



TEXAS

The University of Texas at Austin



MLSys

Light-AI Interaction: Bridging Photonics and AI with Cross-Layer Hardware/Algorithm Co-Design

Jiaqi Gu^{1,2}, Hanqing Zhu¹, Chenghao Feng^{1,3},
Ray T. Chen¹, David Z. Pan¹

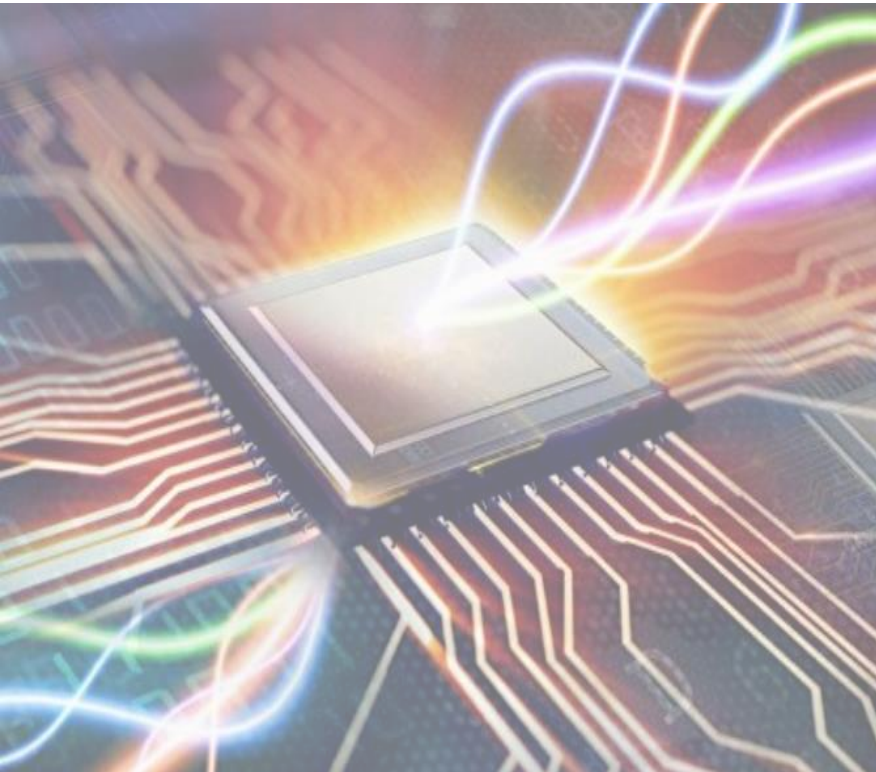
¹The University of Texas at Austin, ²Arizona State University,
³Alpine Optoelectronics

jqgu@utexas.edu; <http://jqgu.net>

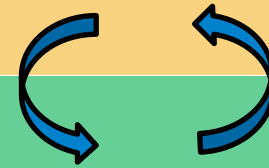
This work was supported in part by AFOSR MURI

June 15, 2023

What is This Talk About: Overview of Optical AI



Design **Integrated Photonics Hardware**
for AI / ML Computing



Apply **AI / ML** for Photonic Hardware
Design Automation

Photonic AI Computing Basics

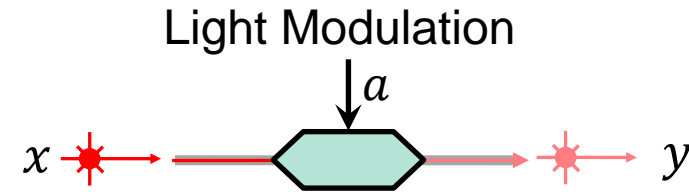
- ◆ Principle: modulation (encode), interference (MVM), photodetection (readout)
- ◆ Good at ultra-fast (10-100ps), parallel linear operations in the analog domain
- ◆ 10 TOPS/W (SoTA) → 1M TOPS/W (potential)

Computing Primitives

Photonic Implementation

Scalar Multiply

$$y = a \cdot x$$

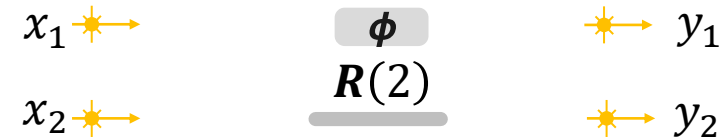


2x2 Unitary Matrix Multiply

$$y = R(2) \times x$$

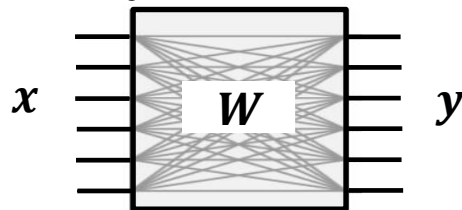
$$R(2) = \begin{pmatrix} \cos \phi & -\sin \phi \\ \sin \phi & \cos \phi \end{pmatrix}$$

Mach-Zehnder Interferometer (MZI)

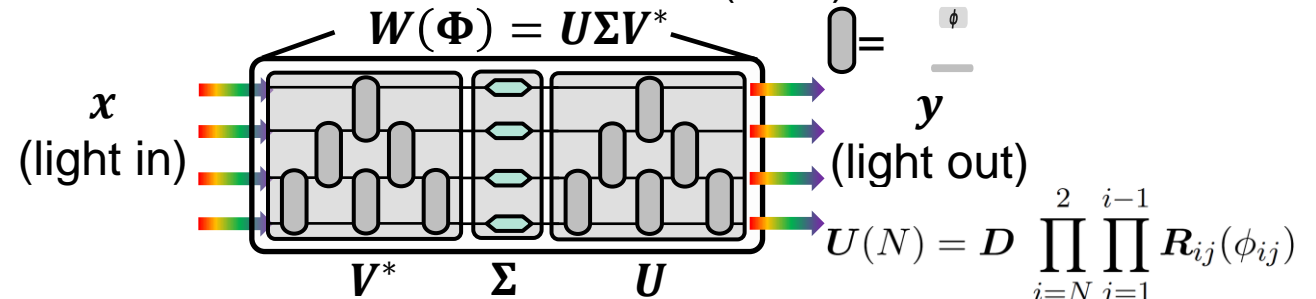


Matrix-Vector Multiply (MVM)

$$y = W \times x$$



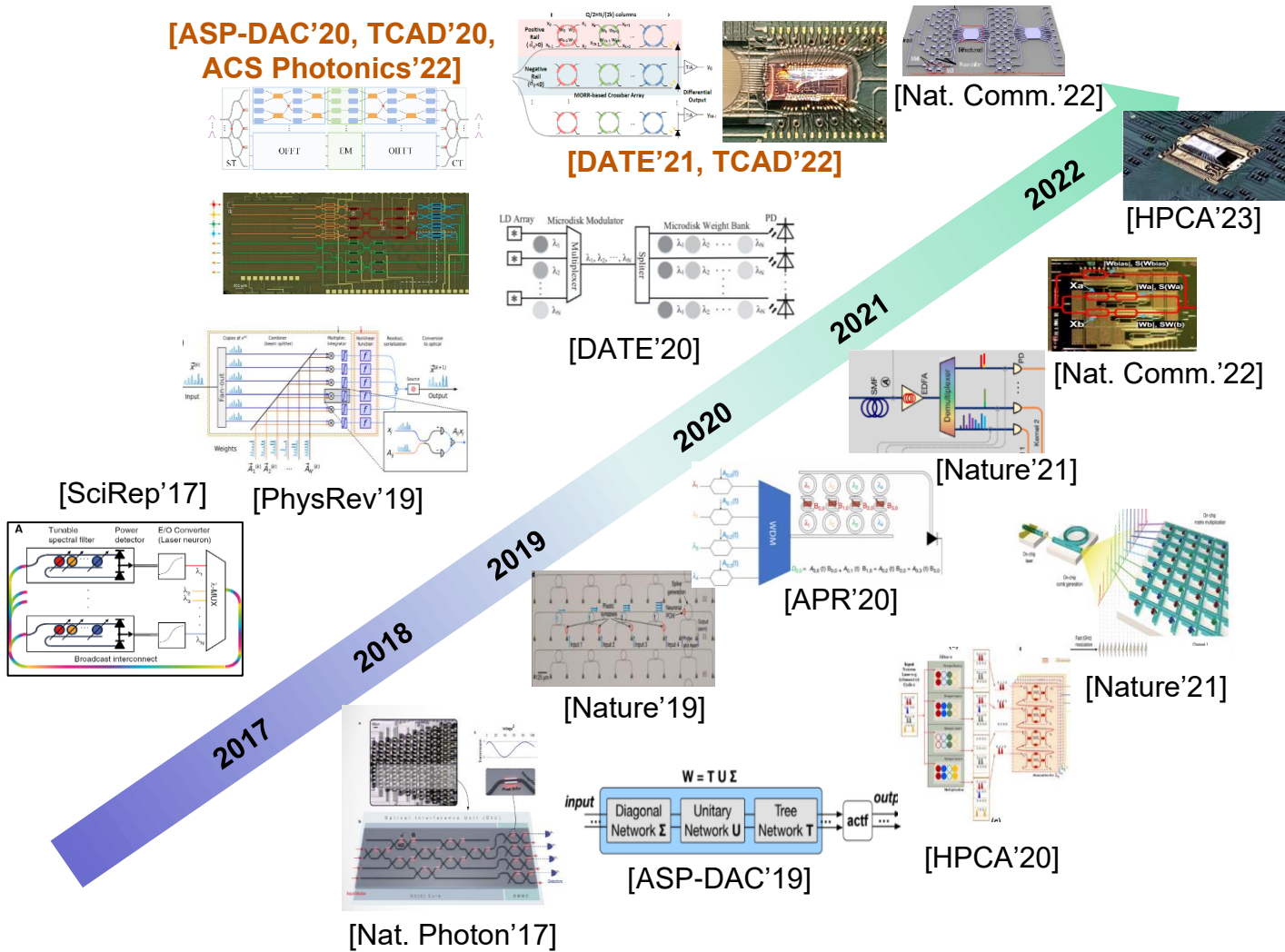
Photonic Tensor Core (PTC)



One-shot computing at speed-of-light!

Photonic AI is Booming

Photonic Neural Network Trends in Academia



Foundry / EPDA Support in Industry

Photonic Computing Chip Designs



LIGHTMATTER

Lighton



XANADU

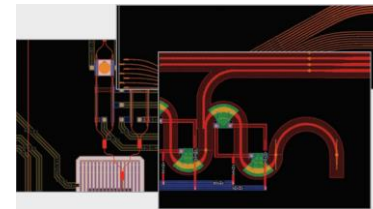
LIGHTELLIGENCE
LUMINOUS

OPTELLIGENCE
OPTIUS
light powered computing

CogniFiber

SALIENCE
LABS

Electronic-Photonic Design Automation Tools



Ansys

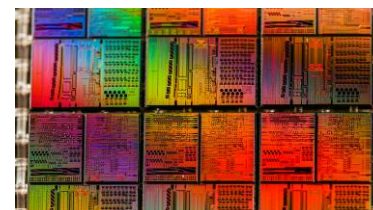
SYNOPTIS[®]

LUMERICAL

PHOTONIC SOLUTIONS

cadence[®]
PHOTONICS

PDK / Tape-out / Packaging Support



amf
ADVANCED
MICRO
FOUNDRY

AIM
photonics

SIEPIC

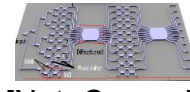
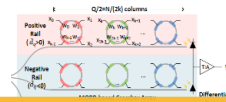
GlobalFoundries[™]

Tower
Semiconductor

Photonic AI is Booming

Photonic Neural Network Trends in Academia

[ASP-DAC'20, TCAD'20, ACS Photonics'22]



[Nat. Comm.'22]

Challenging to Design **Scalable, Robust & Adaptive** Photonic ML Computing Platforms

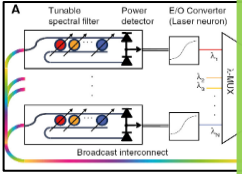
[DATE'20]

[Nat. Comm.'22]

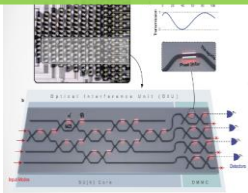
2020

[SciRep'17]

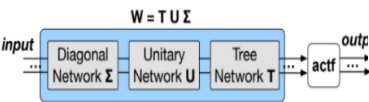
[PhysRev'19]



2017



[Nat. Photon'17]



[ASP-DAC'19]



[HPCA'20]

Foundry / EPDA Support in Industry

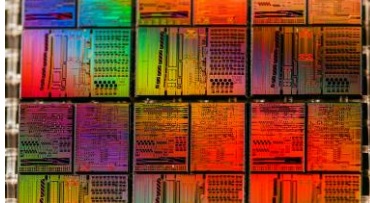
Photonic Computing Chip Designs

Lightmatter
Lighton
XANADU
LIGHTINTELLIGENCE
LUMINOUS
OPTIUS
light powered computing
SALIENCE LABS
Cognifiber

Electronic-Photonic Design Automation Tools

Ansys
SYNOPSYS
LUMERICAL
PHOTONIC SOLUTIONS
cadence

PDK / Tape-out / Packaging Support



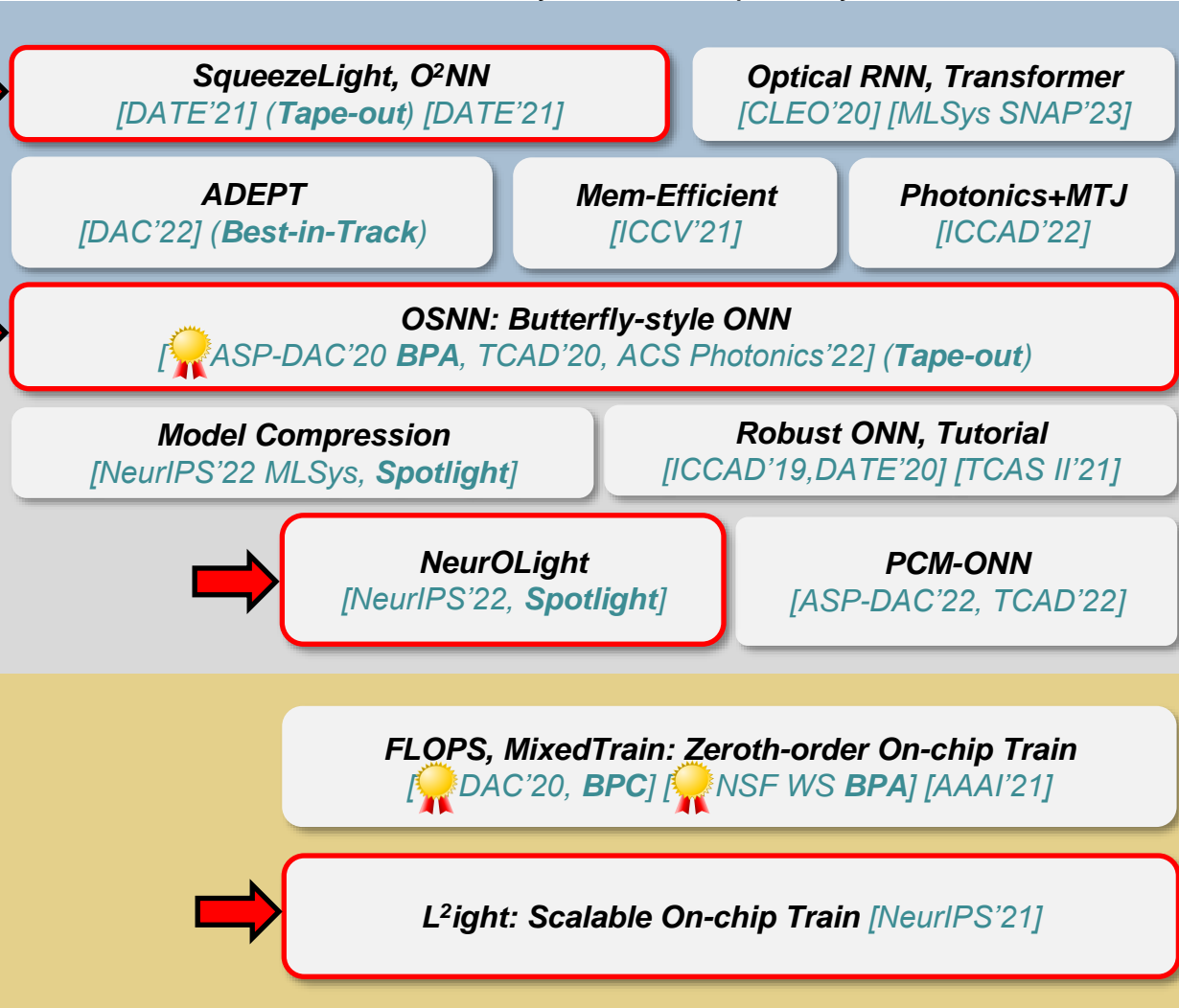
amf ADVANCED MICRO FOUNDRY
AIM photonics
SiEPIC
GlobalFoundries
Tower Semiconductor

Virtuous Cycle: Photonics for AI ↔ AI for Photonics

Work & Contributions

Publications: >30 in CAD, ML, Photonics Communities

Area Efficiency Adaptability Robustness



Open-Source Release

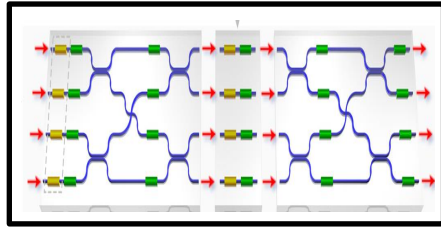


A PyTorch Library for Photonic Integrated Circuit Simulation and Photonic AI Computing

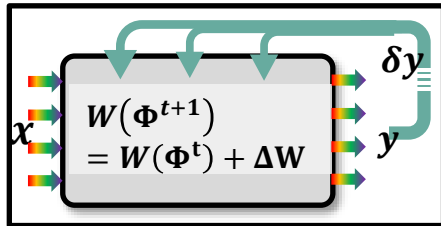
Fork 12 Starred 205



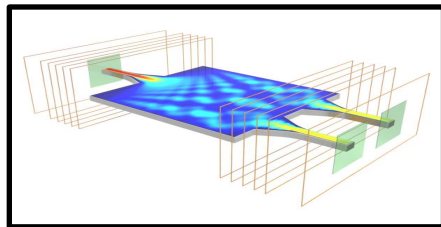
Outline



- ◆ Customized Optical Neural Network Design

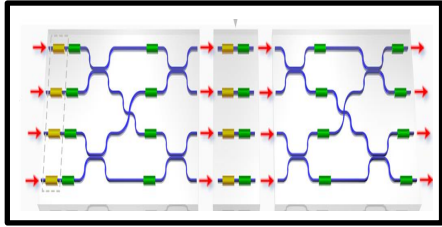


- ◆ ONN On-Chip Training Algorithms

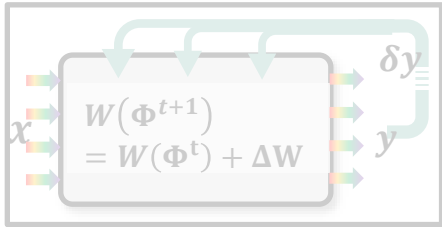


- ◆ ML-Assisted Photonic Design Automation

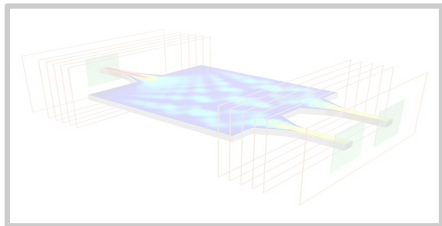
Outline



- ◆ Customized Optical Neural Network Design



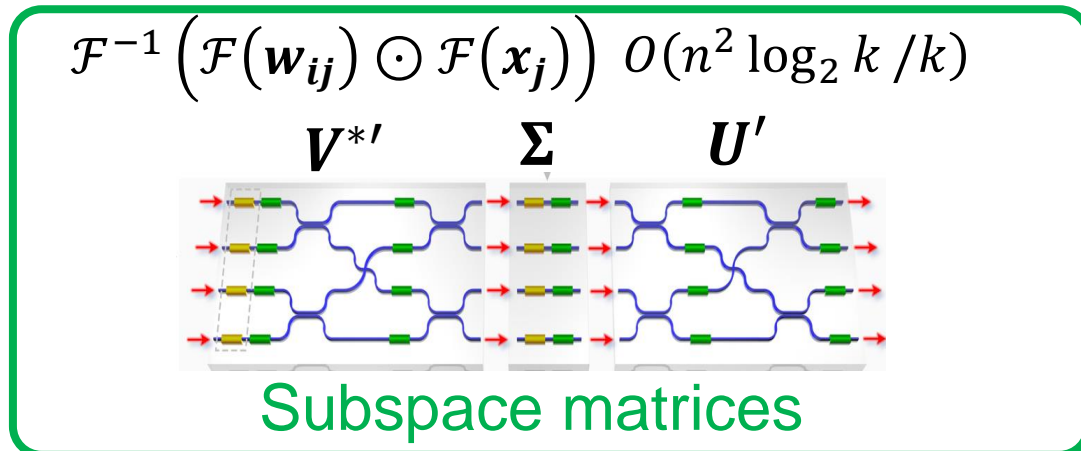
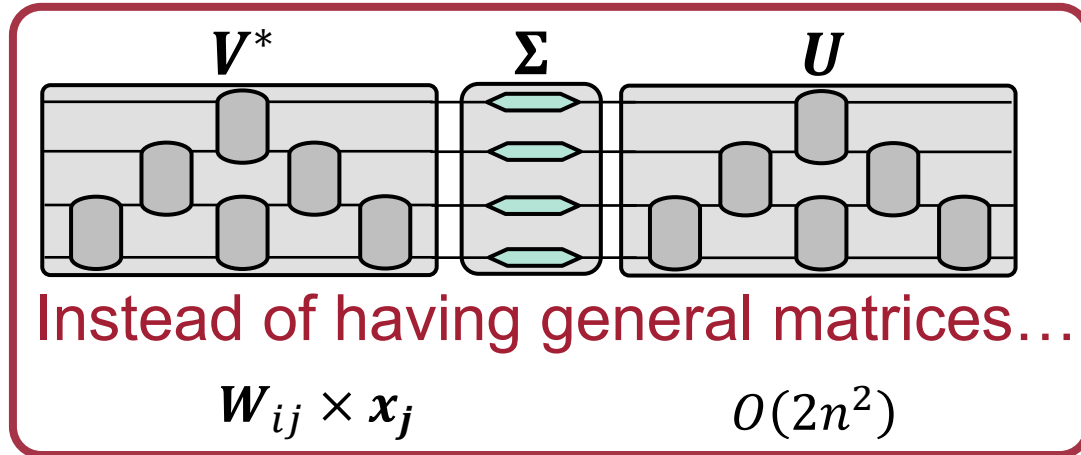
- ◆ ONN On-Chip Training Algorithms



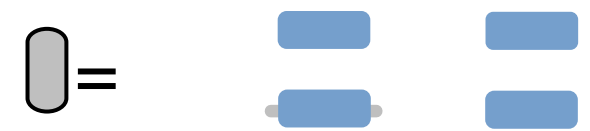
- ◆ ML-Assisted Photonic Design Automation

From GEMM To Specialized Subspace Linear

- ◆ Overparameterized DNN → GEMM is not necessary → “circuit compression”
- ◆ Large universal $U\Sigma V$ MZI array → Compact subspace $U'\Sigma V'$ butterfly mesh



Large MZI array
 $O(n^2)$ MZIs



Compact butterfly mesh
 $O(n \log n)$ basic devices

Break it to basic devices

50/50 coupler

$$\frac{1}{\sqrt{2}} \begin{pmatrix} 1 & j \\ j & 1 \end{pmatrix}$$

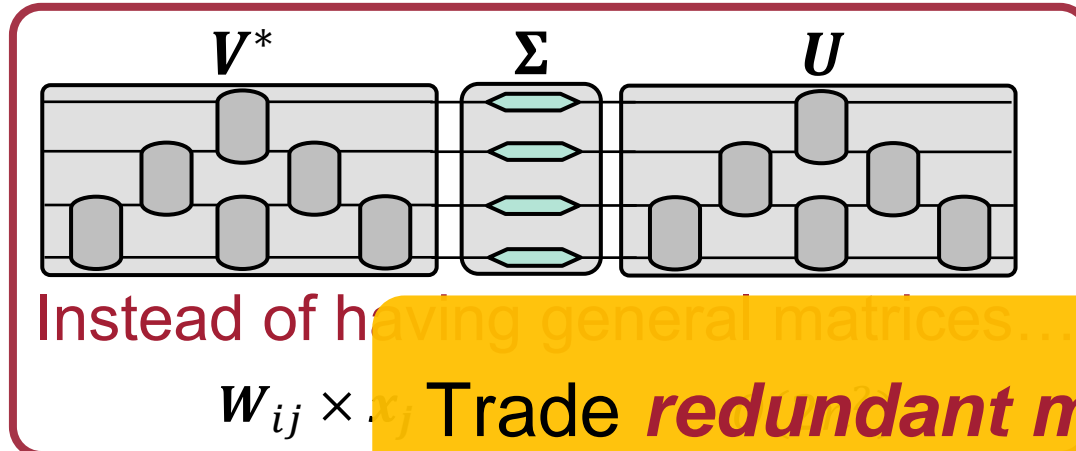
phase shifter

$$e^{j\phi}$$

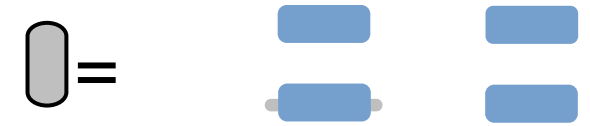


From GEMM To Specialized Subspace Linear

- ◆ Overparameterized DNN → GEMM is not necessary → “circuit compression”
- ◆ Large universal $U\Sigma V$ MZI array → Compact subspace $U'\Sigma V'$ butterfly mesh



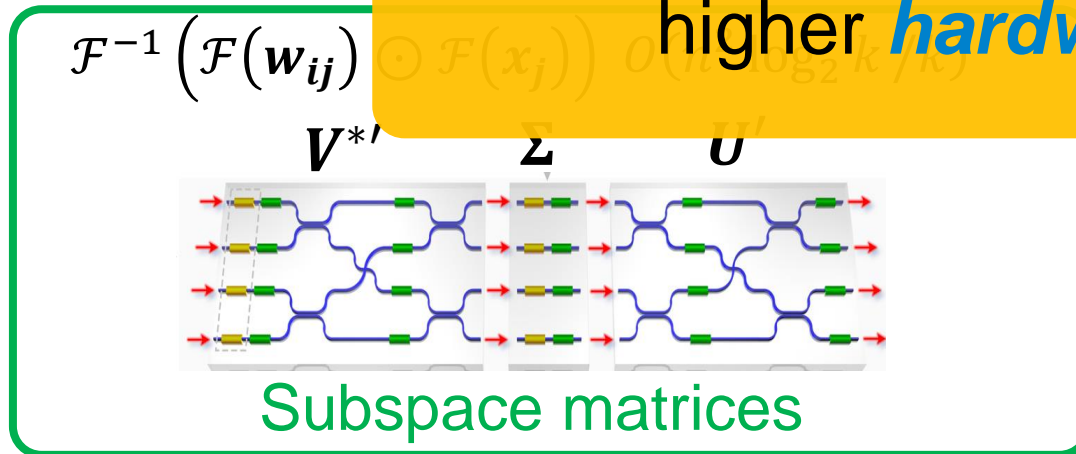
Large MZI array
 $O(n^2)$ MZIs



Instead of having general matrices...

$$W_{ij} \times x_j$$

Trade **redundant matrix expressivity** for higher **hardware efficiency**



Subspace matrices

$O(n \log n)$ basic devices

basic devices

50/50 coupler

$$\frac{1}{\sqrt{2}} \begin{pmatrix} 1 & j \\ j & 1 \end{pmatrix}$$

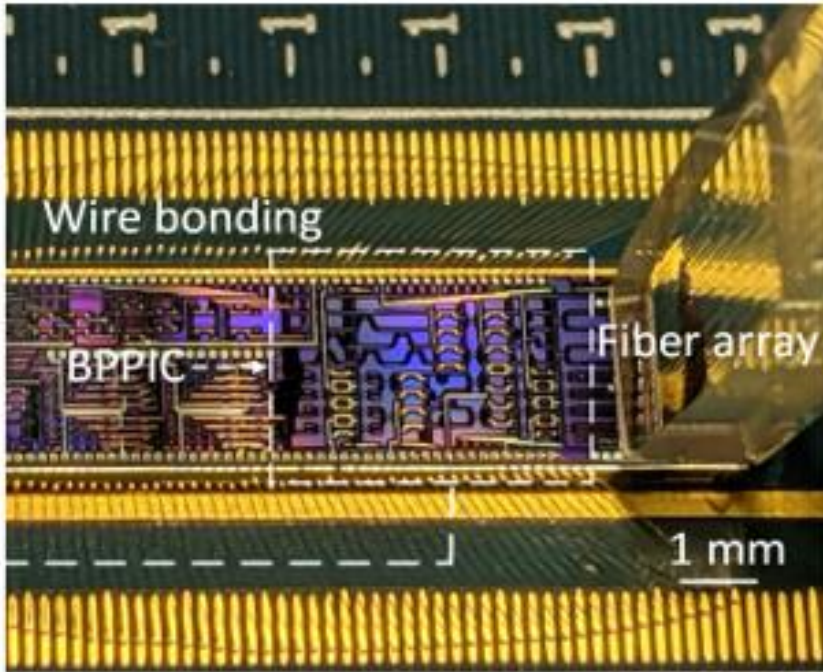
phase shifter

$$e^{j\phi}$$

Photonic Neural Chip Tapeout & Demonstration



4x4 butterfly tensor core



Peak compute density & energy efficiency

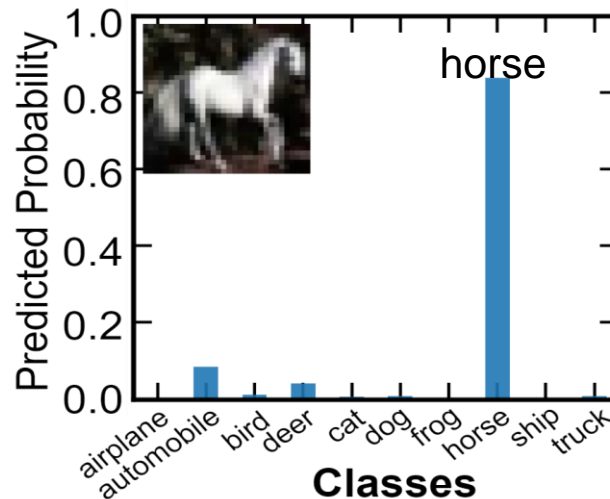
- **Our ONN** (225 TOPS/mm², 9.5 TOPS/W)
- NVIDIA A100 (0.76 TOPS/mm², 1.56 TOPS/W)
- Google TPUv4 (0.69 TOPS/mm², 1.62 TOPS/W)
- 40nm RRAM Accel (0.03 TOPS/mm², 2.20 TOPS/W) [Giordano+, VLSI'21]

>85% accuracy

ResNet-20 (0.27M #param) CIFAR-10

3-bit weight resolution

Fixed butterfly transform

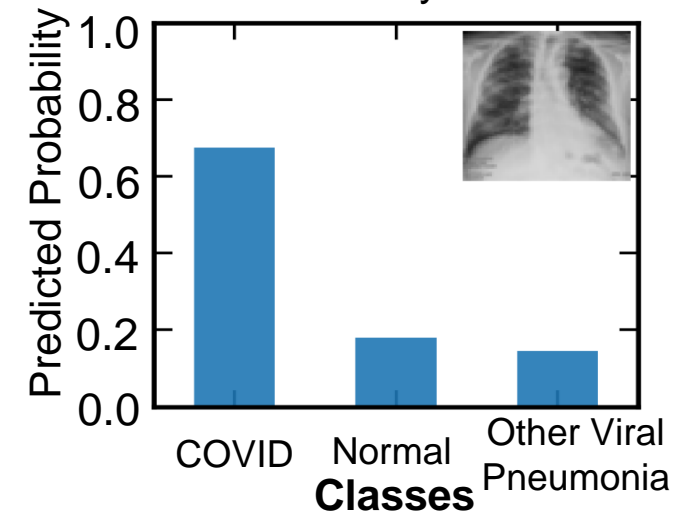


96.5% accuracy

VGG8 (4M #param) COVID Chest X-ray

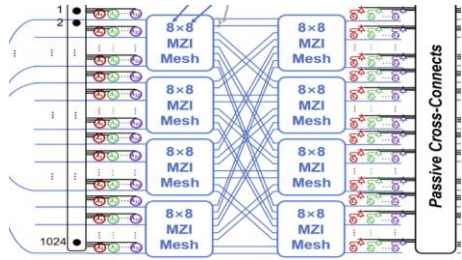
3-bit weight resolution

Fixed butterfly transform

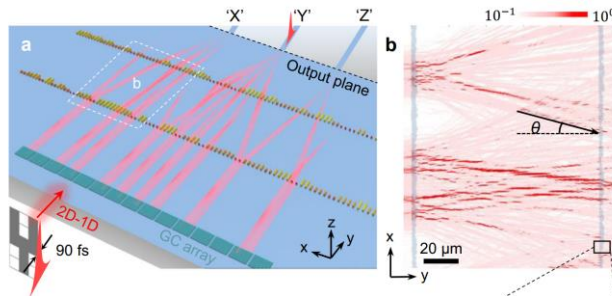


More Customized ONN Designs Beyond GEMM

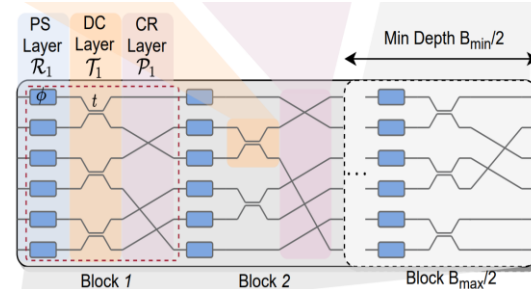
◆ Specialized circuits for hardware-efficient subspace linear op



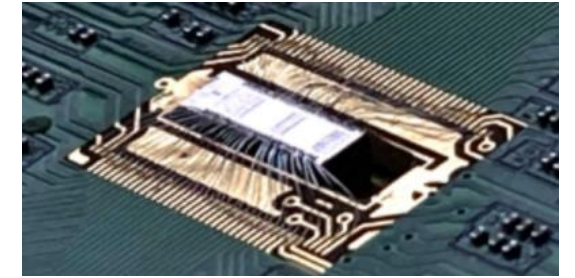
Tensorized MZI-ONN
[Xiao+, *APL Photonics* 2021]



Metalens-based diffractive ONN
[Wang+, *Nat. Commun* 2022]

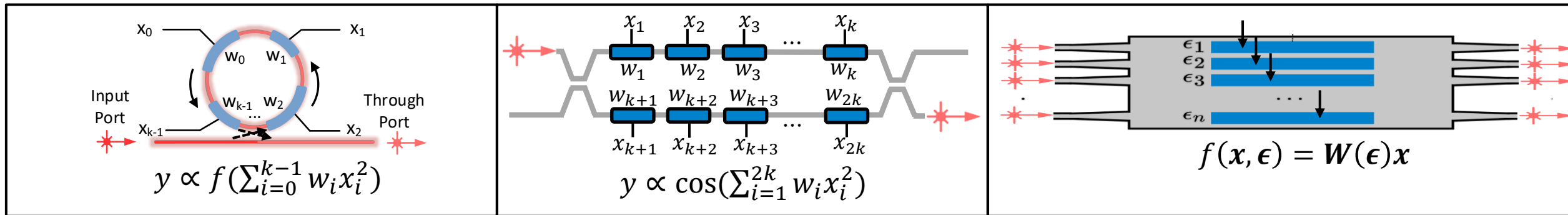


Auto-designed PIC topology
[Gu+, *DAC* 2022]



Fourier lens photonic Conv
[Li+, *HPCA* 2023]

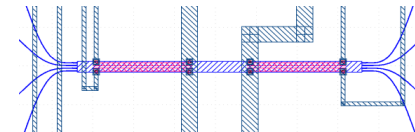
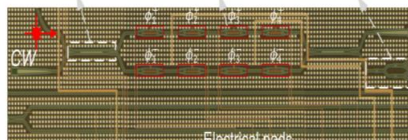
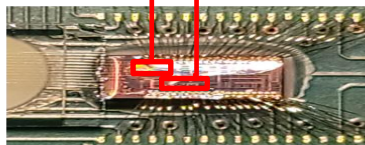
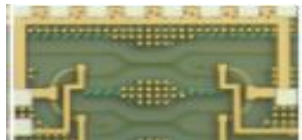
◆ Customized devices beyond 1 MAC/device → single-device vector/MVM unit



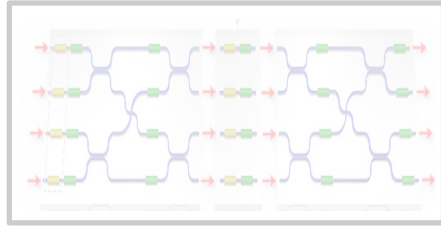
MORR: [Gu+, DATE'21, TCAD'22]

MOMZI: [Feng+, Pho.West'23, arXiv'23]

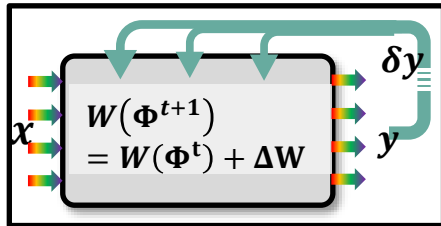
MOMMI: [Gu+, under sub., arXiv'23]



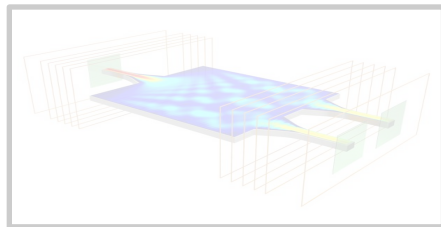
Outline



- ◆ Customized Optical Neural Network Design



- ◆ ONN On-Chip Training Algorithms



- ◆ ML-Assisted Photonic Design Automation

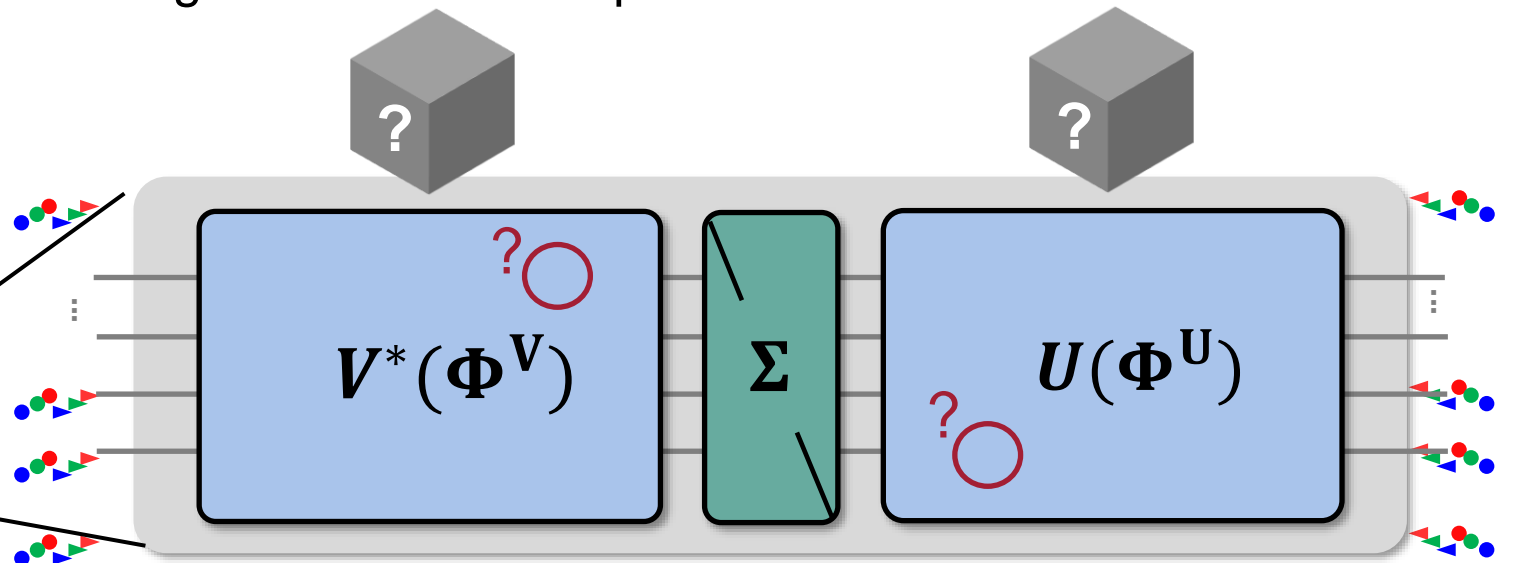
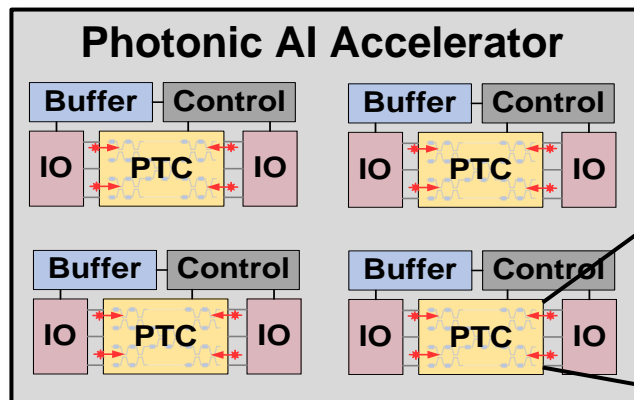
Inference → Training: Self-Learnable AI Engine

- ◆ Why on-chip training? *reliability, adaptability, efficiency, privacy...*



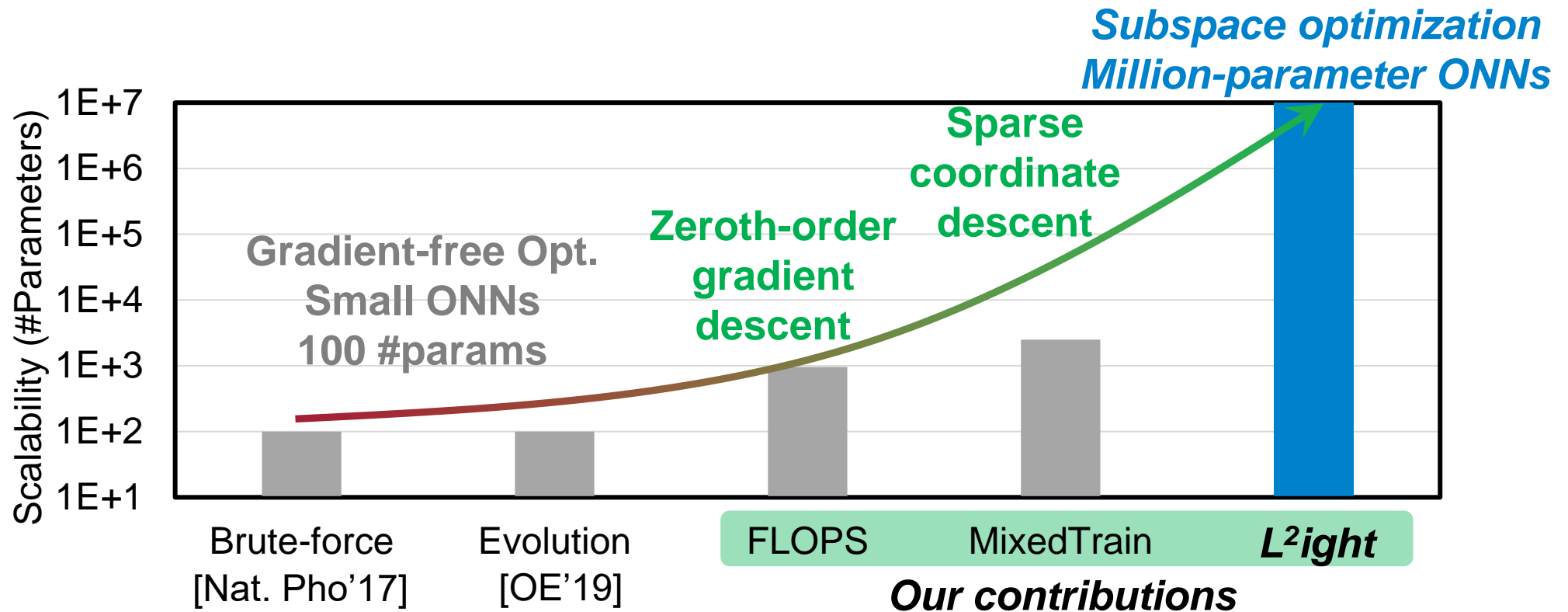
- ◆ Challenges

- › No access to **intermediate states or full gradients** (U/V are blackbox)
- › **Noisy** circuits (randomness)
- › Algorithm must be **simple** enough to be run on chip



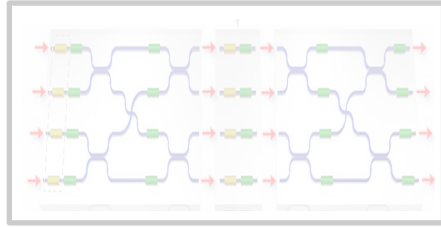
Efficient On-Chip Training Protocols

- ◆ $>10,000\times$ trainability \uparrow $+30\times$ efficiency \uparrow : Customize algorithm for the hardware
- ◆ Utilize optics reciprocity to calculate subspace 1st-order gradients with sparsity

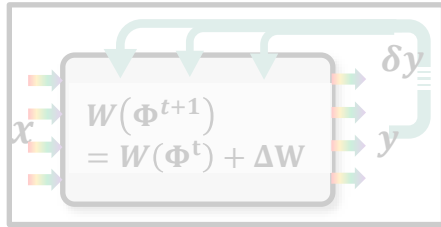


- J. Gu, Z. Zhao, *et al.*, **FLOPS**, ACM/IEEE Design Automation Conference (DAC), 2020 (**Best Paper Finalist**) (**Best Poster Award**)
- J. Gu, C. Feng, *et al.*, **Mixed-Train**, Association for the Advancement of Artificial Intelligence (AAAI), 2021
- J. Gu, H. Zhu, *et al.*, **L2ight**, Conference on Neural Information Processing Systems (NeurIPS), 2021

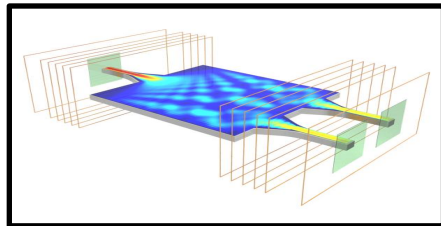
Outline



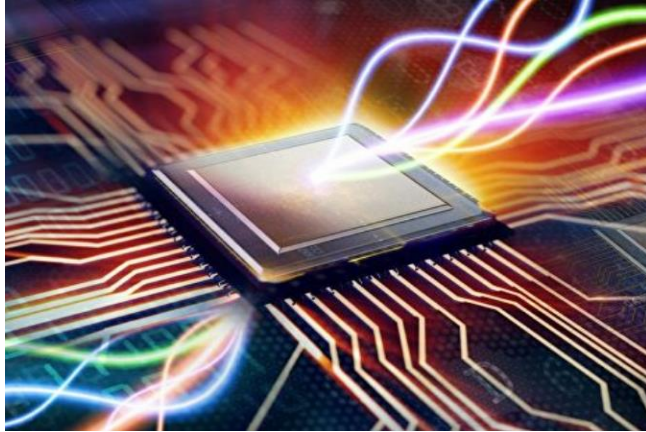
- ◆ Customized Optical Neural Network Design



- ◆ ONN On-Chip Training Algorithms



- ◆ ML-Assisted Photonic Design Automation



Photonics for **AI**



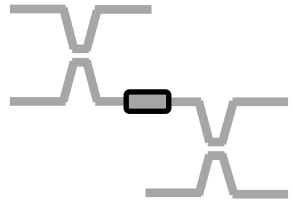
*Light-AI
Virtuous Cycle*



AI for **Photonics**

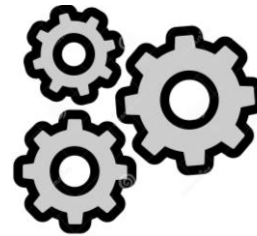
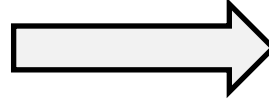


Manual Design

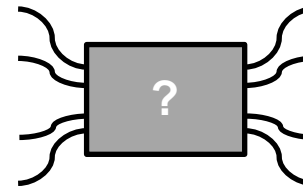


Standard Devices

AI-Enabled



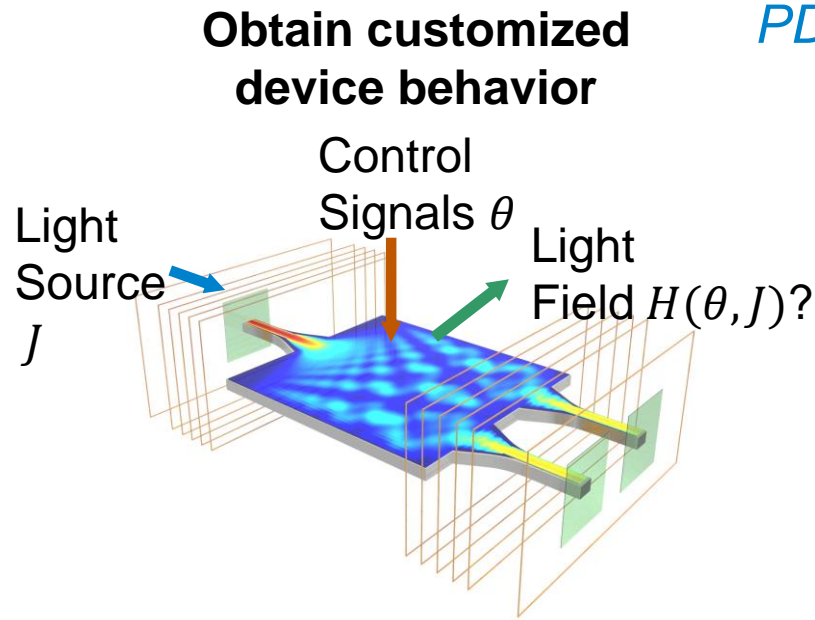
Automated Photonic
IC Design



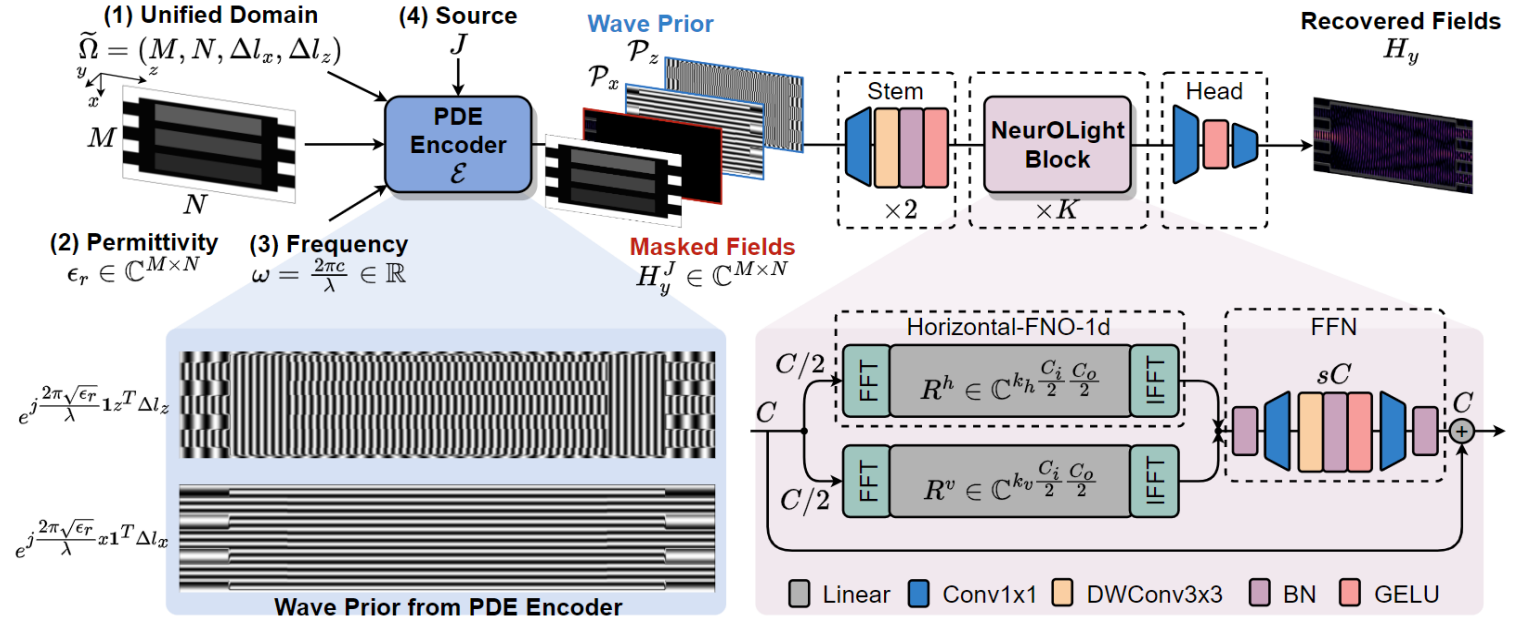
Customized Photonic
Structure

AI for Optical Simulation [NeurOLight, Gu+, NeurIPS'22]

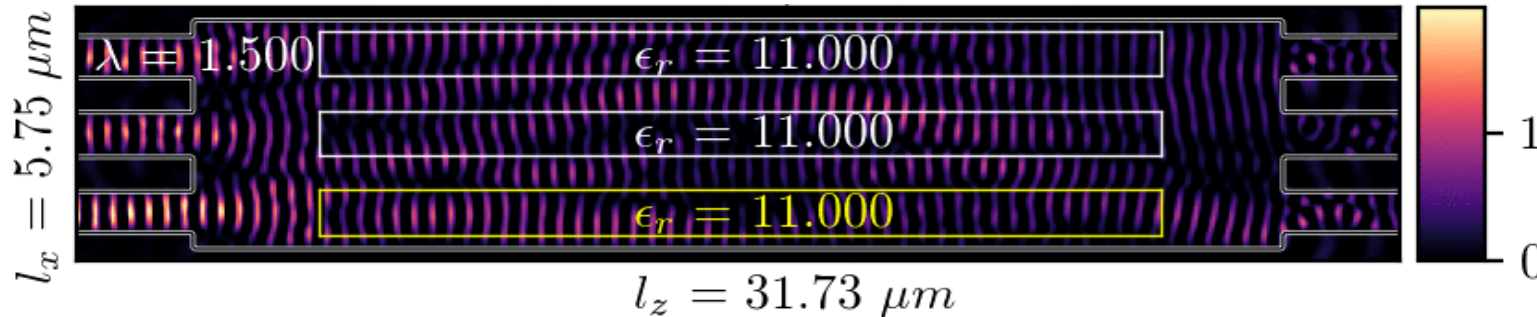
- ◆ Avoid slow simulation in the loop → ML-enabled fast Maxwell equation solving



PDE encoding + *efficient neural operator* + *physics-augmented training*



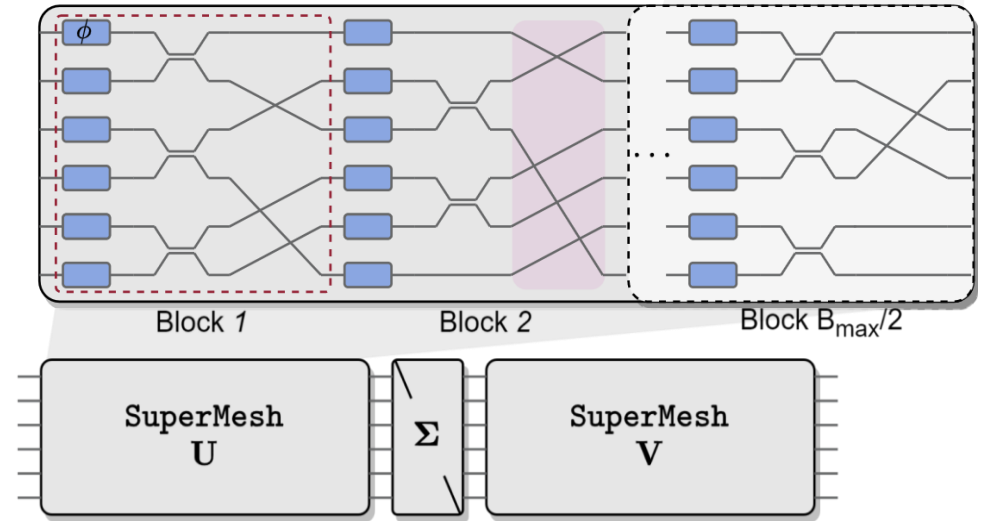
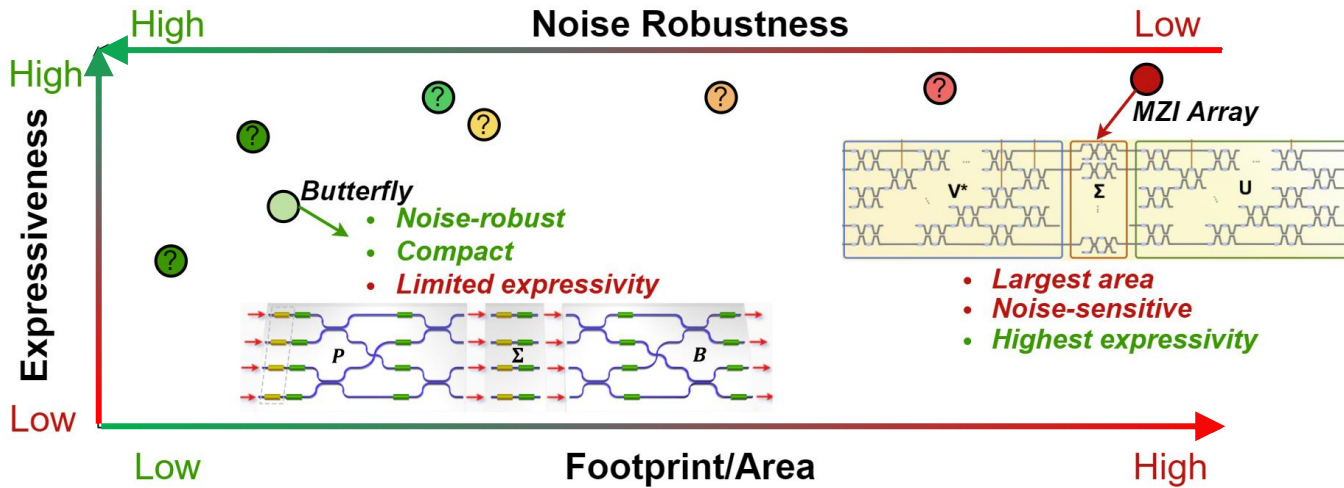
One-shot batched prediction on parametric Maxwell Eq. solutions



>200× speedup:
Fast inference (<10 ms)
120 FPS

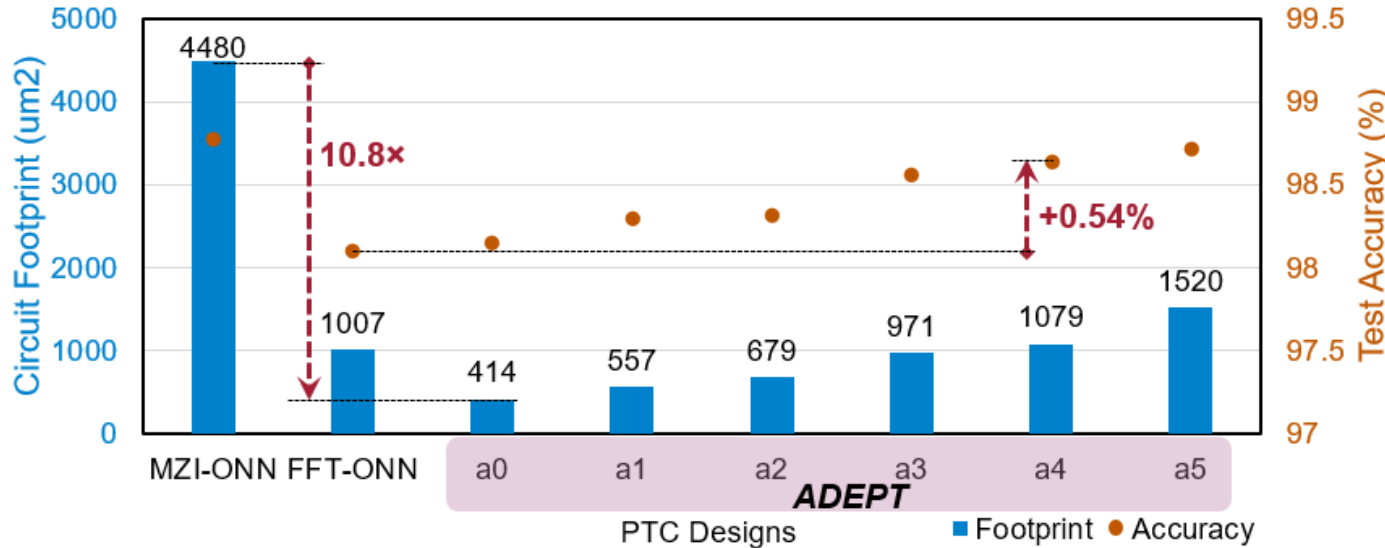
Auto-Design for Photonic Circuits [ADEPT, Gu+, DAC'22]

◆ Inefficient manual/heuristic design → Automated **circuit topology search**



Discrete → *Continuous & Differentiable*
Auto adapt to PDK and chip constraint

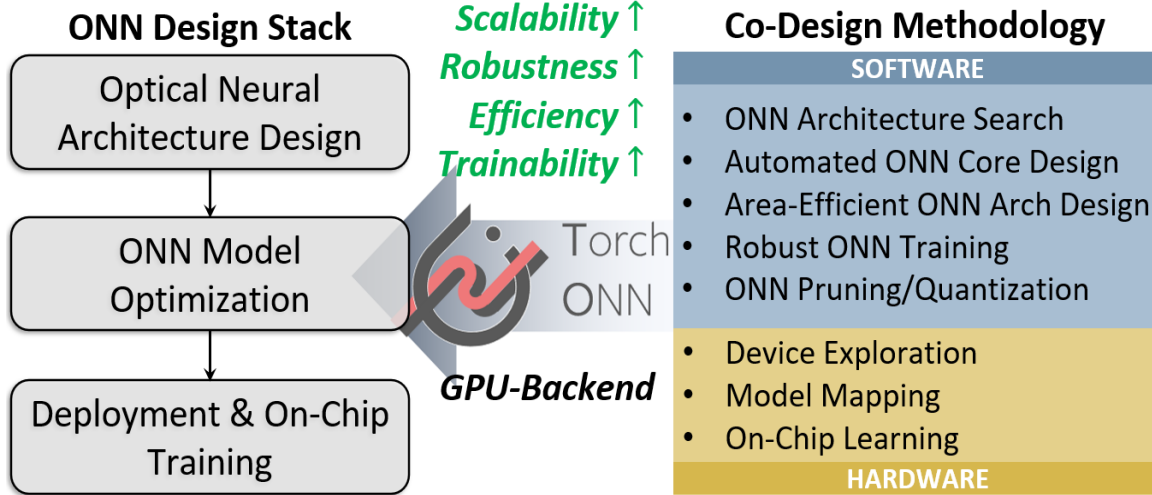
2-30× smaller
More noise-robust



The Future of Photonics ↔ AI is Bright



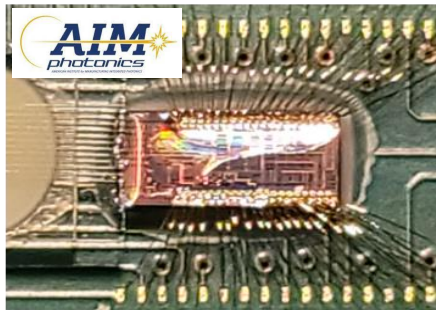
- ◆ HW/SW co-design for optical AI infer. /train + ML for optics
- ◆ Future: mem/arch, system integration, advanced app.



[JeremieMelo/pytorch-onn](https://github.com/JeremieMelo/pytorch-onn)

Photonics for AI AI for Photonics

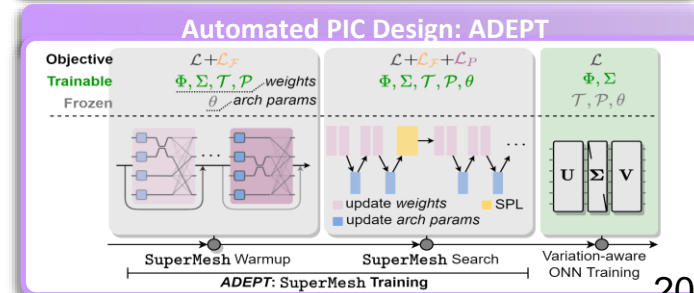
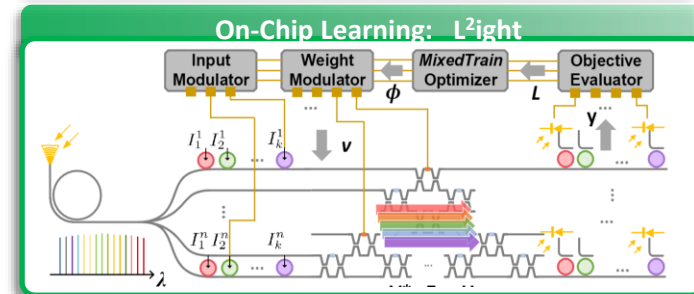
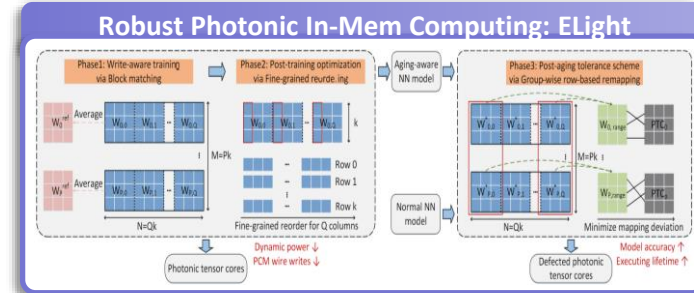
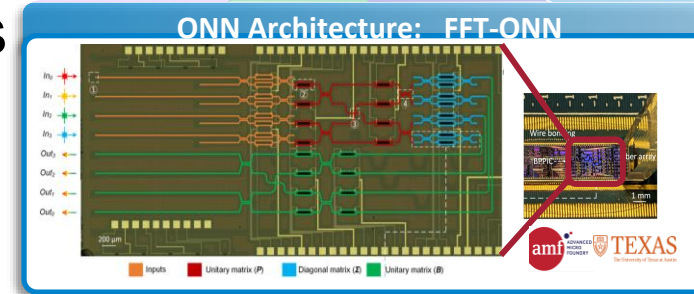
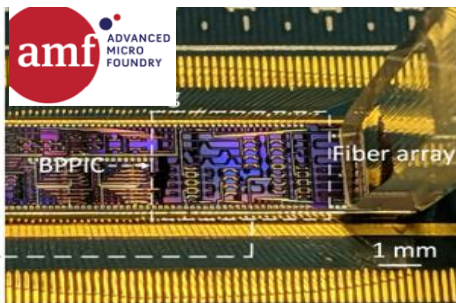
Hands-on Tutorial on TorchONN @
Design Automation Conference (DAC) July'23, Moscone Center



Acknowledgment

Contributors: Jiaqi Gu, Chenghao Feng, Hanqing Zhu, Zheng Zhao, Zhoufeng Ying, Ray T. Chen, David Z. Pan

Funding Agency: AFOSR MURI



Thank You

Q & A