

# Towards Hardware-Efficient Optical Neural Networks: Beyond FFT Architecture via Joint Learnability

Jiaqi Gu *Student Member, IEEE*, Zheng Zhao, Chenghao Feng, Zhoufeng Ying, Mingjie Liu, Ray T. Chen *Fellow, IEEE*, and David Z. Pan *Fellow, IEEE*

**Abstract**— As a promising neuromorphic framework, the optical neural network (ONN) demonstrates ultra-high inference speed with low energy consumption. However, the previous ONN architectures have high area overhead which limits their practicality. In this paper, we propose an area-efficient ONN architecture based on structured neural networks, leveraging optical fast Fourier transform for efficient computation. A two-phase software training flow with structured pruning is proposed to further reduce the optical component utilization. Experimental results demonstrate that the proposed architecture can achieve  $2.2\sim3.7\times$  area cost improvement compared with the previous singular value decomposition-based architecture with comparable inference accuracy. A novel optical microdisk-based convolutional neural network architecture with joint learnability is proposed as an extension to move beyond Fourier transform and multi-layer perception, enabling hardware-aware ONN design space exploration with lower area cost, higher power efficiency, and better noise-robustness.

## I. INTRODUCTION

Deep neural networks (DNNs) have demonstrated superior performance in a variety of intelligent tasks, for example convolutional neural networks on image classification [1] and recurrent neural networks on language translation [2]. Multi-layer perceptrons (MLPs) are among the most fundamental components in modern DNNs, which are typically used as regression layers, classifiers, embedding layers, and attention layers, etc. However, it becomes challenging for traditional electrical digital von Neumann schemes to support escalating computation demands owing to speed and energy inefficiency [3]–[7]. To resolve this issue, significant efforts have been made on hardware design of neuromorphic computing frameworks to improve the computational speed of neural networks, such as electronic architectures [8]–[10] and photonic architectures [11]–[15]. Among extensive neuromorphic computing systems, optical neural networks (ONNs) distinguish themselves by ultra-high bandwidth, ultra-low latency, and near-zero energy consumption. Even though ONNs are currently not competitive in terms of area cost, they still offer a promising alternative approach to microelectronic implementations given the above advantages.

Recently, several works demonstrated that MLP inference can be efficiently performed at the speed of light with optical components, e.g., spike processing [11] and reservoir computing [16]. They claimed a photodetection rate over 100 GHz in photonic networks, with near-zero energy consumption if passive photonic components are used [17]. Based on matrix singular value decomposition (SVD) and unitary matrix parametrization [18], [19], Shen *et al.* [3] designed and fabricated a fully optical neural network that achieves an MLP

The preliminary version has been presented at the ACM/IEEE Asian and South Pacific Design Automation Conference (ASP-DAC) in 2020. This work was supported in part by the Multidisciplinary University Research Initiative (MURI) program through the Air Force Office of Scientific Research (AFOSR), contract No. FA 9550-17-1-0071, monitored by Dr. Gernot S. Pomrenke.

J. Gu, C. Feng, M. Liu, R. T. Chen and D. Z. Pan are with the Department of Electrical and Computer Engineering, The University of Texas at Austin, TX, USA.

Z. Zhao is with Synopsys Inc., CA, USA.

Z. Ying is with Alpine Optoelectronics, CA, USA.

with Mach-zehnder interferometer (MZI) arrays. Once the weight matrices in the MLP are trained and decomposed, thermo-optic phase shifters on the arms of MZIs can be set up accordingly. Since the weight matrices are fixed after training, this fully optical neural network can be completely passive, thus minimizes the total energy consumption. However, this SVD-based architecture is limited by high photonic component utilization and area cost. Considering a single fully-connected layer with an  $m \times n$  weight matrix, the SVD-based ONN architecture requires  $\mathcal{O}(m^2 + n^2)$  MZIs for implementation. Another work [20] proposed a slimmed ONN architecture (TΣU) based on the previous one [3], which substitutes one of the unitary blocks with a sparse tree network. However, its area cost improvement is limited. Therefore, this high hardware complexity of the SVD-based ONN architecture has become the bottleneck of its hardware implementation.

In addition to hardware implementation, recent advances in neural architecture design and network compression techniques have shown significant reduction in computational cost. For example, structured neural networks (SNNs) [21] were proposed to significantly reduce computational complexity and thus, become amenable to hardware. Besides, network pruning offers another powerful approach to slimming down neural networks by cutting off insignificant neuron connections. While non-structured pruning [22] produces random neuron sparsity, group sparsity regularization [23] and structured pruning [9] can lead to better network regularity and hardware efficiency. However, readily-available pruning techniques are rather challenging to be applied to the SVD-based architecture due to some issues, such as accuracy degradation and hardware irregularity. The gap between hardware-aware pruning and the SVD-based architecture gives another motivation for a pruning-friendly ONN architecture.

In this paper, we propose a new ONN architecture that improves area efficiency over previous ONN architectures. It leverages optical fast Fourier transform (OFFT) and its inverse (OIFFT) to implement structured neural networks, achieving lower optical component utilization. It also enables the application of structured pruning given its architectural regularity. The proposed architecture partitions the weight matrices into block-circulant matrices [24] and efficiently performs circulant matrix multiplication through OFFT/OIFFT. We also adopt a two-phase software training flow with structured pruning to further reduce photonic component utilization while maintaining comparable inference accuracy to previous ONN architectures. We extend this architecture to a hardware-efficient optical convolutional neural network design with joint learnability, and demonstrate its superior power efficiency and noise-robustness compared with Fourier transform based design. The main contributions of this work are as follows:

- We propose a novel, area-efficient optical neural network architecture with OFFT/OIFFT, and exploit a two-phase software training flow with structured pruning to learn hardware-friendly sparse neural networks that directly eliminate part of OFFT/OIFFT modules for further area efficiency improvement.
- We experimentally show that pruning is challenging to be ap-

plied to previous ONN architectures due to accuracy loss and retrainability issues.

- We experimentally demonstrate that our proposed architecture can lead to an area saving of  $2.2\sim3.7\times$  compared with the previous SVD-based ONN architecture, with negligible inference accuracy loss.
- We extend our ASP-DAC version of ONN architecture [25] to a novel design for microdisk-based frequency-domain optical convolutional neural networks with high parallelism.
- We propose a trainable frequency-domain transform structure and demonstrate it can be pruned with high sparsity and outperforms traditional Fourier transform with less component count, higher power efficiency, and better noise-robustness.

The remainder of this paper is organized as follows. Section II introduces the background knowledge for our proposed architecture. Section III demonstrates the challenges to apply pruning to SVD-based architectures. Section IV presents details about the proposed ONN architecture and software pruning flow. Section V analytically compares our hardware utilization with the SVD-based architecture. Section VI demonstrates an extension to optical convolutional neural network with trainable transform structures. Section VII reports the experimental results for our proposed ONN architecture and its CNN extension, followed by the conclusion in Section VIII.

## II. PRELIMINARIES

In this section, we introduce the background knowledge for our proposed architecture. We discuss principles of circulant matrix representation and its fast computation algorithms in Section II-A and illustrate structured pruning techniques with Group Lasso regularization in Section II-B.

### A. FFT-based Circulant Matrix Computation

Unlike the SVD-based ONNs which focus on classical MLPs, our proposed architecture is based on structured neural networks (SNNs) with circulant matrix representation. SNNs are a class of neural networks that are specially designed for computational complexity reduction, whose weight matrices are regularized using the composition of structured sub-matrices [21]. Among all structured matrices, circulant matrices are often preferred in recent SNN designs.

As an example, we show an  $n \times n$  circulant matrix  $\mathbf{W}$  as follows,

$$\begin{bmatrix} w_0 & w_{n-1} & \cdots & w_1 \\ w_1 & w_0 & \cdots & w_2 \\ \vdots & \vdots & \ddots & \vdots \\ w_{n-1} & w_{n-2} & \cdots & w_0 \end{bmatrix}.$$

The first column vector  $\mathbf{w} = [w_0, w_1, \dots, w_{n-1}]^T$  represents all independent parameters in  $\mathbf{W}$ , and other columns are just its circulation.

According to [24], circulant matrix-vector multiplication can be efficiently calculated through fast Fourier transform. Specifically, given an  $n \times n$  circulant matrix  $\mathbf{W}$  and a length- $n$  vector  $\mathbf{x}$ ,  $\mathbf{y} = \mathbf{W}\mathbf{x}$  can be efficiently performed with  $\mathcal{O}(n \log n)$  complexity as,

$$\mathbf{y} = \mathcal{F}^{-1}(\mathcal{F}(\mathbf{w}) \odot \mathcal{F}(\mathbf{x})), \quad (1)$$

where  $\mathcal{F}(\cdot)$  represents  $n$ -point real-to-complex fast Fourier transform (FFT),  $\mathcal{F}^{-1}(\cdot)$  represents its inverse (IFFT), and  $\odot$  represents complex vector element-wise multiplication.

SNNs benefit from high computational efficiency while maintaining comparable model expressivity to classical NNs. Theoretical analysis [26] shows that SNNs can approximate arbitrary continuous functions with arbitrary accuracy given enough parameters, and are also capable of achieving the identical error bound to that of classical NNs.

Therefore, based on SNNs with circulant matrix representation, the proposed architecture features low computational complexity and comparable model expressivity.

### B. Structured Pruning with Group Lasso Penalty

The proposed ONN architecture enables the application of structured pruning to further save optical components while maintaining accuracy and structural regularity. Structured pruning trims the neuron connections in NNs to mitigate computational complexity. Unlike  $\ell_1$  or  $\ell_2$  norm regularization, which produces arbitrarily-appearing zero elements, structured pruning with Group Lasso regularization [9], [27] leads to zero entries in groups. This coarse-grained sparsity is more friendly to hardware implementation than non-structured sparsity. The formulation of Group Lasso regularization term is given as follows,

$$L_{GL} = \sum_{g=0}^G \sqrt{1/p_g} \|\beta_g\|_2, \quad (2)$$

where  $G$  is the total number of parameter groups,  $\beta_g$  is the parameter vector in the  $g$ -th group,  $\|\cdot\|_2$  represents  $\ell_2$  norm,  $p_g$  represents the vector length of  $\beta_g$ , which accounts for the varying group sizes. Intuitively, the  $\ell_2$  norm penalty  $\|\beta_g\|_2$  encourages all elements in the  $g$ -th group to converge to 0, and the group-wise summation operation is equivalent to group-level  $\ell_1$  norm regularization, which contributes to the coarse-grained sparsity. Leveraging the structured pruning together with Group Lasso regularization, our proposed architecture can save even more photonic components.

## III. CHALLENGES IN PRUNING SVD-BASED ARCHITECTURE

In this section, we demonstrate that the network pruning is experimentally challenging in the SVD-based architecture. As far as we know, it is hard to find any pruning method that can be directly applied to sparsifying MZI arrays.

In the SVD-based architecture, an  $m \times n$  weight matrix  $\mathbf{W}$  can be decomposed into  $\mathbf{W} = \mathbf{U}\Sigma\mathbf{V}$  using singular value decomposition. Unitary matrices  $\mathbf{U}$  and  $\mathbf{V}$  can be further parametrized [18] into the product of planar rotation matrices  $\mathbf{U} = \mathbf{D} \cdot \prod_{i=m}^2 \prod_{j=1}^{i-1} \mathbf{R}_{ij}$ , where  $\mathbf{D}$  is a diagonal matrix and each unitary rotation  $\mathbf{R}_{ij}$  can be represented by an angle or phase  $\phi$ . Each unitary rotation matrix with phase  $\phi$  can be implemented with an MZI. We denote all phases after parametrization as  $\Phi$ . Phases with particular values, i.e.,  $0, \pi/2, \pi$ , and  $-\pi/2$ , can physically eliminate the use of the corresponding MZIs. We refer to these particular phases as sparse phases. One of the methods to perform pruning is to train a sparse weight matrix, but the sparsity can barely maintain after decomposition and parametrization. Another straight-forward method is post-training phase pruning. It directly clamps sparse phases but could cause significant accuracy degradation due to its unretrainability.

We experimentally illustrate the correlation between inference accuracy and phase sparsity. Phase sparsity for MZI-based ONN is defined as the percentage of prunable phases in all phases, i.e.,  $|\{\phi | \phi = 0, \pi/2, \pi, -\pi/2\}| / |\Phi|$ . Concretely, we clamp  $\Phi$  with a threshold  $\epsilon$  to get  $\widehat{\Phi}$ . Then we evaluate the inference accuracy with reconstructed weight matrix  $\widehat{\mathbf{W}}$ . Figure 1 shows that, in different network configurations, on average no more than 15% phases can be pruned to achieve negligible ( $\sim 0.5\%$ ) absolute accuracy degradation. With over 20% phases being pruned, its model expressivity will be severely harmed with significant accuracy loss ( $> 1\%$ ). This accuracy loss partially attributes to the difficulty of retraining the weight matrices while maintaining phase sparsity. Another challenge derives from the hardware irregularity caused by non-structured phase pruning, which further limits the area improvement it can achieve.

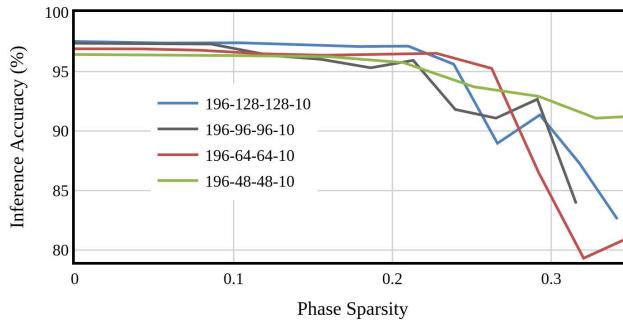


Fig. 1: Correlation between phase sparsity and inference accuracy with different network configurations based on MNIST [28] dataset. Phase sparsity indicates the proportion of prunable phases in the MZI-based ONN. 196-128-128-10 represents an MLP configuration with 196 inputs ( $16 \times 16$ ), two hidden layers with 128 neurons in each layer, and 10 output logits for the last layer.

The above limitations also apply to the TΣU-based architecture [20] as it has a similar architectural design to the SVD-based one. This unsatisfying incompatibility between previous ONN architectures and pruning techniques offers a strong motivation for us to propose a new architecture to better leverage pruning techniques.

#### IV. PROPOSED ARCHITECTURE

In this section, we will discuss details about the proposed architecture and pruning method. In the first part, we illustrate five stages of our proposed architecture. In the second part, we focus on the two-phase software training flow with structured pruning.

##### A. Proposed Architecture

Based on structured neural networks, our proposed architecture implements a structured version of MLPs with circulant matrix representation. A single layer in the proposed architecture performs linear transformation via block-circulant matrix multiplication  $\mathbf{y} = \mathbf{W}\mathbf{x}$ . Consider an  $n$ -input,  $m$ -output layer, the weight matrix  $\mathbf{W} \in \mathbb{R}^{m \times n}$  is partitioned into  $p \times q$  sub-matrices, each being a  $k \times k$  circulant matrix. To perform tiled matrix multiplication, the input  $\mathbf{x}$  is also partitioned into  $q$  segments  $\mathbf{x} = (\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{q-1})$ . Thus  $\mathbf{y} = \mathbf{W}\mathbf{x}$  can be performed in a tiled way,

$$\mathbf{y} = \begin{pmatrix} \mathbf{y}_0 \\ \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_{p-1} \end{pmatrix} = \begin{pmatrix} \sum_{j=0}^{q-1} \mathbf{W}_{0j} \mathbf{x}_j \\ \sum_{j=0}^{q-1} \mathbf{W}_{1j} \mathbf{x}_j \\ \vdots \\ \sum_{j=0}^{q-1} \mathbf{W}_{p-1j} \mathbf{x}_j \end{pmatrix}. \quad (3)$$

The  $i$ th segment  $\mathbf{y}_i = \sum_{j=0}^{q-1} \mathbf{W}_{ij} \mathbf{x}_j$  is the accumulation of  $q$  independent circulant matrix multiplications. Each  $\mathbf{W}_{ij} \mathbf{x}_j$  can be efficiently calculated using the fast computation algorithm mentioned in Eq. (1). Based on the aforementioned equations, we realize block-circulant matrix multiplication  $\mathbf{y} = \mathbf{W}\mathbf{x}$  in five stages: 1) Splitter tree (ST) stage to split input optical signals for reuse; 2) OFFT stage to calculate  $\mathcal{F}(\mathbf{x})$ ; 3) element-wise multiplication (EM) stage to calculate  $\mathcal{F}(\mathbf{w}_{ij}) \odot \mathcal{F}(\mathbf{x}_j)$  as described in Eq. (1); 4) OIFFT stage to calculate  $\mathcal{F}^{-1}(\cdot)$ ; 5) combiner tree (CT) stage to accumulate partial multiplications to form the final results.  $\mathcal{F}(\mathbf{w}_{ij})$  can be precomputed and encoded into optical components, thus there is no extra stage to physically perform it. The schematic diagram of our proposed architecture is shown in Fig. 2. Details of the above five stages will be discussed in the rest of this section.

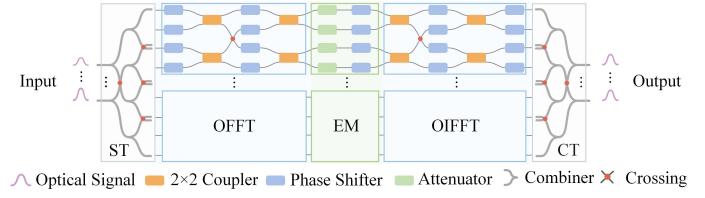


Fig. 2: Schematic diagram of a single layer of the proposed architecture. All adjacent phase shifters on the same waveguide are already merged into one phase shifter.

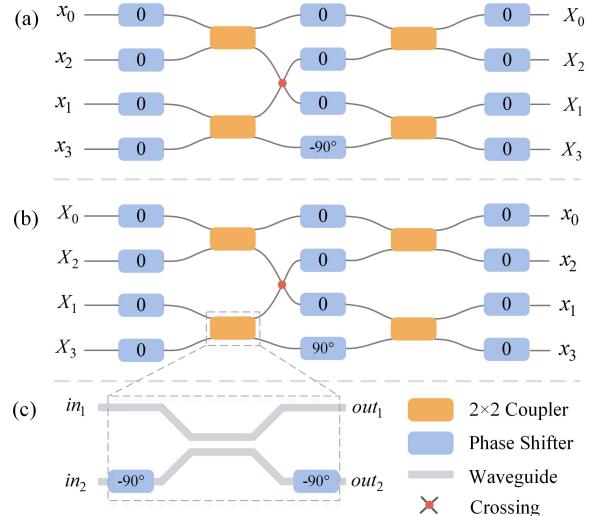


Fig. 3: Schematics of (a) 4-point OFFT, (b) 4-point OIFFT, and (c) 2  $\times$  2 coupler. Note that phase shifters shown above are not merged for structural completeness consideration.

1) OFFT/OIFFT Stages: To better model the optical components used to implement the OFFT/OIFFT stages, we introduce a unitary FFT as,

$$\mathbf{X}_k = \frac{1}{\sqrt{N}} \sum_{n=0}^{N-1} \mathbf{x}_n e^{-i \frac{2\pi k n}{N}} \quad k = 0, 1, \dots, N-1. \quad (4)$$

We denote this special operation as  $\widehat{\mathcal{F}}(\cdot)$  and its inverse as  $\widehat{\mathcal{F}}^{-1}(\cdot)$ , to distinguish from the original FFT/IFFT operations. Equivalently, we re-write the circulant matrix multiplication with the above new operations,

$$\mathbf{y} = \widehat{\mathcal{F}}^{-1}(\mathcal{F}(\mathbf{w}) \odot \widehat{\mathcal{F}}(\mathbf{x})). \quad (5)$$

This unitary FFT operation can be realized with optical components. We first give a simple example for the optical implementation of a 2-point unitary FFT. As shown in Eq. (7), the transformation matrix of a 2-point unitary FFT can be decomposed into three transform matrices. They can be directly mapped to a 3-dB directional coupler with two  $-\pi/2$  phase shifters on its lower input/output ports. The transfer matrix of a 50/50 optical directional coupler is given by,

$$\frac{1}{\sqrt{2}} \begin{pmatrix} 1 & j \\ j & 1 \end{pmatrix}. \quad (6)$$

The transfer function of a phase shifter is  $out = in \cdot e^{j\phi}$ . For brevity, we refer to this cascaded structure as a 2  $\times$  2 coupler, which is shown

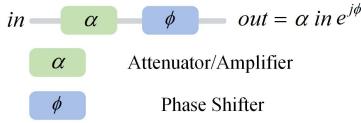


Fig. 4: Complex number multiplication realized by cascaded attenuator/amplifier and phase shifter.

in Fig. 3(c).

$$\begin{pmatrix} \text{out}_1 \\ \text{out}_2 \end{pmatrix} = \frac{1}{\sqrt{2}} \begin{pmatrix} \text{in}_1 + \text{in}_2 \\ \text{in}_1 - \text{in}_2 \end{pmatrix} \\ = \underbrace{\begin{pmatrix} 1 & 0 \\ 0 & -j \end{pmatrix}}_{\text{output phase shifter}} \underbrace{\frac{1}{\sqrt{2}} \begin{pmatrix} 1 & j \\ j & 1 \end{pmatrix}}_{\text{directional coupler}} \underbrace{\begin{pmatrix} 1 & 0 \\ 0 & -j \end{pmatrix}}_{\text{input phase shifter}} \begin{pmatrix} \text{in}_1 \\ \text{in}_2 \end{pmatrix} \quad (7)$$

Based on  $2 \times 2$  couplers and phase shifters, larger-sized OFFT/OIFFT can be constructed with a butterfly structure. The schematics of a simple 4-point OFFT and OIFFT are shown in Fig. 3(a) and Fig. 3(b). Extra 0-degree phase shifters are inserted for phase tuning purpose.

This butterfly-structured OFFT may have scalability issues because the number of waveguide crossings (CR) will increase rapidly when the number of point gets larger. However, this unsatisfying scalability will not limit our proposed architecture for two reasons. First, only small values of  $k$ , e.g., 2, 4, 8, will be adopted to balance hardware efficiency and model expressivity. Second, input and output sequences can be reordered to avoid unnecessary waveguide crossings, as shown in Fig. 3.

2) *EM Stage*: In the EM stage, complex vector element-wise multiplications will be performed in the Fourier domain as  $\alpha e^\phi \cdot I_{in} e^{\phi_{in}} = \alpha I_{in} e^{\phi_{in} + \phi}$ , where  $I_{in}$  and  $\phi_{in}$  are magnitude and phase of input Fourier light signals respectively. Leveraging the polarization of light, we use optical attenuators (AT) or amplification materials/optical on-chip amplifiers with a scaling factor  $\alpha$  to realize modulus multiplication  $\alpha \cdot I_{in}$  and phase shifters with  $\phi$  phase shift for argument addition  $e^{j(\phi + \phi_{in})}$ , which is shown in Fig. 4.

3) *ST/CT Stage*: We introduce tree-structured splitter/combiner networks to realize input signal splitting and output signal accumulation, respectively. To reuse input segments  $\mathbf{x}_j$  in multiple blocks, optical splitters (SP) are used to split optical signals. Similarly, to accumulate partial multiplication results, i.e.,  $\mathbf{y}_i = \sum_{j=0}^{q-1} \mathbf{W}_{ij} \mathbf{x}_j$ , we adopt optical combiners (CB) for signal addition. Given that optical splitters can be realized by using combiners in an inverted direction, we will focus on the combiner tree structure for brevity.

The transfer function of an  $N$ -to-1 optical combiner is,

$$\text{out} = \frac{1}{\sqrt{N}} \sum_{l=0}^{N-1} \text{in}_l. \quad (8)$$

Accumulating  $q$  length- $k$  vectors by simply using  $k$   $q$ -to-1 combiners introduces a huge number of waveguide crossings which may cause intractable implementation difficulty. Also, combiners with more than two ports are still challenging for manufacturing. In order to alleviate this problem, we adopt a tree-structured combiner network, shown in Fig. 5. This combiner tree consists of  $k(k-1)$  combiners and reduces the number of waveguide crossings to  $k(k-1)(q-1)/2$ . Given that combiners will cause optical intensity loss by a factor of  $1/\sqrt{N}$  as shown in Eq. (8), we assume there will be optical amplifiers added to the end to compensate this loss.

In terms of cascading multiple layers, our proposed FFT-based MLP is fully optical, such that the output optical signals can be

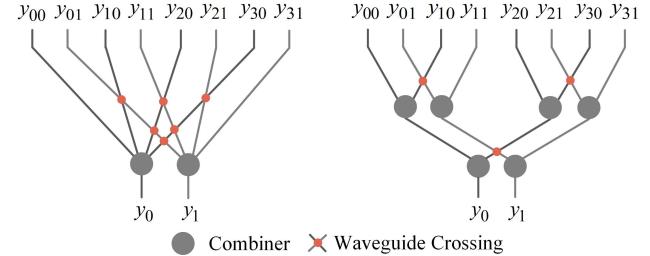


Fig. 5: Comparison between direct combining (left) and combiner tree (right) with 4 length-2 vectors accumulated.

directly fed into the next layer without optical-electrical-optical (O-E-O) conversion. At the end of the last layer, photo-detection is used for signal readout, and the phase information of the outputs are removed, which can be fully modeled during our training process without causing any accuracy loss.

### B. Two-phase Training Flow with Structured Pruning

Structured pruning can be applied to our proposed architecture during training given its architectural regularity. As described in Alg. 1, we exploit a two-phase software training flow with structured pruning to train a more compact NN model with fewer redundancies and negligible accuracy loss. Lines 2-4 perform the first initial training phase with Group Lasso regularization term added to our loss function.

$$\mathbf{L} = \mathbf{L}_{base} + \lambda \mathbf{L}_{GL}, \quad (9)$$

where  $\mathbf{L}_{base}$  is the basic loss function, e.g., cross-entropy loss if targeted at classification tasks, and  $\lambda$  is a hyper-parameter used to weigh the regularization term  $\mathbf{L}_{GL}$  given in Eq. 2. The initial training phase explores a good local minimum in the full parameter space to get rough convergence. This is designed to provide a good initial model for the subsequent pruning. Line 5 enters the structured pruning phase. The pruning mask  $\mathbf{M}$  is generated to mark  $w_{ij}$  whose  $\ell_2$  norm falls below a threshold  $T$ . Those marked weight groups will be forced to zero. Hence, the corresponding hardware modules can be completely eliminated. As training and pruning are alternately performed, the network sparsity will incrementally improve. Line 12 applies a smooth function, e.g., polynomial or Tanh, to gradually increase pruning threshold to avoid accuracy degradation caused by aggressive pruning.

### V. THEORETICAL ANALYSIS ON PROPOSED ARCHITECTURE

In this section, we analyze the hardware utilization and compare with previous architectures.

We derive a theoretical estimation of hardware utilization of the proposed architecture, the SVD-based architecture [3], and the slimmed  $T\Sigma U$ -based architecture [20]. By comparing the hardware component utilization, we show that theoretically our proposed architecture costs fewer optical components than the SVD-based architecture and  $T\Sigma U$ -based architecture. The comparison results are summarized in Table I for clear demonstration.

For simplicity, we convert all area-costly components, i.e.,  $2 \times 2$  couplers, MZIs, and attenuators, to 3-dB directional couplers (DCs) and phase shifters (PSs). Specifically, one  $2 \times 2$  coupler can be taken as one DC and two PSs, and one MZI can be taken as two DCs and one PS. Since an attenuator can be achieved by a single-input directional coupler with appropriate transfer factor, we count one attenuator as one DC.

**Algorithm 1** Two-Phase Training Flow with Structured Pruning

```

Input: Initial parameter  $\mathbf{w}^0 \in \mathbb{R}^{p \times q \times k}$ , pruning threshold  $T$ , initial training timestep  $t_{init}$ , and learning rate  $\alpha$ ;
Output: Converged parameter  $\mathbf{w}^t$  and a pruning mask  $M \in \mathbb{Z}^{p \times q}$ ;
1:  $M \leftarrow 1$  ▷ Initialize pruning mask to all 1
2: for  $t \leftarrow 1, \dots, t_{init}$  do ▷ Phase 1: Initial training
3:    $L^t(\mathbf{w}^{t-1}) \leftarrow L_{base}^t(\mathbf{w}^{t-1}) + \lambda \cdot L_{GL}^t(\mathbf{w}^{t-1})$ 
4:    $\mathbf{w}^t \leftarrow \mathbf{w}^{t-1} - \alpha \cdot \nabla_{\mathbf{w}} L^t(\mathbf{w}^{t-1})$ 
5: end for
6: while  $\mathbf{w}^t$  not converged do ▷ Phase 2: Structured pruning
7:   for all  $\mathbf{w}_{i,j}^{t-1} \in \mathbf{w}^{t-1}$  do
8:     if  $\|\mathbf{w}_{i,j}^{t-1}\|_2 < T$  then
9:        $M[i, j] \leftarrow 0$  ▷ Update pruning mask
10:    end if
11:   end for
12:   ApplyDropMask( $M, \mathbf{w}^{t-1}$ )
13:    $L^t(\mathbf{w}^{t-1}) \leftarrow L_{base}^t(\mathbf{w}^{t-1}) + \lambda \cdot L_{GL}^t(\mathbf{w}^{t-1})$ 
14:    $\mathbf{w}^t \leftarrow \mathbf{w}^{t-1} - \alpha \cdot \nabla_{\mathbf{w}} L^t(\mathbf{w}^{t-1})$ 
15:   UpdateThreshold( $T$ ) ▷ Smoothly increase threshold
16: end while

```

TABLE I: Summary of hardware component cost on an  $m \times n$  layer in SVD-based ONN and our proposed architecture (size- $k$  circulant blocks). Most area-consuming components are considered. PS and DC represent phase shifter and directional coupler.

	#DC	#PS
SVD-ONN	$m(m-1) + n(n-1) + \max(m, n)$	$\frac{m(m-1)+n(n-1)}{2}$
$T\Sigma U$ -ONN	$m(m-1) + 2n + \max(m, n)$	$\frac{m(m-1)+2n}{2}$
Our ONN	$\frac{mn(\log_2 k+1)}{k}$	$\frac{mn(2\log_2 k+1)}{k}$

Given an  $n$ -input,  $m$ -output layer, the SVD-based implementation requires  $m(m-1)/2 + n(n-1)/2$  MZIs and  $\max(m, n)$  attenuators to realize the weight matrix. Therefore, with the aforementioned assumption, the total number of components it costs is given by,

$$\begin{aligned} \#DC_{SVD} &= m(m-1) + n(n-1) + \max(m, n) \\ \#PS_{SVD} &= m(m-1)/2 + n(n-1)/2. \end{aligned} \quad (10)$$

For the slimmed  $T\Sigma U$ -based ONN architecture [20], one unitary matrix is replaced by a compact sparse tree network consisting of  $n$  MZIs. Therefore, the component utilization of  $T\Sigma U$ -based ONN is given by,

$$\begin{aligned} \#DC_{T\Sigma U} &= m(m-1) + 2n + \max(m, n) \\ \#PS_{T\Sigma U} &= m(m-1)/2 + n. \end{aligned} \quad (11)$$

For our architecture, each  $k \times k$  circulant matrix costs  $k$  attenuators and corresponding components required by  $k$ -point OFFT/OIFFT. The following formulation gives the number of components for a  $k$ -point OFFT/OIFFT.

$$\begin{aligned} \#DC_{OFFT}(k) &= 2 \times \#DC_{OFFT}(k/2) + k/2 = \frac{k}{2} \log_2 k \\ \#PS_{OFFT}(k) &= k(\log_2 k + 1) \end{aligned} \quad (12)$$

A phase shift is physically meaningful only when it is within  $(-2\pi, 0]$  as phases can wrap around. Hence, multiple successive phase shifters on the same segment of a waveguide can be merged as one phase shifter, which can be seen when comparing Fig. 2 and Fig. 3. Then the total number of components used in our design to implement an  $m \times n$  weight matrix with size- $k$  circulant sub-matrices is given by,

$$\begin{aligned} \#DC_{Ours}(k) &= \frac{m}{k} \times \frac{n}{k} \times (2 \times \#DC_{OFFT}(k) + k) \\ &= \frac{mn}{k} (\log_2 k + 1) \\ \#PS_{Ours}(k) &= \frac{m}{k} \times \frac{n}{k} \times (2 \times \#PS_{OFFT}(k) - k) \\ &= \frac{mn}{k} (2 \log_2 k + 1). \end{aligned} \quad (13)$$

In practical cases,  $k$  will be set to small values, such as 2, 4, and 8. Given arbitrary values of  $m$  and  $n$ , the proposed architecture costs theoretically fewer optical components than the SVD-based architecture.

We also give a qualitative comparison with incoherent microring resonator-based ONNs (MRR-ONN). There are two MRR-ONN variants. The first one is based on all-pass microring (MR) resonators [29]. The second one proposed later is based on the differential add-drop MR resonators [30]. We assume an  $M \times N$  matrix multiplication in the following tasks. Since the physical dimensions of MRs are smaller than couplers and phase shifters in general, thus a lower area cost can be expected for MRR-ONNs compared with ours. However, in terms of model expressivity, all-pass MRR-ONN is much less than the other two, since it only supports positive weights. Add-drop MRR-ONN and our architecture can support a full-weight range without positive limitation. In terms of robustness, MRR-ONNs are less robust since the MR resonators are more sensitive to device variations and environmental changes than phase shifters. Especially for add-drop MRR-ONN, its differential structure amplifies the noise on the MR transmission factor by 2 times on its represented weight. Thus less robustness can be expected for MRR-ONNs. Furthermore, in terms of power consumption, our architecture can benefit from structured sparsity to obtain a much lower power, which will be shown in later Experimental Results sections. In contrast, for MRR-ONNs, even though a group of weights get pruned to zero values, the corresponding MR resonators are not idle [29], [30], which means its power consumption can barely benefit from pruning techniques. Therefore, from the above qualitative analysis, though our architecture demonstrates a relatively larger footprint than MRR-ONNs, we outperform them in terms of model expressivity, robustness, and power.

## VI. EXTENSION TO OPTICAL CNN WITH LEARNABLE TRANSFORMATIONS

To demonstrate the applicability of the proposed architecture, we extend this architecture to a compact frequency-domain microdisk (MD)-based optical convolutional neural network (CNN) with joint learnability, where the convolutional kernels and frequency-domain transforms are jointly optimized during hardware-aware training.

### A. Microdisk-based Frequency-domain CNN Architecture

Given the two-dimensional (2-D) nature of photonic integrated chips (PICs), currently we only demonstrate optical designs for MLPs. Previous solutions to accelerate convolutional neural networks (CNNs) are based on kernel sliding, convolution unrolling, and time multiplexing [31], [32]. At each time step, the input feature chunks and corresponding convolutional kernels are flattened as a one-dimensional vector and fed into the ONNs to perform vector dot-product. Another solution to solve this is to use *im2col* algorithm [29], [33], that transforms convolution to general matrix multiplication (GEMM). Convolutional kernels and input features are re-shaped as matrix-matrix multiplication, which can be directly mapped on ONNs. Such implementation is inherently inefficient as overlapped convolutional patterns will create a huge amount of data redundancy in the unrolled feature maps. In this work, we proposed to achieve CNNs with a new ONN architecture equipped with learnable transformation structures. Figure 6 demonstrates our proposed optical MD-based CNN architecture featured by kernel sharing, learnable transformation, and augmented frequency-domain kernel techniques. Multi-channel input feature maps are encoded onto multiple wavelengths and input into the learnable frequency-domain transforms, then split into multiple branches through the fanout network for parallel multi-kernel

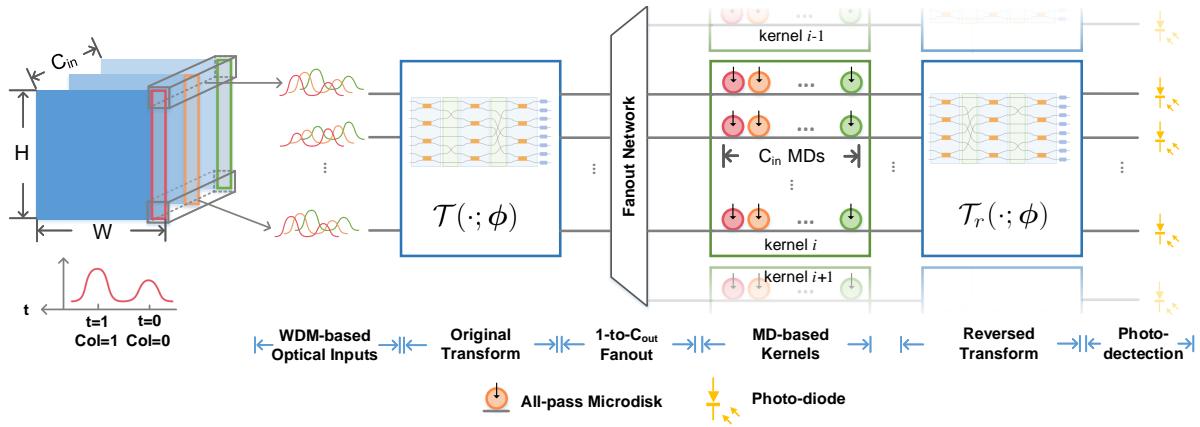


Fig. 6: Architecture of an MD-based optical convolutional layer with trainable frequency-domain transforms. Columns of input features are fed into the architecture in different time steps. Multiple kernels are implemented with multiple photonic chiplets to achieve higher parallelism.

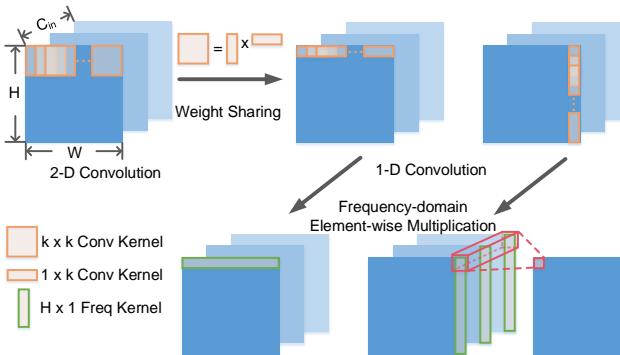


Fig. 7: 2-D convolutional kernel decomposition using weight sharing and frequency-domain transformation.

processing. Frequency-domain convolution is performed in the MD-based kernel banks and the final results are transformed back to the real domain via the reversed transforms. Note that we do not include a detailed discussion on the pooling operations since they are not the computationally-intensive parts in NNs. For example, optical comparators can be used to achieve max-pooling. Average-pooling can be implemented by a fixed-weight convolution engine based on combiner-tree networks. Multiple layers can be cascaded through O-E-O conversion. The phase information loss during photo-detection can be fully modeled during training without harming the model expressivity, which is actually a competitive substitute for ReLU activation in the complex NN domain [34]. All of our experiments in later sections model this phase removal during training, which shows that this non-ideality induced by photo-detection does not cause any accuracy loss. We will introduce details of the principles of the designed optical CNN in the following section.

### B. Kernel Weight Sharing

Modern CNN architectures, e.g., inception architecture [35], adopts weight sharing to reduce the number of parameters in the convolutional layers. For example, a  $5 \times 5$  2-D convolution involves 25 parameters. It can be replaced by two cascaded lightweight  $1 \times 5$  and  $5 \times 1$  convolutions, which only contain 10 unique variables. Such a strategy trains a low-rank convolutional kernel and can benefit its photonic implements as it can be directly applicable to 2-D PICs, which is visualized in Fig. 7.

### C. Learnable Frequency-domain Convolution

Spatial domain convolution requires to slide the receptive field of convolutional kernels across the input features. This could induce hardware implementation difficulty and inefficiency as time multiplexing increases the latency and control complexity of photonic convolution. we solve this issue by a parametrized frequency-domain convolution method. As mentioned before, we decompose the 2-D convolution as row-wise and column-wise 1-D convolutions through weight sharing. For brevity, we focus on the column-wise frequency-domain convolution in the following discussion. The same principle also applies to the row-wise convolution. The column-wise convolution can be formulated as,

$$\mathbf{w} * \mathbf{x} = \mathcal{T}^{-1}(\mathcal{T}(\mathbf{w}; \phi) \odot \mathcal{T}(\mathbf{x}; \phi); \phi), \quad (14)$$

where  $\mathcal{T}(\cdot; \phi)$  is the learnable frequency-domain projection, and  $\phi$  represents the trainable parameters in it. This parametrized transformation enlarges the parameter space to compensate for the model expressiveness degradation induced by kernel weight sharing. Considering the learnable transform as a high-dimensional unitary rotation, it is not necessary to adopt an inverse transform pair to limit the exploration space. To enable the maximum learnability of our trainable transform structure, we relax the inverse transform to a reversed transform,

$$\mathbf{w} * \mathbf{x} = \mathcal{T}_r(\mathcal{T}(\mathbf{w}; \phi) \odot \mathcal{T}(\mathbf{x}; \phi); \phi_r), \quad (15)$$

where  $\mathcal{T}_r$  has a reversed butterfly structure but is not constrained to be the inverse of  $\mathcal{T}$ .

We now discuss how our proposed trainable transform structures can move beyond Fourier transform, thus enable hardware-aware learnability. Fourier transform is a complex domain transformation that is mathematically designed for frequency component extraction. However, the Fourier transform is not necessary to be the best-performed transformation that can be used in CNNs. Other manually designed unitary transforms are also experimentally demonstrated to have a similar ability for signal integration and extraction [36]. Hence, we upgrade the fixed transformation structure to an adaptive structure where all phase shifters are trainable. As mentioned in the Section V, phase shifters in the same segment of waveguide can be merged into one phase shifter. Therefore, to avoid redundant trainable phase shifters, we re-design the learnable basic block, as shown in Fig. 8. For the original transformation, two phase shifters  $\phi_1$  and  $\phi_2$  are placed

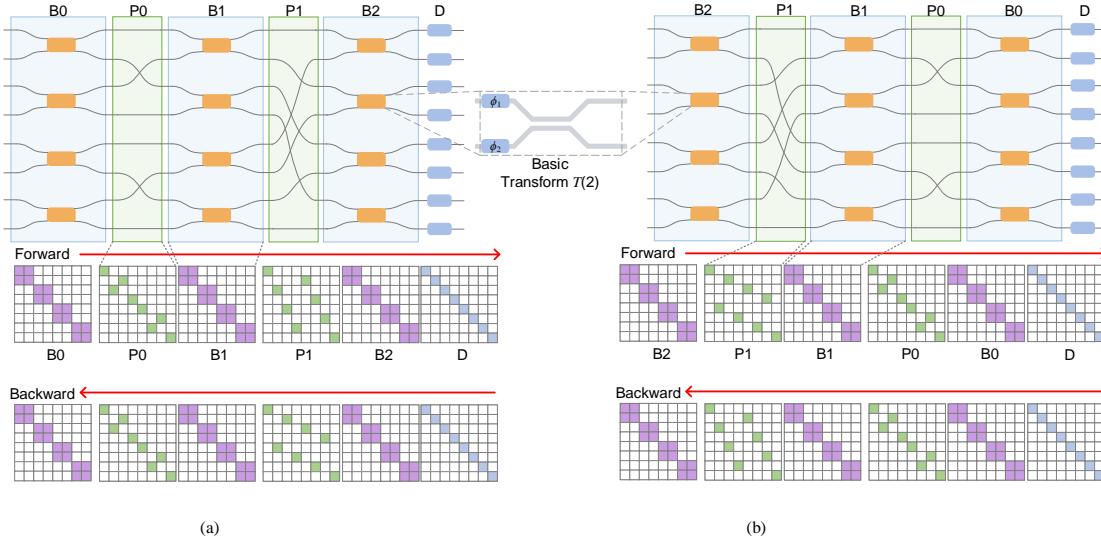


Fig. 8: (a) The original learnable frequency-domain transformation structure. (b) The reversed learnable transformation structure.

on the input port of the directional coupler. The transfer function of a learned basic block can be formulated as,

$$\begin{aligned} \mathcal{T}(2) &= \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & j \\ j & 1 \end{pmatrix} \begin{pmatrix} e^{j\phi_1} & 0 \\ 0 & e^{j\phi_2} \end{pmatrix} \\ &= \frac{1}{\sqrt{2}} \begin{pmatrix} \cos \phi_1 + j \sin \phi_1 & -\sin \phi_1 + j \cos \phi_1 \\ -\sin \phi_2 + j \cos \phi_2 & \cos \phi_2 + j \sin \phi_2 \end{pmatrix}. \end{aligned} \quad (16)$$

In the reversed transformation structure, the basic block is the same as used in the original transforms since the inverse basic block requires a conjugate transposed transfer function which is not implementable with this basic block. Based on this basic block, we recursively build a trainable  $N$ -length transform with a butterfly structure, which can be described as  $\log_2 N$  stages of projection,  $\log_2 N - 1$  stages of permutation, and a final extra group of phase shifters. The original transformation, shown in Fig. 8(a), can be formulated as,

$$\mathcal{T}(N) = \mathcal{D} \mathcal{B}_{\log_2 N-1}(N) \prod_{i=0}^{\log_2 N-2} \mathcal{P}_i(N) \mathcal{B}_i(N), \quad (17)$$

where  $\mathcal{B}_i(N)$  the  $i$ -th stage of butterfly projection,  $\mathcal{P}_i(N)$  is the  $i$ -th stage signal permutation, and the diagonal matrix  $\mathcal{D}$  represents the final extra column of phase shifters. The butterfly projection operator  $\mathcal{B}(N)$  is a diagonal matrix with a series of  $\mathcal{T}(2)$  as its diagonal sub-matrices,

$$\mathcal{B}(N) = \begin{pmatrix} \mathcal{T}_0(2) & 0 & \dots & 0 \\ 0 & \mathcal{T}_1(2) & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & \mathcal{T}_{N/2-1}(2) \end{pmatrix} \quad (18)$$

The index permutation operator  $\mathcal{P}_i(N)$  can be expressed as a size- $N$  identity matrix with reordered rows. As shown in  $\mathcal{P}_0$  and  $\mathcal{P}_1$  in Fig. 8, the green entries represent 1, and other blank entries represent 0. Note that the permutation operators in the reversed structure is simply the reversed counterparts in the original structure, i.e.,  $\mathcal{P}_{i,ori}(N) = \mathcal{P}_{i,rev}^T(N)$ . The reversed learnable transformation, shown in Fig. 8(b), is designed to have reversed butterfly structure which can be derived as follows,

$$\mathcal{T}_r(N) = \mathcal{D} \left( \prod_{i=0}^{\log_2 N-2} \mathcal{B}_{r,i}(N) \mathcal{P}_{r,i}(N) \right) \mathcal{B}_{r,\log_2 N-1}(N). \quad (19)$$

Note that the reversed transform is not guaranteed to be inverse to the original transform, which requires particular phase configurations discussed later.

Compared with its MZI-based counterparts, this trainable butterfly transformation structure has a constrained projection capability as only a limited set of unitary matrices can be implemented by it [37], [38]. As shown in unitary group parametrization, a full  $N$ -dimensional unitary space  $\mathcal{U}(N)$  has  $N(N - 1)/2$  independent parameters, while the butterfly structure substitutes part of parametrized unitary matrices with fixed permutation operators. Hence, based on full two-dimensional unitary matrices  $\mathcal{U}(2)$ , the butterfly structure has  $2N \log_2 N$  independent parameters. Our proposed learnable block  $\mathcal{T}(2)$  is a reduced version of  $\mathcal{U}(2)$ , as it only covers half of the full 2-D planar rotation space. The pruned transform space  $\mathcal{T}^*(2)$  can be expressed as the conjugate transpose of  $\mathcal{T}(2)$ , which is not implementable without waveguide crossings.

$$\mathcal{T}^*(2) = \frac{1}{\sqrt{2}} \begin{pmatrix} 0 & -j \\ -j & 0 \end{pmatrix} \begin{pmatrix} 1 & j \\ j & 1 \end{pmatrix} \begin{pmatrix} e^{j\phi_1} & 0 \\ 0 & e^{j\phi_2} \end{pmatrix} \quad (20)$$

Equivalently, our learnable transformation structure has  $N \log_2 N$  free parameters.

#### D. Microdisk-based Augmented Kernels

To enable highly-parallel CNN architecture with reinforced model expressiveness, we propose MD-based augmented convolutional kernels with multi-level parallelism across input features, input channels and output channels.

In our design, each 2-D convolutional layer consists of two cascaded 1-D frequency-domain convolutions along columns and rows. We will focus on the column-wise convolution, and the same architecture applies to its row-wise counterpart with an extra matrix transposition operation. We denote the input feature map as  $\mathbf{I} \in \mathbb{R}^{C_{in} \times H \times W}$ , which  $C_{in}$ ,  $H$ ,  $W$  represent the number of input channel, spatial height, and spatial width, respectively. At time step  $t$ , The corresponding column  $\mathbf{I}_{:,t,:} \in \mathbb{R}^{C_{in} \times H \times 1}$  will be input into the optical CNN. Different input channels are encoded by different wavelengths  $\{\lambda_0, \lambda_1, \dots, \lambda_{C_{in}-1}\}$ . Through the wide-band learnable transformation structure, we obtain the frequency-domain features  $\mathcal{T}(\mathbf{I}_{:,t,:}; \phi)$ . This stage enables parallel transformation across the input channels. Then the optical signals carrying those features will be split into  $C_{out}$  planes for data reuse. Such a multi-dimensional ONN design can be supported by state-of-the-art

integration technology with multiple photonic chiplets [39]. In the MD-based convolution stage,  $C_{out} \times C_{in} \times H$  all-pass MDs are used to implement the frequency-domain kernels  $\mathbf{W} \in \mathbb{R}^{C_{out} \times C_{in} \times H}$ . Given that the working principle of MD is primarily optical signal magnitude modulation, our augmented kernels are trainable only in the magnitude space without phase modulation. Each convolutional core is designed to perform the convolution of one output channel. This MD-based convolution is different from the previous EM stage consisting of attenuators and phase shifters. First, all pass MDs can only perform configurable magnitude modulation of the input signals with fixed phase responses, which means the augmented kernels will not expand over the entire complex space. Here we give the transfer function of an MD,

$$\begin{aligned} I_{out} &= W \cdot I_{in} \\ \cos \theta &= \frac{a^2 + r^2 - W(1 + r^2 a^2)}{2(1 - W)ar} \\ \phi_{out} &= \pi + \theta + \arctan \frac{r \sin \theta - 2r^2 a \sin \theta \cos \theta + ra^2 \sin \theta}{(a - r \sin \theta)(1 - ra \cos \theta)}, \end{aligned} \quad (21)$$

where  $I_{in}$  is the magnitude of the input light,  $I_{out}, \phi_{out}$  are magnitude and phase of the output optical signal,  $\theta, a, r$  are the phase, self-coupling coefficient, and coupling loss factor of an MD, respectively.  $W$  is the transmittivity of the MD which corresponds to the trained augmented kernel weight. Typically, parameter  $a$  and  $r$  are very close to 1. Our proposed architecture enables another level of parallelism across output channels. Given that different convolutional kernels share the same input features, multiple MD convolution cores and reversed transform structures will share one original transform structure for hardware reuse and highly-parallel convolution.

A higher modeling capacity is enabled by our augmented kernel technique. Instead of training spatial kernels  $\mathbf{w}$ , we explicitly train the latent weights  $\mathbf{W}$  in the frequency domain without performing  $\mathcal{T}(\mathbf{w}; \phi)$  during training. The augmented latent weights  $\mathbf{W}$  will not meet the conjugate symmetry constraint as its spatial-domain counterparts are not real-valued. Hence, this enables a potentially infinite solution space in the spatial kernel space with various kernel sizes and shapes.

We briefly discuss the scalability of this WDM-based highly-parallel architecture. WDM plays an important role in the high parallelism of our proposed frequency-domain optical CNN. Currently, the widely acknowledged maximum number of wavelength in the single-mode dense-WDM (DWDM) is over 200 [40]–[42]. This means in our architecture, the number of input channels  $C_{in}$  that can be processed in parallel is over 200 if a single mode is adopted, which can support most modern CNN architectures. Since different modes of optical signals can also propagate through a multi-mode waveguide independently, the capacity of DWDM can be further extended in another dimension. If  $M$  different modes of optical signals are adopted, the number of parallel input channels can be extended by another  $M$  times, where  $M$  can be up to 10 given the current technology. Therefore, the potential input-channel-wise parallelism that we can provide is enough for most modern CNN applications.

#### E. Discussion: Exploring Inverse Transform Pairs in Constrained Unitary Space

In manually designed frequency-domain convolution algorithms, domain transformation will be designed to be inverse, e.g., FFT and IFFT. This implies an inverse constraint between two mutually-reversed transform structures  $\mathcal{T}$  and  $\mathcal{T}_r$ . To be able to realize trainable inverse transform pairs, we add unitary constraints to our learnable transform structures,

$$\mathcal{T}_r(\cdot, \phi_r) = \mathcal{T}^{-1}(\cdot; \phi). \quad (22)$$

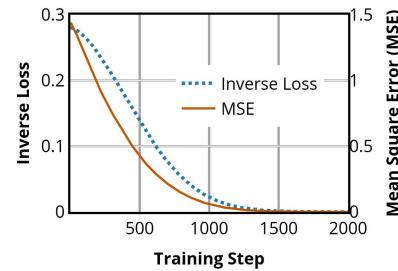


Fig. 9: Training curve of inverse loss  $\mathcal{L}_{inv}$  and mean square error between trained phase configurations and theoretical 4-point OFFT settings.

Inverse constraints typically can be addressed via adding a regularization term in training,

$$\mathcal{L}_{inv} = \|\mathcal{T}_r(\mathcal{T}(\mathbf{e})) - \mathbf{e}\|_2, \quad \mathbf{e} \in \mathbb{C}^N. \quad (23)$$

However, this requires explicit transfer matrices of  $\mathcal{T}$  and  $\mathcal{T}_r$  to compute this regularization term [43], which is memory-intensive and computational expensive as indicated by Eq. (18), Eq. (19). We propose an efficient regularization method to exert inverse constraint.

$$\mathcal{L}_{inv} = \|\mathcal{T}_r(\mathcal{T}(\mathbf{e})) - \mathbf{e}\|_2, \quad \mathbf{e} \in \mathbb{C}^N, \quad (24)$$

where  $\mathbf{e}$  is the orthonormal bases of  $N$ -dimensional complex space. Notice that if  $\mathcal{T}_r(\mathcal{T}(\mathbf{e})) = \mathbf{e}$ , then for any  $\mathbf{x} = \boldsymbol{\alpha}^T \mathbf{e}$  the following statement holds,

$$\mathcal{T}_r(\mathcal{T}(\mathbf{x})) = \mathcal{T}_r(\mathcal{T}(\boldsymbol{\alpha}^T \mathbf{e})) = \boldsymbol{\alpha}^T \mathcal{T}_r(\mathcal{T}(\mathbf{e})) = \mathbf{x}. \quad (25)$$

Thus transforms  $\mathcal{T}$  and  $\mathcal{T}_r$  are inverse transforms once the regularization loss reaches 0. This surrogate method reduce the computation complexity from  $\mathcal{O}(N^2 \log_2 N)$  in Eq. (17) to  $\mathcal{O}(N \log_2 N)$ , where diagonal matrix multiplication with  $\mathcal{B}(N)$  is simplified by  $2 \times 2$  submatrix multiplication with  $\mathcal{T}(2)$ .

Using our proposed inverse pair regularization method, we show that our trainable transform  $\mathcal{T}$  can efficiently learn Fourier transform by setting  $\mathcal{T}_r$  as OIFFT. Figure 9 demonstrates that the trainable transform will quickly converge to the theoretical OFFT as the mean square error between trained phase settings and target phase shifter settings reduces to 0 when the loss converges.

#### F. Discussion: Hardware-aware Pruning for Trainable Transforms

In this section, we demonstrate that our proposed trainable transform has excellent compatibility with hardware-aware pruning techniques. Compared to the fixed manual design of frequency-domain transforms, e.g., OFFT, our pruned trainable transform can potentially improve power consumption and noise-robustness by eliminating a proportion of configurable devices in a hardware-aware way. In Fig. 8, a column of phase shifters are located in an aligned placement style to guarantee light path coherency. Those thermo-optically configurable phase shifters contribute to a large proportion of on-chip energy consumption and nearly half of the chip area [44], [45]. In the FFT transform, a majority of phase shifters are non-zero degrees; thus the power consumption is under-optimized. Moreover, Even though there could be several arbitrary zero degree phase shifters in FFT transforms, they are arbitrarily located in the structure, the removal of which barely saves any chip area considering the aligned device placement. Therefore, pruning an entire column of phase shifters will physically improve both area and energy cost. We adopt a phase-wrapping Group Lasso regularization similar to Eq. (2) together with incremental pruning technique to slim the trainable transforms targeted

at lower area cost and lower power consumption. The proposed phase-wrapping Group Lasso (PhaseGL) is formulated as,

$$L_{PhaseGL} = \sum_{g=0}^G \sqrt{1/p_g} \|\phi_g - \phi_g^*\|_2, \quad (26)$$

$$\phi_{g,i}^* = \begin{cases} 0, & \phi_{g,i} \in [0, \pi), 0 \leq i < p_g \\ 2\pi, & \phi_{g,i} \in [\pi, 2\pi), 0 \leq i < p_g, \end{cases}$$

where  $\phi_g$  is a column of phase shifters and this regularization term encourages phases towards their corresponding prunable targets  $\phi_g^*$ .  $G$  is the total columns of phase shifters, which is  $(\log_2 N + 1)$  for a length- $N$  transform. Once the group lasso of a column falls below a threshold  $T_T$ , the entire column of phase shifters are pruned. The ratio of pruned columns to all phase shifter columns is called transform sparsity ( $T$  sparsity), defined as,

$$s_T = \frac{|\{\phi_g | \sqrt{1/p_g} \|\phi_g - \phi_g^*\| < T_T\}|}{G}.$$

Our proposed regularization and pruning strategy improves area cost as an entire column of phase shifters are pruned to save chip area in the actual layout. Furthermore, power consumption can also be improved as the total power consumption for trainable transform structures can be estimated as  $P_T \propto \sum \phi$ . This phase-wrapping pruning also improves the tolerance to noise for two aspects, considering there are manufacturing variations in the phase shifter [45], [46]. First, fewer noise sources eventually introduce smaller error when device-level noises are considered in the physical implementation [20], [46]. Second, given the quadratic relationship between phase shift and voltage controls, expressed as  $\phi = \gamma v^2$ , smaller phases will be less sensitive to  $\gamma$  noise [45].

#### G. Discussion: Hardware Cost of the Proposed MD-based Optical CNN

We give a summary on the hardware component usage of the proposed MD-based optical CNN architecture in Table. II. Our architecture shares the original transform among multiple kernels to save area. Our proposed pruning technique can regularly sparsify the transform structures for further area reduction. The MD-based convolution stage is very compact since the footprint of an MD is two-order-of-magnitude smaller than a DC. In contrast, the SVD-based ONN costs  $H(C_{out}^2 + C_{in}^2 \times K^4)$  DCs and  $H(C_{out}^2/2 + C_{in}^2 \times K^4/2)$  PSs to achieve the same latency with our architecture, i.e.,  $H$  forwards to finish a convolutional layer, where  $K$  is the spatial kernel size. For example, if we set  $H=64$ ,  $C_{in}=C_{out}=32$ ,  $K=3$ ,  $s_T=0.5$ , our architecture uses  $>370\times$  fewer DCs and  $>180\times$  fewer PSs than the single-wavelength SVD-based ONN. If SVD-based ONNs also use WDM techniques for higher parallelism with the same number of wavelength as ours, i.e., 32, we still outperform theirs by  $11.6\times$  fewer DCs and  $5.6\times$  fewer PSs. Hence, our frequency-domain CNN architecture outperforms previous MZI-ONNs with higher computational efficiency and better scalability by a large margin.

## VII. EXPERIMENTAL RESULTS

We conduct numerical simulations for functionality validation and evaluate our proposed architecture on the hand-written digit recognition dataset (MNIST) [28] with various network configurations. Quantitative evaluation shows that our proposed architecture outperforms the SVD-based and TΣU-based ONN architectures in terms of area cost without any accuracy degradation. We further evaluate our proposed MD-based optical CNN architecture and demonstrates its superior power reduction and robustness improvement on MNIST and FashionMNIST [50] dataset.

TABLE II: Hardware cost summary on the proposed MD-based optical CNN architecture. The input feature map is of size  $H \times W \times C_{in}$ , the number of output channels is  $C_{out}$ , and the sparsity of the learnable transforms is  $s_T \in [0, 1]$ . For simplicity, we assume  $H = W$ , which is a widely used configuration for most CNNs. Given the ultra-compact footprint of an MD, e.g.,  $5 \times 5 \mu m^2$  [47], we count 100 MDs as one DC in the area estimation. The row-wise and column-wise convolutions are both counted in this table.

Structure	Hardware Cost
$T$	$H \log_2 H$ DCs + $2s_T H(1 + \log_2 H)$ PSs
Kernel	$2HC_{in}C_{out}$ MDs $\approx \frac{H}{50}C_{in}C_{out}$ DCs
$T_r$	$H \log_2 HC_{out}$ DCs + $2s_T H(1 + \log_2 H)C_{out}$ PSs
Total	$\approx H(\log_2 H + \frac{C_{in}}{50})C_{out}$ DCs + $2s_T H(1 + \log_2 H)C_{out}$ PSs

TABLE III: Optical component sizes used in the area estimation.

Optical Component	Length ( $\mu m$ )	Width ( $\mu m$ )
3-dB Directional Coupler [3]	54.4	40.3
Thermo-optic Phase Shifter [44]	60.16	0.50
2-to-1 Optical Combiner [48]	20.00	3.65
Waveguide Crossing [49]	5.9	5.9

#### A. Simulation Validation

To validate the functionality of our proposed architecture, we conduct optical simulations on a  $4 \times 4$  circulant matrix-vector multiplication module using Lumerical INTERCONNECT tools. First, we encode a  $4 \times 4$  identity weight matrix into our architecture and input 4 parallel optical signals to validate its functionality. For brevity, we plot several different representative cases in Fig. 10a. It shows that our designed architecture can correctly realize identity projection. Further, we randomly generate a length-4 real-valued weight vector  $\mathbf{w} = (0.2, -0.1, 0.24, -0.15)$  to represent a circulant matrix, and encode  $\mathcal{F}(\mathbf{w}) = (0.19e^{0j}, 0.064e^{-2.246j}, 0.69e^{0j}, 0.064e^{2.246j})$  into attenuators and phase shifters in the EM stage. The simulation results in Fig. 10b show good fidelity ( $< 1.2\%$  maximum relative error) to the ground truth results.

#### B. Comparison Experiments on FFT-based ONNs

To evaluate our proposed ONN architecture, we conduct a comparison experiment on a machine learning dataset MNIST [28], and compare the hardware utilization, model expressivity among four architectures: 1) SVD-based architecture [3]; 2) TΣU-based architecture [20]; 3) Ours without pruning; 4) Ours with pruning.

For the SVD-based ONN, we simply train an original MLP since this architecture directly implements matrix multiplication in the fully-connected layer. In the TΣU-based architecture, training is performed on a sparse matrix  $T$  designed for dimensionality matching, a diagonal matrix  $\Sigma$ , and a pseudo-unitary matrix  $U$  with unitary regularization. This architecture eliminates one of the area-cost unitary matrices and adopts a sparse-tree  $T$  to match dimensionality. The orthogonality constraint of  $U$  is first relaxed with unitary regularization in the training and satisfied by post-training unitary projection [20].

We implement the proposed architecture with different configurations in PyTorch and test the inference accuracy on a machine with an Intel Core i9-7900X CPU and an NVIDIA TitanXp GPU. We set  $\lambda$  to 0.3 for the Group Lasso regularization term, initialize all trainable weights with a Kaiming-Normal initializer [51], adopt the Adam optimizer [52] with initial learning rate= $1 \times 10^{-3}$  and a step-wise exponential-decay learning rate schedule with decay rate=0.9. We use the ideal rectified linear units (ReLU) activation function as nonlinearity. All NN models are trained for 40 epochs with a mini-batch size of 32 till fully converged. Figure 11 plots the test accuracy

TABLE IV: Comparison of inference accuracy and hardware utilization on MNIST dataset with different configurations. For example, configuration (28×28)-1024(8)-10(2) indicates a 2-layer neural network, where the first layer has 784 input channels, 1024 output channels with size-8 circulant matrices, and so on.

Network Configurations		Block Sparsity	#Parameter	Accuracy	#DC	#PS	Area ( $cm^2$ )
Model 1	SVD [3]: (28×28)-400-10	0.00	318 K	98.49%	934 K	467 K	20.62
	TΣU [20]: (28×28)-400-10	0.00	318 K	98.49%	777 K	388 K	17.15
	Ours w/o Prune: (28×28)-1024(8)-10(2)	0.00	105 K	98.32%	412 K	718 K	9.33
	Ours w/ Prune: (28×28)-1024(8)-10(2)	0.40	63 K	98.26%	244 K	425 K	5.53
Model 2	SVD [3]: (14×14)-70-10	0.00	14 K	96.93%	48 K	24 K	1.07
	TΣU [20]: (14×14)-70-10	0.00	14 K	96.93%	44 K	22 K	0.97
	Ours w/o Prune: (14×14)-256(4)-10(2)	0.00	14 K	96.93%	40 K	67 K	0.90
	Ours w/ Prune: (14×14)-256(4)-10(2)	0.45	8 K	96.91%	22 K	36 K	0.49
Model 3	SVD [3]: (28×28)-400-128-10	0.00	366 K	98.58%	967 K	483 K	21.35
	TΣU [20]: (28×28)-400-128-10	0.00	366 K	98.58%	794 K	396 K	17.52
	Ours w/o Prune: (28×28)-1024(8)-128(4)-10(2)	0.00	134 K	98.53%	501 K	868 K	11.34
	Ours w/ Prune: (28×28)-1024(8)-128(4)-10(2)	0.39	81 K	98.43%	289 K	517 K	6.77
Model 4	SVD [3]: (14×14)-160-160-10	0.00	59 K	97.67%	141 K	70 K	3.10
	TΣU [20]: (14×14)-160-160-10	0.00	59 K	97.67%	91 K	45 K	2.00
	Ours w/o Prune: (14×14)-256(4)-256(8)-10(2)	0.00	22 K	97.67%	73 K	123 K	1.64
	Ours w/ Prune: (14×14)-256(4)-256(8)-10(2)	0.37	14 K	97.52%	47 K	79 K	1.05

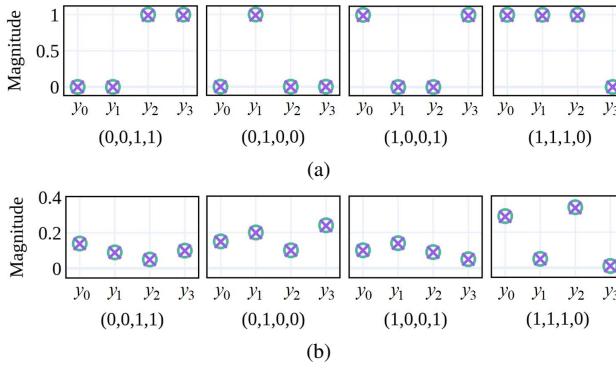


Fig. 10: (a) Simulated output intensities (crosses) and ground truth (circles) of a 4×4 identity circulant matrix-vector multiplication. (b) Simulated output intensities (crosses) and ground truth (circles) of a 4×4 circulant matrix-vector multiplication, with  $w=(0.2, -0.1, 0.24, -0.15)$ . E.g., (0,0,1,1) is the input signal.

and sparsity curves as training proceeds. The first 5 epochs are the initial training phase. After that, the model has roughly converged. In the subsequent structured pruning phase, we apply our incremental pruning strategy. In each iteration, we set any weights in the EM stage to zero once their Group Lasso falls below the threshold. Then, we adopt a step-wise function to smoothly increase the threshold  $T$  after each epoch finishes.

The structured sparsity for our proposed FFT-based MLP is defined as the percentage of pruned parameters in all parameters, i.e.,  $\|\{w\| \|\omega_{ij}\|_2 < T\}\| / \|w\|$ . We call it block sparsity. We can see that every time sparsity increases, the test accuracy decreases accordingly and then fully recovers during the next training epoch. This alternate pruning and re-training mechanism incrementally improves sparsity while minimizing accuracy loss.

For a fair comparison, all architectures are trained with the same hyper-parameters and have similar test accuracy in each experiment configuration. To estimate the component utilization and area cost, we adopt exactly the same type of photonic devices in all architectures, as listed in Table III, and accumulate the area of each optical component for approximation. Placement or routing information is not considered

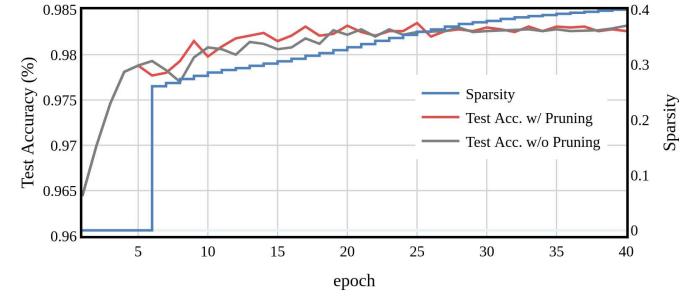


Fig. 11: Training curve of the proposed architecture with setup of (28×28)-1024(8)-10(2).

in our estimation.

In Table IV, the first column indicates different neural network configurations. For example, (14×14)-256(4)-10(2) describes a 2-layer network, with 196 input channels, 256 output channels in the first layer ( $k=4$ ), and 10 output channels in the second layer ( $k=2$ ). The TΣU-based architecture adopts a unique training methodology and claims to have small accuracy degradation (< 1%) [20], thus we assume it has approximately the same accuracy as the SVD-based architecture. In the TΣU-based architecture, the total number of MZIs used to implement an  $m \times n$  weight matrix is bounded by  $n(n + 1)/2$ .

Among various network configurations, our proposed architecture outperforms the SVD-based architecture and the TΣU-based architecture with lower optical component utilization and better area cost. We normalize all areas to our architecture with pruning applied and show the normalized area comparison in Fig. 12. Consistent with analytical formulations in Section V, the experimental results show that, as the difference between input and output channels for each layer in the original MLPs gets larger, our proposed architecture can save a larger proportion of optical components. Specifically, fully-connected layers with highly-imbalanced input and output channels, e.g., ((28×28)-400-10), could benefit the most by using our proposed architecture. For MLPs with nearly-square weight matrices, e.g., ((14×14)-160-160-10), although the area gap is decreasing, our proposed architecture still shows superior area efficiency. This is because FFT-based structured matrix multiplications can reduce many parameter redundancies and save components while still maintaining model expressivity. Fur-

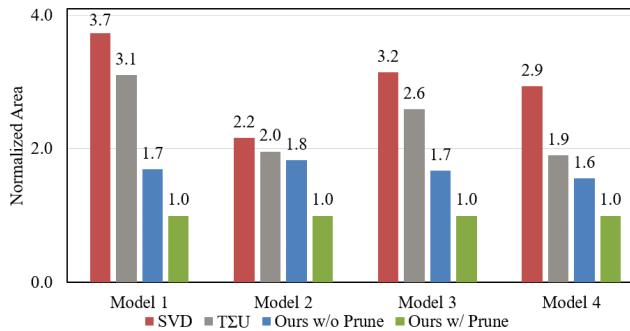


Fig. 12: Normalized area comparison with different model configurations. *Model 1-4* refer to Table IV. *SVD* refers to [3] and *TΣU* refers to [20].

thermore, ablation experiments on our structured pruning method validate the effectiveness of the proposed two-phase training flow. It can save an extra 30-50% optical components with negligible model expressivity loss. Even though the area cost is generally not where optical neuromorphic systems excel, their ultra-low latency and low power consumption make them very promising NN inference accelerators, e.g., in data centers. Therefore, by introducing an area-efficient and pruning-compatible ONN architecture, our work enables more compact ONN implementations without accuracy degradation.

### C. Comparison Among Different Trainable Transform Settings

As mentioned in previous sections, we extend our ONN architecture to MD-based CNNs with trainable frequency-domain transforms. We will demonstrate several experimental evaluations on our proposed MD-based CNN architecture.

First, we discuss how different transform settings impact the CNN performance. Recall that our frequency-domain optical CNN architecture decomposes 2-D convolution into row-wise and column-wise 1-D frequency-domain convolutions. Each 1-D convolution involves original and reversed transforms, thus total four transforms are trainable for each convolutional layer, denoted as  $\mathcal{T}_{row}$ ,  $\mathcal{T}_{row,r}$ ,  $\mathcal{T}_{col}$ ,  $\mathcal{T}_{col,r}$ . Motivated by manual designs of frequency-domain transform, we observe that the row-wise and column-wise 1-D transforms typically share the same transform pair such that the combination of them will be equivalent to their 2-D counterparts. Moreover, the transfer between spatial and frequency domains are designed to be inverse pairs, e.g., FFT and IFFT. Those two successful design priors motivate us to investigate whether transform sharing and inverse pair constraint will benefit frequency CNN performance.

Therefore we evaluate the performance of four different transform settings on MNIST dataset: (1) four transforms are trained independently (*AllFree*); (2) Column-wise and row-wise convolutions share the same transform as  $\mathcal{T}_{row} = \mathcal{T}_{col}$ ,  $\mathcal{T}_{row,r} = \mathcal{T}_{col,r}$  (*Shared*); (3) Reversed transforms are constrained to be close to the inverse transform as  $\mathcal{T}_{row,r} \approx \mathcal{T}_{row}^{-1}$ ,  $\mathcal{T}_{col,r} \approx \mathcal{T}_{col}^{-1}$  (*Inverse*); (4) Transforms are shared between column-wise and row-wise convolutions and the inverse constraints are applied (*InvShared*).

Our CNN configuration is  $16 \times 16 - C16 - BN - ReLU - MaxPool - F32 - ReLU - F10$ , where  $C16$  represents convolution with 16 output channels,  $BN$  represents batch normalization,  $F32$  represents fully connected layers with 32 neurons. Feature maps are reduced to  $5 \times 5$  after Maxpooling, and input images are downsampled to  $16 \times 16$ . Table V shows the comparison results.

Based on the results, we observe that the inverse constraint and shared transform produces no benefits in terms of inference accuracy. Training the original and reversed transforms across row-wise and column-wise convolutions independently offers the best results. Thus, we will use *AllFree* transform settings for our experiments.

TABLE V: Accuracy comparison among four trainable transform settings.

Settings	AllFree	Shared	Inverse	InvShared
Test Accuracy	96.88%	96.13%	96.41%	96.40%

TABLE VI: Transform sparsity ( $\mathcal{T}$  sparsity) and power consumption comparison among optical FFT and our trainable transform with hardware-aware pruning on MNIST and FashionMNIST dataset.  $\mathcal{T}$  sparsity represents how many columns of phase shifters are pruned in our trainable frequency-domain transforms. The power consumption assumes maximum parallelism across output channels, thus 1 original transform and  $C_{out}$  reversed transforms are counted for each layer. For the MNIST dataset, we adopt the ONN configuration as  $16 \times 16 - C16 - BN - ReLU - MaxPool5 - F32 - ReLU - F10$ , and for the FashionMNIST dataset we set the ONN configuration as  $16 \times 16 - C24 - BN - ReLU - MaxPool6 - F64 - ReLU - F10$ . The power consumption is estimated by the sum of phase shifts given that the phase shift is proportional to the thermal tuning power, i.e.,  $\phi \propto v^2$ . Other power consumption sources, e.g., insertion loss, are not considered for simplicity.

Dataset	Transform	OFFT	Trainable (Pruned)
MNIST [28]	$\mathcal{T}$ Sparsity	0%	88.2%
	Normalized Power	100%	18.4%
FashionMNIST [50]	$\mathcal{T}$ Sparsity	0%	88.4%
	Normalized Power	100%	15.5%

### D. Comparison with Hardware-aware Transform Pruning

To jointly optimize classification accuracy and hardware cost in terms of area, power, and robustness, we perform hardware-aware pruning assisted by phase-wrapping Group Lasso regularization to our proposed trainable transforms. The weight for  $L_{PhaseGL}$  is 0.05, and we set 10 epochs for the first pretraining phase and 40 epochs for incremental structured pruning.

1) *Power Consumption Evaluation*: We calculate the energy cost by summing all phase shifts as they are proportional to power consumption, and show the energy saved by our pruned transforms in Table VI. We gradually increase the pruning threshold from both 0 degree and  $2\pi$  degree sides, and pruned 88.2% columns of phase shifters with maximum output channel parallelism. Through this pruned transform, we save 81.6% power consumption compared with optical FFT structure. We also evaluate the power consumption by applying pruned trainable transform in our block-circulant matrix based MLP architecture. The block sparsity, transform sparsity  $\mathcal{T}$  sparsity, power consumption, and area cost are estimated in Table VII. Therefore, our energy-saving and area-efficient ONN architecture is more suitable for resource-constrained applications, e.g., edge computing and online learning tasks [53], [54].

2) *Variation-Robustness Evaluation*: To evaluate the noise-robustness of the frequency-domain transform, we inject device-level variations into phase shifters to introduce phase programming errors and demonstrate the accuracy and its variance under different noise intensities  $\sigma$  on MNIST and FashionMNIST dataset. Specifically, we inject Gaussian noise  $\Delta\gamma \sim \mathcal{N}(0, \sigma^2)$  into the  $\gamma$  coefficient of each phase shifter to perturb its phase response  $\phi_n = (\gamma + \Delta\gamma)v^2$ , where  $\gamma$  is calculated by the voltage that can produce  $\pi$  phase shift as  $\gamma = \pi/v_\pi^2$  and we adopt 4.36V as the typical value of  $v_\pi$  [3], [45]. Figure 13 shows that  $\sim 80\%$  structured sparsity can be achieved by our phase-wrapping pruning method, and our pruned trainable transform outperforms the OFFT structure with over 80% power reduction and much better robustness under various noise intensities.

We also evaluate the robustness on our circulant-matrix-based MLP architecture. Figure 14 compares the phase shifter gamma noise robustness among 1) SVD-based ONNs, 2) optical FFT-based

TABLE VII: Comparison of block sparsity, frequency-domain transform ( $\mathcal{T}$ ) sparsity, normalized power consumption, and estimated area ( $cm^2$ ) among 1) SVD-based ONN, 2) $T\Sigma U$ -based ONN, 3) optical FFT, 4) our trainable transform without pruning transforms, and 5) our trainable transform with hardware-aware pruning on MNIST dataset. SVD-based and  $T\Sigma U$ -based ONN configuration is  $28 \times 28 - 400 - 10$ , and ours is  $28 \times 28 - 1024(8) - 10(2)$ . All ONNs have a similar inference accuracy with a 0.5% accuracy discrepancy among all architectures. Block sparsity is for pruned circulant blocks.  $\mathcal{T}$  sparsity is for pruned trainable frequency-domain transforms. The power consumption is normalized to SVD-based ONN, which is estimated by the sum of all phase shifts given that the phase shift is proportional to the thermal tuning power, i.e.,  $\phi \propto v^2$ .

Architecture	Block Sparsity	$\mathcal{T}$ Sparsity	Power	Area ( $cm^2$ )
SVD-based [3]	-	-	100%	20.62
$T\Sigma U$ -based [20]	-	-	83.1%	17.15
Ours-OFFT [25]	0.40	0.00	98.9%	5.53
Ours-Trainable	0.71	0.00	79.9%	2.54
Ours-Trainable	0.66	0.96	9.9%	2.99

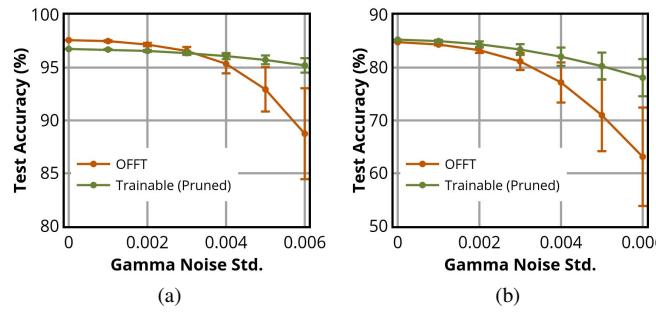


Fig. 13: Robustness comparison among OFFT and pruned trainable transform on MNIST and FashionMNIST dataset. Error bar is drawn to show the  $\pm 1\sigma$  accuracy variance from 20 runs. For MNIST dataset, we adopt the ONN configuration as  $16 \times 16$ -C16-BN-ReLU-MaxPool5-F32-ReLU-F10, and for FashionMNIST dataset we set the ONN configuration as  $16 \times 16$ -C24-BN-ReLU-MaxPool6-F64-ReLU-F10.

architecture, and 3) optical trainable transform-based architecture with transform pruning. The SVD-based architecture shows severe accuracy loss due to phase error amplification effects within its MZI arrays [3], [45], [46]. We do not show accuracy below 90% for clear demonstration. Our FFT-based architecture and trainable transform based architecture both benefit from superior noise robustness due to their structured sparsity and blocking design. The high sparsity removes a large proportion of noise sources, i.e., thermo-optic phase shifters, and the block-circulant structure partitions the weight matrix to avoid correlated error propagation among different weight blocks [45].

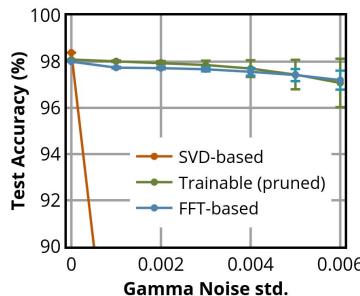


Fig. 14: Robustness comparison on MNIST among SVD-based ONN, optical FFT-based architecture, and trainable transform based architecture with transform pruning. Error bar is drawn to show the  $\pm 1\sigma$  accuracy variance from 20 runs. We adopt the SVD-based ONN configuration as  $28 \times 28 - 400 - 10$  and our architecture as  $28 \times 28 - 1024(8) - 10(2)$ .

## VIII. CONCLUSION

In this work, we propose a hardware-efficient optical neural network architecture. Our proposed ONN architecture leverages block-circulant matrix representation and efficiently realizes matrix-vector multiplication via optical fast Fourier transform. This architecture consists of five stages, including splitter tree, OFFT, element-wise multiplication, OIFFT, and combiner tree. Compared with the previous SVD-based and  $T\Sigma U$ -based ONN architectures, the new design cuts down the optical component utilization and improves area cost by  $2.2 \sim 3.7 \times$  among various network configurations. Our proposed two-phase training flow performs structured pruning to our architecture and further improves hardware efficiency with negligible accuracy degradation. We extend the proposed architecture to an optical microdisk-based frequency-domain CNN, and propose a trainable transform structure to enable a larger design space exploration. We demonstrate structured pruning to our trainable transform structures and it achieves less component usage, over 80% power reduction in CNNs, over 90% power reduction in MLPs, and much better variation-robustness under device-level noises than prior work.

## ACKNOWLEDGMENT

The authors acknowledge the Multidisciplinary University Research Initiative (MURI) program through the Air Force Office of Scientific Research (AFOSR), contract No. FA 9550-17-1-0071, monitored by Dr. Gernot S. Pomrenke.

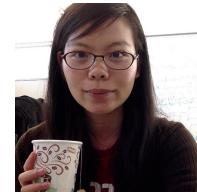
## REFERENCES

- [1] A. Krizhevsky, I. Sutskever, and G. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Proc. NIPS*, 2012.
- [2] T. Mikolov, M. Karafiat, L. Burget *et al.*, “Recurrent neural network based language model,” in *INTERSPEECH*, 2010.
- [3] Y. Shen, N. C. Harris, S. Skirla *et al.*, “Deep learning with coherent nanophotonic circuits,” *Nature Photonics*, 2017.
- [4] Z. Ying, C. Feng, Z. Zhao, S. Dhar *et al.*, “Electronic-photonic arithmetic logic unit for high-speed computing,” *Nature Communications*, 2020.
- [5] C. Feng, Z. Ying, Z. Zhao, J. Gu *et al.*, “Wavelength-division-multiplexing (WDM)-based integrated electronicphotonic switching network (EPSN) for high-speed data processing and transportation,” *Nanophotonics*, 2020.
- [6] C. Feng, Z. Ying, Z. Zhao, J. Gu *et al.*, “Integrated WDM-based Optical Comparator for High-speed Computing,” in *Proc. CLEO*, 2020.
- [7] M. Miscuglio, Z. Hu, S. Li, J. Gu *et al.*, “Million-channel parallelism Fourier-optic convolutional filter and neural network processor,” in *Proc. CLEO*, 2020.
- [8] S. K. Esser, P. A. Merolla, J. V. Arthur *et al.*, “Convolutional networks for fast, energy-efficient neuromorphic computing,” *PNAS*, 2016.
- [9] Y. Wang, W. Wen, B. Liu *et al.*, “Group scissor: Scaling neuromorphic computing design to large neural networks,” in *Proc. DAC*, 2017.
- [10] Y. Zhang, X. Wang, and E. G. Friedman, “Memristor-based circuit design for multilayer neural networks,” *IEEE TCAS I*, 2018.
- [11] A. N. Tait, M. A. Nahmias, B. J. Shastri *et al.*, “Broadcast and weight: An integrated network for scalable photonic spike processing,” *J. Light. Technol.*, 2014.
- [12] J. Bueno, S. Maktoobi, L. Froehly *et al.*, “Reinforcement learning in a large-scale photonic recurrent neural network,” *Optica*, 2018.
- [13] C. Feng, Z. Zhao, Z. Ying, J. Gu, D. Z. Pan, and R. T. Chen, “Compact design of On-chip Elman Optical Recurrent Neural Network,” in *Proc. CLEO*, 2020.
- [14] F. Zokaei, Q. Lou, N. Youngblood *et al.*, “LightBulb: A Photonic-Nonvolatile-Memory-based Accelerator for Binarized Convolutional Neural Networks,” in *Proc. DATE*, 2020.
- [15] M. Miscuglio and V. J. Sorger, “Photonic tensor cores for machine learning,” *Applied Physics Review*, 2020.
- [16] D. Brunner, M. C. Soriano, C. R. Mirasso *et al.*, “Parallel photonic information processing at gigabyte per second data rates using transient states,” *Nature Communications*, 2013.
- [17] L. Vivien, A. Polzer, D. Marris-Morini *et al.*, “Zero-bias 40gbit/s germanium waveguide photodetector on silicon,” *Opt. Express*, 2012.
- [18] M. Reck, A. Zeilinger, H. Bernstein *et al.*, “Experimental realization of any discrete unitary operator,” *Physical review letters*, 1994.

- [19] A. Ribeiro, A. Ruocco, L. Vanacker *et al.*, "Demonstration of a  $4 \times 4$ -port universal linear circuit," *Optica*, 2016.
- [20] Z. Zhao, D. Liu, M. Li *et al.*, "Hardware-software co-design of slimmed optical neural networks," in *Proc. ASPDAC*, 2019.
- [21] Z. Li, S. Wang, C. Ding *et al.*, "Efficient recurrent neural networks using structured matrices in fpgas," in *ICLR Workshop*, 2018.
- [22] S. Han, J. Pool, J. Tran, and W. J. Dally, "Learning both weights and connections for efficient neural networks," in *Proc. NIPS*, 2015.
- [23] N. Simon, J. Friedman, T. Hastie, and R. Tibshirani, "A sparse-group lasso," *Journal of Computational and Graphical Statistics*, 2013.
- [24] O. Grandstrand, *Innovation and Intellectual Property Rights*. Oxford University Press, 2004.
- [25] J. Gu, Z. Zhao, C. Feng *et al.*, "Towards area-efficient optical neural networks: an FFT-based architecture," in *Proc. ASPDAC*, 2020.
- [26] L. Zhao, S. Liao, Y. Wang, Z. Li, J. Tang, and B. Yuan, "Theoretical properties for neural networks with weight matrices of low displacement rank," in *Proc. ICML*, 2017.
- [27] J. Friedman, T. Hastie, and R. Tibshirani, "A note on the group lasso and a sparse group lasso," *arXiv preprint arXiv:1001.0736*, 2010.
- [28] Y. LeCun, "The MNIST database of handwritten digits," <http://yann.lecun.com/exdb/mnist/>, 1998.
- [29] W. Liu, W. Liu, Y. Ye, Q. Lou, Y. Xie, and L. Jiang, "Holylight: A nanophotonic accelerator for deep learning in data centers," in *Proc. DATE*, 2019.
- [30] A. N. Tait, T. F. de Lima, E. Zhou *et al.*, "Neuromorphic photonic networks using silicon photonic weight banks," *Sci. Rep.*, 2017.
- [31] V. Bangari, B. A. Marquez, H. Miller, A. N. Tait, M. A. Nahmias, T. F. de Lima, H. Peng, P. R. Prucnal, and B. J. Shastri, "Digital electronics and analog photonics for convolutional neural networks (DEAP-CNNs)," *IEEE JSTQE*, 2020.
- [32] H. Bagherian, S. A. Skirlo, Y. Shen *et al.*, "On-chip optical convolutional neural networks," *ArXiv*, vol. abs/1808.03303, 2018.
- [33] S. Xu, J. Wang, R. Wang, J. Chen, and W. Zou, "High-accuracy optical convolution unit architecture for convolutional neural networks by cascaded acousto-optical modulator arrays," *Opt. Express*, 2019.
- [34] W. Uijens, "Activating frequencies: Exploring non-linearities in the fourier domain," 2018.
- [35] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proc. CVPR*, 2016.
- [36] R. Zhao, Y. Hu, J. Dotzel, C. D. Sa, and Z. Zhang, "Building efficient deep neural networks with unitary group convolutions," in *Proc. CVPR*, 2019.
- [37] L. Jing, Y. Shen, T. Dubcek, J. Peurifoy, S. Skirlo, Y. LeCun, M. Tegmark, and M. Soljačić, "Tunable efficient unitary neural networks (EUNN) and their application to RNNs," in *Proc. ICML*, 2017.
- [38] M. Y.-S. Fang, S. Manipatruni, C. Wierzynski, A. Khosrowshahi, and M. R. DeWeese, "Design of optical neural networks with component imprecisions," *Opt. Express*, 2019.
- [39] R. Meade, S. Ardalan, M. Davenport, J. Fini, C. Sun, M. Wade, A. Wright-Gladstein, and C. Zhang, "TeraPHY: A high-density electronic-photonic chiplet for optical i/o from a multi-chip module," in *Proc. IEEE OFC*, 2019.
- [40] D. T. H. Tan, A. Grieco, and Y. Fainman, "Towards 100 channel dense wavelength division multiplexing with 100ghz spacing on silicon," *Opt. Express*, 2014.
- [41] C. Feng, Z. Ying, Z. Zhao *et al.*, "Wavelength-division-multiplexing-based electronic-photonic network for high-speed computing," in *Proc. SPIE, Smart Photonic and Optoelectronic Integrated Circuits XXII*, 2020.
- [42] J. Yu and X. Zhou, "Ultra-high-capacity dwdm transmission system for 100g and beyond," *IEEE Communications Magazine*, 2010.
- [43] T. Dao, A. Gu, M. Eichhorn, A. Rudra, and C. Ré, "Learning fast algorithms for linear transforms using butterfly factorizations," in *Proc. ICML*, 2019.
- [44] N. C. Harris *et al.*, "Efficient, compact and low loss thermo-optic phase shifter in silicon," *Opt. Express*, 2014.
- [45] J. Gu, Z. Zhao, C. Feng, H. Zhu, R. T. Chen, and D. Z. Pan, "ROQ: A noise-aware quantization scheme towards robust optical neural networks with low-bit controls," in *Proc. DATE*, 2020.
- [46] Z. Zhao, J. Gu, Z. Ying *et al.*, "Design technology for scalable and robust photonic integrated circuits," in *Proc. ICCAD*, 2019.
- [47] E. Timurdogan, Z. Su, C. V. Poulton *et al.*, "AIM Process Design Kit (AIMPDKv2.0): Silicon Photonics Passive and Active Component Libraries on a 300mm Wafer," in *Optical Fiber Communication Conference*, 2018.
- [48] Z. Sheng, Z. Wang, C. Qiu, L. Li, A. Pang, A. Wu, X. Wang, S. Zou, and F. Gan, "A compact and low-loss mmi coupler fabricated with cmos technology," *IEEE Photonics Journal*, 2012.
- [49] Y. Zhang, A. Hosseini, X. Xu, D. Kwong, and R. T. Chen, "Ultralow-loss silicon waveguide crossing using bloch modes in index-engineered cascaded multimode-interference couplers," *Opt. Lett.*, 2013.
- [50] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms," *Arxiv*, 2017.
- [51] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *ICCV*, 2015.
- [52] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Proc. ICLR*, 2015.
- [53] J. Gu, Z. Zhao, C. Feng, W. Li, R. T. Chen, and D. Z. Pan, "FLOPS: Efficient On-Chip Learning for Optical Neural Networks Through Stochastic Zeroth-Order Optimization," in *Proc. DAC*, 2020.
- [54] T. W. Hughes, M. Minkov, Y. Shi, and S. Fan, "Training of photonic neural networks through in situ backpropagation and gradient measurement," *Optica*, 2018.



**Jiaqi Gu** received the B.E. degree in Microelectronic Science and Engineering from Fudan University, Shanghai, China in 2018. He is currently a post-graduate student studying for his Ph.D. degree in the Department of Electrical and Computer Engineering, The University of Texas at Austin, Austin, TX, USA, under the supervision of Prof. David Z. Pan. His current research interests include machine learning, algorithm and architecture design, optical neuromorphic computing for AI acceleration, and GPU acceleration for VLSI physical design automation. He has received the Best Paper Reward at ASP-DAC'20 and the Best Paper Finalist at DAC'20.



**Zheng Zhao** received the B.S. degree in automation from Tongji University, Shanghai, China, in 2012, and the M.S. degree in electrical and computer engineering from Shanghai Jiao Tong University, Shanghai, China, in 2015. She is currently working toward the Ph.D. degree with the Department of Electrical and Computer Engineering, University of Texas at Austin, Austin, TX, USA, under the supervision of Prof. D. Z. Pan. Her research interests include optical computing/interconnect, neuromorphic computing and logic synthesis.



**Chenghao Feng** received the B.S. degree in physics from Nanjing University, Nanjing, China, in 2018. He is currently working toward the Ph.D. degree with the Department of Electrical and Computer Engineering, University of Texas at Austin, Austin, TX, USA. His research interests include silicon photonics devices and system design for optical computing and interconnect in integrated photonics.



**Zhoufeng Ying** received the B.E. and M.E. degrees in optical engineering from Nanjing University, Nanjing, China, in 2014 and 2016, respectively. He is currently working toward the Ph.D. degree with the Department of Electrical and Computer Engineering, University of Texas at Austin, Austin, TX, USA. His research interests include optical computing and interconnect in integrated photonics.



**Mingjie Liu** Mingjie Liu received his B.S degree from Peking University and M.S. degree from the University of Michigan, Ann Arbor in 2016 and 2018, respectively. He is currently pursuing his Ph.D. degree in Electrical and Computer Engineering at The University of Texas at Austin. His current research interests include applied machine learning for design automation, and physical design automation for analog and mixed-signal integrated circuits.



**David Z. Pan** (S'97-M'00-SM'06-F'14) received his B.S. degree from Peking University, and his M.S. and Ph.D. degrees from University of California, Los Angeles (UCLA). From 2000 to 2003, he was a Research Staff Member with IBM T. J. Watson Research Center. He is currently Engineering Foundation Professor at the Department of Electrical and Computer Engineering, The University of Texas at Austin, Austin, TX, USA and holds the Silicon Laboratories Endowed Chair in Electrical Engineering. His research interests include cross-layer nanometer

IC design for manufacturability, reliability, security, machine learning and hardware acceleration, design/CAD for analog/mixed signal designs and emerging technologies. He has published over 360 journal articles and refereed conference papers, and is the holder of 8 U.S. patents.

He has served as a Senior Associate Editor for ACM Transactions on Design Automation of Electronic Systems (TODAES), an Associate Editor for IEEE Transactions on Computer Aided Design of Integrated Circuits and Systems (TCAD), IEEE Transactions on Very Large Scale Integration Systems (TVLSI), IEEE Transactions on Circuits and Systems PART I (TCAS-I), IEEE Transactions on Circuits and Systems PART II (TCAS-II), IEEE Design & Test, Science China Information Sciences, Journal of Computer Science and Technology, IEEE CAS Society Newsletter, etc. He has served in the Executive and Program Committees of many major conferences, including DAC, ICCAD, ASPDAC, and ISPD. He is the ASPDAC 2017 Program Chair, ICCAD 2018 Program Chair, DAC 2014 Tutorial Chair, and ISPD 2008 General Chair

He has received a number of prestigious awards for his research contributions, including the SRC Technical Excellence Award in 2013, DAC Top 10 Author in Fifth Decade, DAC Prolific Author Award, ASP-DAC Frequently Cited Author Award, 19 Best Paper Awards at premier venues (ASPDAC 2020, ISPD 2020, DAC 2019, GLSVLSI 2018, VLSI Integration 2018, HOST 2017, SPIE 2016, ISPD 2014, ICCAD 2013, ASPDAC 2012, ISPD 2011, IBM Research 2010 Pat Goldberg Memorial Best Paper Award, ASPDAC 2010, DATE 2009, ICICDT 2009, SRC Techcon in 1998, 2007, 2012 and 2015) and 15 additional Best Paper Award finalists, Communications of the ACM Research Highlights (2014), UT Austin RAISE Faculty Excellence Award (2014), and many international CAD contest awards, among others. He is a Fellow of IEEE and SPIE.



**Ray T. Chen** (M91SM98F04) received the B.S. degree in physics from the National Tsing Hua University, Hsinchu, Taiwan, in 1980, and the M.S. degree in physics and Ph.D. degree in electrical engineering from the University of California, in 1983 and 1988, respectively. He is the Keys and Joan Curry/Cullen Trust Endowed Chair with the University of Texas at Austin (UT Austin), Austin, TX, USA. He is the Director of the Nanophotonics and Optical Interconnects Research Lab, Microelectronics Research Center. He is also the Director of the AFOSR MURI-Center for Silicon Nanomembrane involving faculty from Stanford, UIUC, Rutgers, and UT Austin. In 1992, he joined the UT Austin to start the optical interconnect research program. From 1988 to 1992, he worked as a Research Scientist, Manager, and Director of the Department of Electro-Optic Engineering, Physical Optics Corporation, Torrance, CA, USA. From 2000 to 2001, he served as the CTO, founder, and Chairman of the Board of Radianit Research, Inc., where he raised 18 million dollars A-Round funding to commercialize polymer-based photonic devices involving more than 20 patents, which were acquired by Finisar in 2002, a publicly traded company in the Silicon Valley (NASDAQ:FNSR). He has been also serving as the founder and Chairman of the Board of Omega Optics Inc., Austin, TX, USA, since its initiation in 2001. Omega Optics has received more than five million dollars in research funding. His research work has been awarded more than 135 research grants and contracts from sponsors such as Army, Navy, Air Force, DARPA, MDA, NSA, NSF, DOE, EPA, NIST, NIH, NASA, the State of Texas, and private industry. The research topics are focused on four main subjects: nanophotonic passive and active devices for bio- and EM-wave sensing and interconnect applications; thin film guided-wave optical interconnection and packaging for 2-D and 3-D laser beam routing and steering; true time delay wide band phased array antenna; and 3-D printed microelectronics and photonics. Experiences garnered through these programs are pivotal elements for his research and further commercialization. His group at UT Austin has reported its research findings in more than 900 published papers, including more than 100 invited papers. He holds more than 60 patents. He has supervised and graduated 51 Ph.D. students from his research group at UT Austin. Many of them are currently professors in the major research universities in USA and abroad. Dr. Chen has chaired or been a program-committee member for more than 120 domestic and international conferences organized by IEEE, SPIE (The International Society of Optical Engineering), OSA, and PSC. He has served as an Editor, Co-editor, or co-author for more than 25 books. He has also served as a Consultant for various federal agencies and private companies and delivered numerous invited/plenary talks to professional societies. He is a Fellow of OSA and SPIE. He was the recipient of the 1987 UC Regents Dissertation Fellowship and the 1999 UT Engineering Foundation Faculty Award, for his contributions in re- search, teaching, and services. He was the recipient of the Honorary Citizenship Award in 2003 from the Austin city council for his contribution in community service. He was also the recipient of the 2008 IEEE Teaching Award, and the 2010 IEEE HKN Loudest Professor Award, and 2013 NASA Certified Technical Achievement Award for contribution on moon surveillance conformable phased array antenna. During his undergraduate years at the National Tsing Hua University, he led the 1979 university debate team to the Championship of the Taiwan College-Cup Debate Contest.