





### **ADEPT-Z: Zero-Shot Automated Circuit Topology Search for Pareto-Optimal Photonic Tensor Cores**

Ziyang Jiang<sup>1</sup>, Pingchuan Ma<sup>1</sup>, Meng Zhang<sup>2</sup>, Rena Huang<sup>2</sup>, Jiaqi Gu<sup>1</sup> <sup>1</sup>Arizona State University, <sup>2</sup>Rensselaer Polytechnic Institute <sup>1</sup>School of Electrical, Computer and Energy Engineering *zjian124@asu.edu jiaqiqu@asu.edu* | *scopex-asu.qithub.io* 



### **Photonic ML Accelerators**

• Evolve from <u>electronics</u> to *integrated photonics* 



Source: Mitchell A. Nahmias, Bhavin J. Shastri, Alexander N. Tait, Thomas Ferreira de Lima and Paul R. Prucnal, "Neuromorphic photonics," Optics & Photonics News, Jan 2018.

# **PTC:** <u>Photonic</u> <u>Tensor</u> <u>Core</u>

#### **Coherent PTCs: Complex-valued Matrix Multiplication**



PTC with MMI Devices [Sci. Rep.'24, Zelaya+]

### Modular Coherent PTC Design

#### **Repeated block structures**

- Phase shifter arrays ( $\Phi's$ )
- Coupler arrays  $(\mathcal{T}'s)$
- Crossing arrays  $(\mathcal{P}'s)$

Matrix:  $\mathbf{U} = \mathcal{P}_B \mathcal{T}_B \Phi_B \dots \mathcal{P}_2 \mathcal{T}_2 \Phi_2 \mathcal{P}_1 \mathcal{T}_1 \Phi_1$ 

### **Basic Optical Components**



3

### **Manual PTC Design** $\rightarrow$ **Auto-Searched Design?**



# **Prior Work on Auto PTC Design**

- Relax as <u>continuous</u> optimization
- Gradient descent to search (similar to DARTS in <u>one-shot</u> neural architecture search)

#### Challenges in *Differentiable* ADEPT

- Limited design space: fixed slots for 2x2 coupler, <u>cannot handle multi-port couplers</u>
- Optimization difficulty: hard to handle multiple non-differentiable objectives / constraints Relaxation is <u>inaccurate and hard to converge</u>.
- High search cost: need to re-train for each search specification with <u>high training cost</u>.

### **Differentiably auto-searched PTC**

**ADEPT** [DAC'22]



### **Proposed <u>Zero-Shot</u> PTC Search: ADEPT-Z**

- Extended search space
  - <u>Arbitrary</u> placement of <u>multi-port</u> couplers
- Support multiple objectives / constraints
  - <u>Accurate</u> hardware cost modeling without surrogate model or relaxation
  - > e.g., power, area, latency,
- Efficient zero-shot search
  - > <u>Training-free</u> during search
  - > Efficient accuracy proxy



### **Extend Architecture Design Space (** $\alpha \in A$ **)**

#### ♦ ADEPT [DAC'22]

Only support 2-input couplers with fixed placements

> **Design space:** 
$$O\left(\left(\frac{K}{2} \cdot K!\right)^{B_{max}}\right)$$



#### ♦ ADEPT-Z

- > Support multi-port couplers with arbitrary placements
- > **Design space**  $\uparrow: O\left(\left(2^{K-1} \cdot K!\right)^{B_{max}}\right)$
- Dense signal interaction  $\rightarrow$  reduce circuit depth

This <u>exponential</u>, <u>discrete</u> design space cannot be searched by gradient descent ADEPT-Z uses evolutionary algorithm to explore this space

# **Overview of Zero-Shot PTC Topology Search**

### Formulation

- $\begin{array}{l} \max_{\alpha \in \mathcal{A}} \left( \text{Efficiency, Density, Accuracy} \right) \\ \text{s.t., Pmin < Power}(\alpha) < \text{Pmax} \\ \text{Amin < Area}(\alpha) < \text{Amax} \\ \tau_{min} < \text{Latency}(\alpha) < \tau_{max} \end{array} \right)$
- <u>Multi-objective</u>: NSGA-II Pareto front search
- <u>Accuracy proxy</u>: fast training-free acc. estimation
- <u>Hardware constraints</u>: posed by population filtering



One-shot Search: 6-10 hours to get a single solution

Zero-shot Search: 3 hours to get multiple pareto-optimal solutions



### **Multi-Objective Evolution: Gene Representation**

#### Compact Gene-to-Circuit mapping

- #activeBlk-DC-CR-DC-CR-....
- > DC: port count list, e.g., [2,3,1], sum to K
- > CR: output port indices, e.g., [3,4,1,2,6,5]

Compact representation of multi-port couplers Easy to manipulate in evolution



9

### **Proposed Zero-Shot PTC Topology Search**



# **Exploration in Large Design Space: Mutation**

- Customized mutation operators: <u>legal move</u> + <u>sufficient change</u> + <u>good stability</u>
- Step 1 Global Search: Global Jump + Local Jump

Step 2 Local Search: Local Jump only

Balance between Exploration and Exploitation



# **Exploitation in Large Design Space: Crossover**

- Guarantee gene legality
- Maintain good property of parents (interpolation)
- DC crossover
  - Genes → Boarder-Aligned Segments
  - Swap aligned segments randomly
- CR crossover
  - Select complementary subset of indices from parents
  - Cross-insert and preserve the relative order in parent genes
- Block crossover
  - Swap two active blocks randomly



### **Proposed Zero-Shot PTC Topology Search**



### **Explicit Hardware Constraints in the Loop**





Consider actual chip layout, practical spacing

More **accurate** estimation than summing optical device **footprints** 



### **Explicit Hardware Constraints in the Loop**



Copyright © 2022 Arizona Board of Regents

### **Proposed Zero-Shot PTC Topology Search**



### Hardware Cost & Training-Free Acc. Estimation

- Energy efficiency (TOPS):  $EE(\alpha) = 2K^2/(Power(\alpha) * Latency(\alpha))$
- <u>Compute density</u> (TOPS/mm<sup>2</sup>):  $CD(\alpha) = 2K^2/(Area(\alpha) * Latency(\alpha))$
- <u>Accuracy</u>: Find 3 easy-tocompute accuracy indicators
  - High rank correlation with true test accuracy (spearman >0.94)
  - Comprehensive accuracy Score S<sub>1</sub>
    w/ regressed coefficients.



ZiCo [Li, ICLR'23]: grad property (trainability, generalization)

Training v.s. Accuracy Estimation >30min  $\rightarrow$  **10 Seconds** 

### **Auto-Searched Pareto-Optimal PTCs on GF PDK**

- Search 16×16 PTC (2/8-port DC) on GlobalFoundries PDK
- Pareto front covers all three manual designs
- Better than Random search
- One-time search, 40 Pareto-optimal solutions within 2.7 hours
  - > 100× faster than ADEPT (6-10 hours per one-shot differentiable search)



18

### **Auto-Searched Pareto-Optimal PTCs on GF PDK**

- Pareto-optimal PTC ADEPT-Z-a0 with core size of 8, 16, 32
- 2.47× higher accuracy-weighted area-energy efficiency product (AAEE) than MZI [Nat. Photon'17] and MMI [Sci. Rep.'24] arrays
- 1.03× higher AAEE than Butterfly Mesh [ASP-DAC'20]



### **Customized PDK Adaptation**

- Adapt to customized foundry PDK (with smaller PSs than GF, different PDs)
- 8.26× higher AAEE than MZI [Nat. Photon'17] and MMI [Sci. Rep.'24] arrays
- 1.04× higher AAEE than Butterfly Mesh [ASP-DAC'20]

#### 4 variants (small-to-large) w/ customized foundry PDK

Metrics	MZI [1]	Butterfly [4]	MMI [10]	ADEPT-Z-a0	ADEPT-Z-a1	ADEPT-Z-a2	ADEPT-Z-a3
Area(O+E)	7.51+0.33	$1.25 \pm 0.33$	16.90 + 0.33	1.19 + 0.33	$1.22 \pm 0.33$	1.21+0.33	1.71 + 0.33
Power	223.59	219.93	219.45	218.37	218.54	218.30	218.93
Latency	100.00	100.00	100.00	100.00	100.00	100.00	100.00
Accuracy	98.74	98.16	<b>98.5</b> 8	97.67	98.02	98.08	98.18
CD	0.653	3.242	0.297	3.370	3.301	3.322	2.507
EE	22.899	23.280	23.330	23.446	23.429	23.454	23.387
AEE	2.919	14.744	1.354	15.434	15.105	15.215	11.450
AAEE	2.882	14.473	1.336	15.074	14.806	14.923	11.242

### **Generalize to Different Tasks/Models**

- ◆ Search PTC on 2-layer CNN + MNIST → Train ONN on different tasks/models
- Average of 1.6× higher AAEE than manual PTC designs

					<u></u>
Model	Dataset	MZI [1]	Butterfly [4]	MMI [10]	ADEPT-Z-a0
CNN	FMNIST	0.472	0.852	0.432	0.853
VGG8	CIFAR10	0.429	0.744	0.382	0.769
ResNet20	SVHN	0.491	0.884	0.445	0.890
					·/

**AAEE**↑

- Future direction in CAD for optical computing
  - Search beyond this layer-wise coherent PTC design template
  - Architecture, interconnect, and PTC topology co-search





Open-Source TorchONN Toolchain

Automating optical AI hardware design toward productivity

# Thank you! Q & A?

#### ADEPT-Z: Zero-Shot Automated Circuit Topology Search for Pareto-Optimal Photonic Tensor Cores

Ziyang Jiang<sup>1</sup>, Pingchuan Ma<sup>1</sup>, Meng Zhang<sup>2</sup>, Rena Huang<sup>2</sup>, Jiaqi Gu<sup>1\*</sup> <sup>1</sup>Arizona State University, <sup>2</sup>Rensselaer Polytechnic Institute

Contact Jiaqi Gu: jiaqigu@asu.edu



ASU Center for Semiconductor Microelectronics (ACME)



