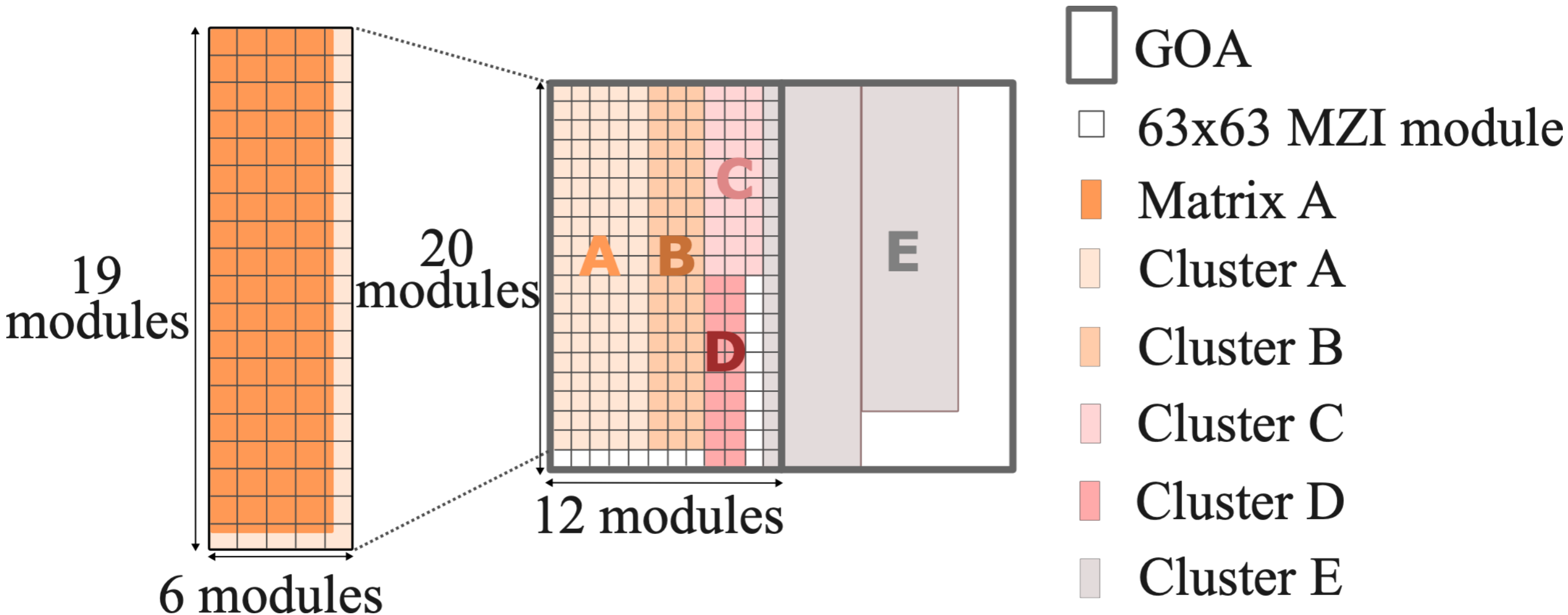

■ MZI needed   ■ MZI affected
■ MZI unused

Represent a 4x4 matrix using an 8x8 GOA
(18+4)/28 MZIs are affected and unused

# **GOA** Architecture



Diagonal Unitary Matrix Matrix

$$w = U_{svd}\Sigma_{svd}V^*_{svd}$$

$$w' = U\Sigma$$

k 1-n Splitters

$\lambda_1$

$\lambda_2$

$\lambda_3$

mk DACs

mk EOMs

WDM

Column of MZI modules

nk PDs

nk TIAs

nk ADCs

Digital Circuits

$\Sigma_1 = I$    $U_1$    $\Sigma_2$    $U_2$

$$w = U\Sigma V^* = U_2\Sigma_2 U_1 I$$

Restored SVD form

→ different
→ wavelengths
→

● MRR

■ kxk MZI module

M MZI

**Higher utiliztiom efficiency and low mapping effort with independent MZI modules**
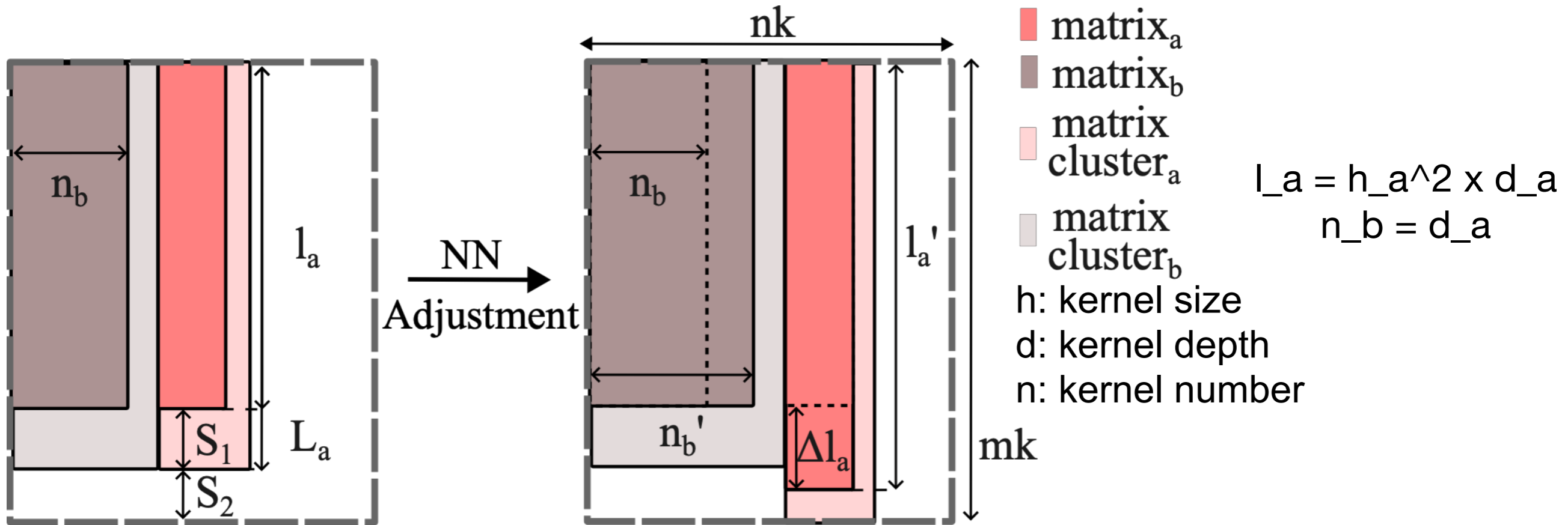
5

# Mapping NNs and Determining **GOA** Parameters



Genetic algorithm with metrics:
- Mapping cost - necessary mappings for one NN
- Area cost
- Power      } MZIs, MRRs and peripheral devices
- E/O conversions

# NN adjustment and Hardware-aware Training



$l\_a = h\_a^2 \times d\_a$

$n\_b = d\_a$

matrix$_a$

matrix$_b$

matrix cluster$_a$

matrix cluster$_b$

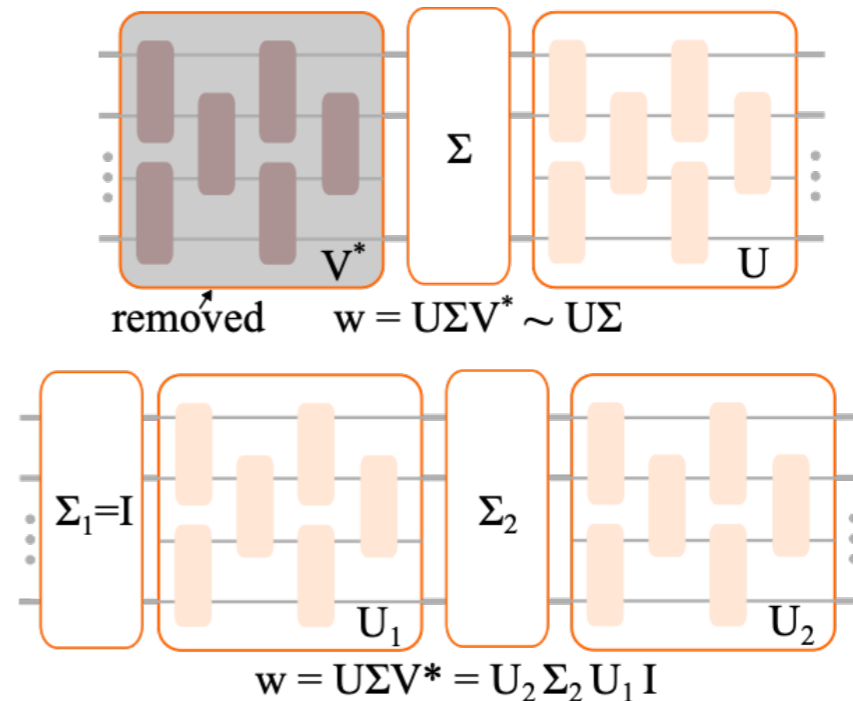h: kernel size
d: kernel depth
n: kernel number

• Training

$$w' = U\Sigma$$

$$U = U_{svd}V^*_{svd}, \text{ where } w = U_{svd}\Sigma_{svd}V^*_{svd}$$

$$\Sigma = \begin{pmatrix} \sigma_{1,1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_{k,k} \end{pmatrix}.$$
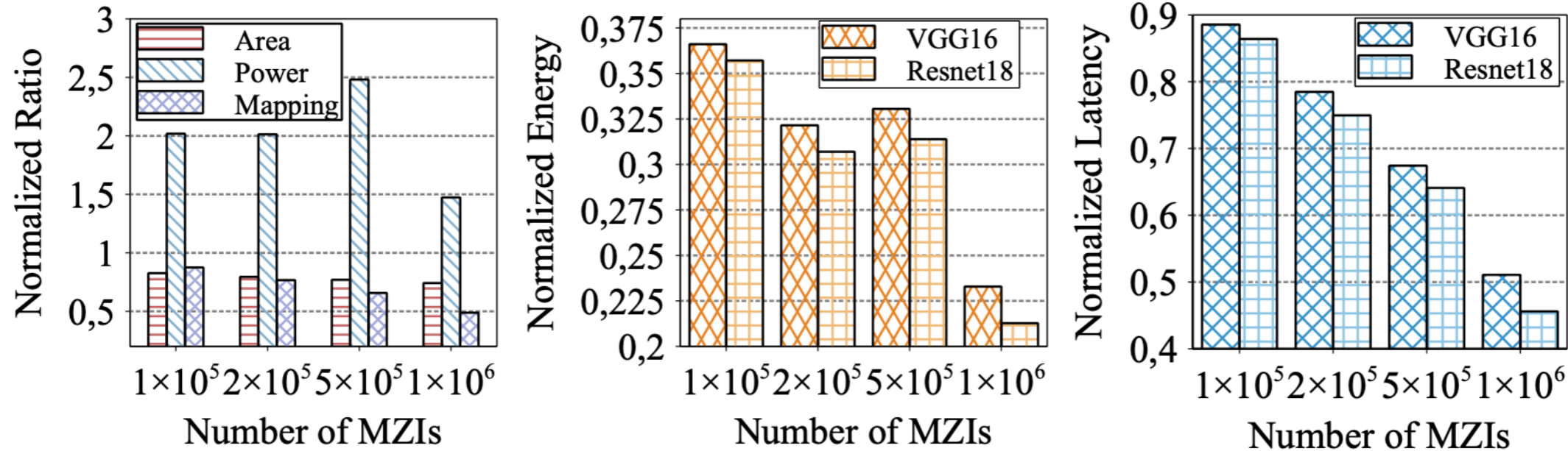
$$\sigma_{k,k} = \underset{\sigma}{\mathrm{argmin}}(\|w_k - \sigma_{k,k}U_k\|_2)$$

• Restoring by columns

removed   $w = U\Sigma V^* \sim U\Sigma$
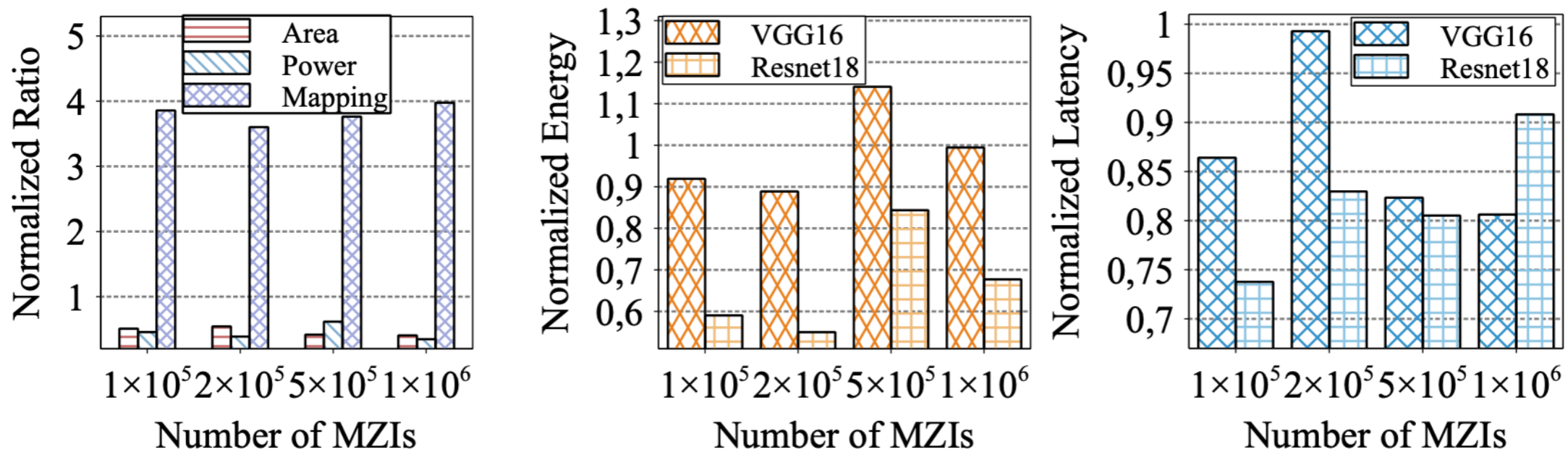
$w = U\Sigma V^* = U_2\Sigma_2 U_1 I$

# Experimental Results-Mapping/Energy/Latency Analysis

- Compared With SVD accelerators [6]
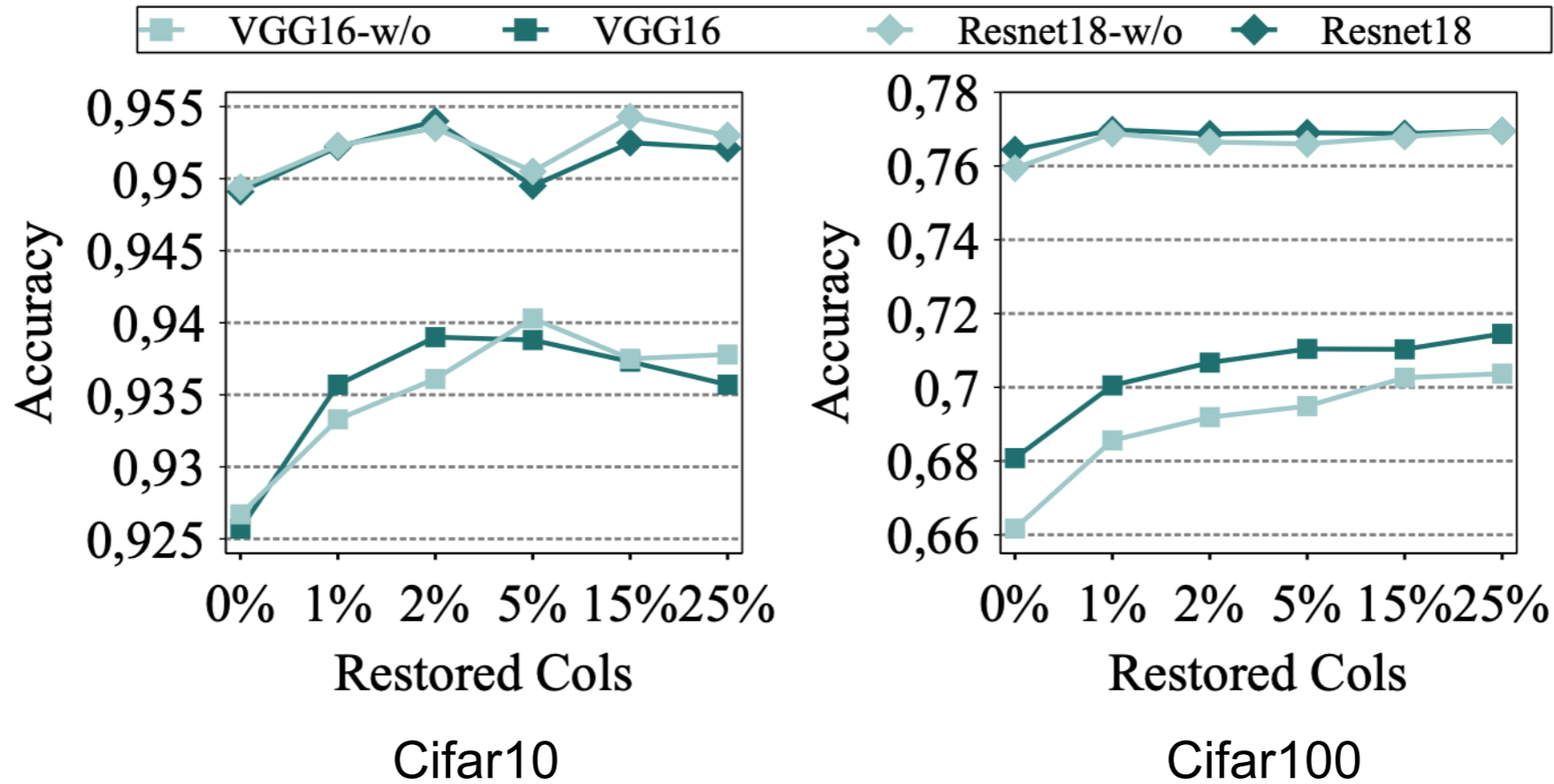


- Compared With Adept accelerators [12]

[6] YichenShen, NicholasCHarris, Skirlo, etal. Deep learning with coherent nanophotonic circuits. Nature photonics, 11(7):441–446, 2017.

[12] Jiaqi Gu, Hanqing Zhu, Chenghao Feng, et al. Adept: Automatic differentiable design of photonic tensor cores. In Design Automation Conference (DAC), 2022.

# Experimental Results-Accuracy Analysis

**Table 2: Results of the proposed framework. MZI number constraint: 20000, $m$, $n$, $k$ = 6,3,44**

| Neural Networks Dataset | Performance Improvement | | | Accuracy | | | | Restored Cols |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Mapping Reduction | Energy Reduction | Latency Reduction | Baseline | This Work w/o adjustment | This Work w/o restoration | This Work | |
| VGG16-Cifar10 | 21.87% | 67.96% | 21.85% | 93.55% | 92.57% | 92.67% | 93.57% | 1% |
| VGG16-Cifar100 | 21.20% | 67.71% | 21.19% | 70.16% | 67.12% | 68.35% | 70.67% | 2% |
| Resnet18-Cifar10 | 24.69% | 69.13% | 24.61% | 94.93% | 94.91% | 94.94% | 95.22% | 1% |
| Resnet18-Cifar100 | 25.52% | 69.47% | 25.45% | 75.79% | 75.94% | 76.44% | 76.44% | 0% |



Cifar10                    Cifar100

# Conclusion

- To reduce mapping effort：

  - a GOA architecture of independent MZI modules is proposed

  - #params of GOA is determined by balancing the area cost, power, mapping cost, and E/O conversions

  - NN adjustments, hardware-aware training, restoration of weight matrices are performed to ensure accuracy

- Mapping efficiency improved up to 25.52%, energy saved up to 67%, latency saved up to 21%, compared to the SVD accelerator

# Thank you!