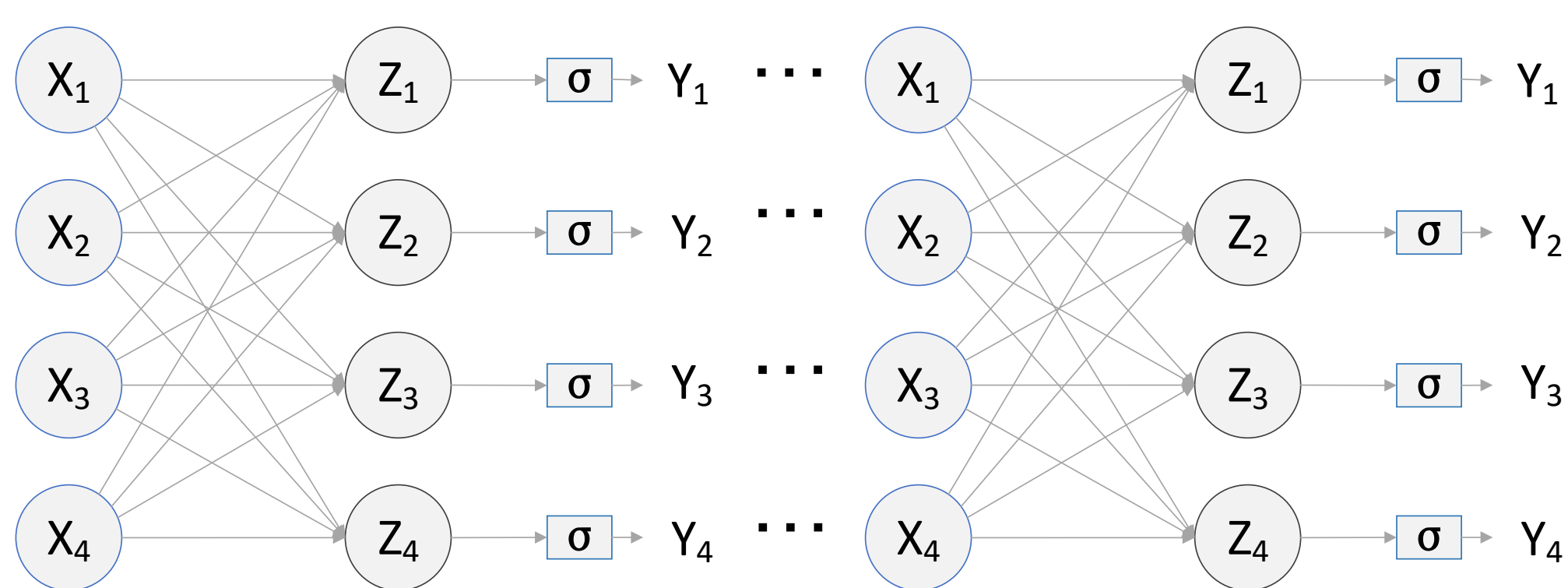# Towards Area-Efficient Optical Neural Networks:
# An FFT-based Architecture

Jiaqi Gu[1], Zheng Zhao[1], Chenghao Feng[1], Mingjie Liu[1], Ray T. Chen[1], David Z. Pan[1]
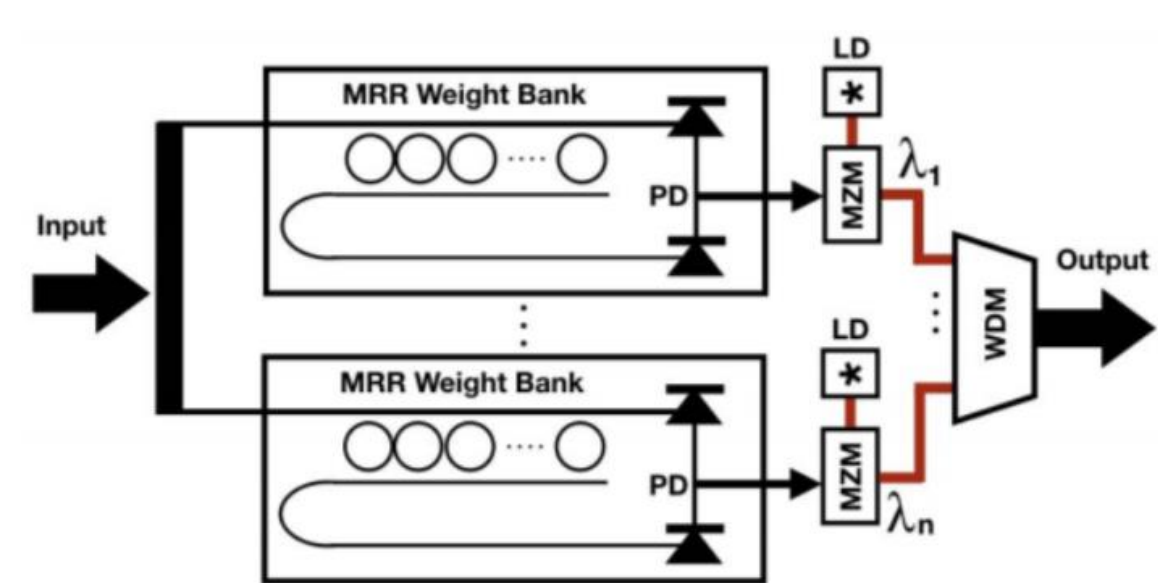
[1]University of Texas at Austin

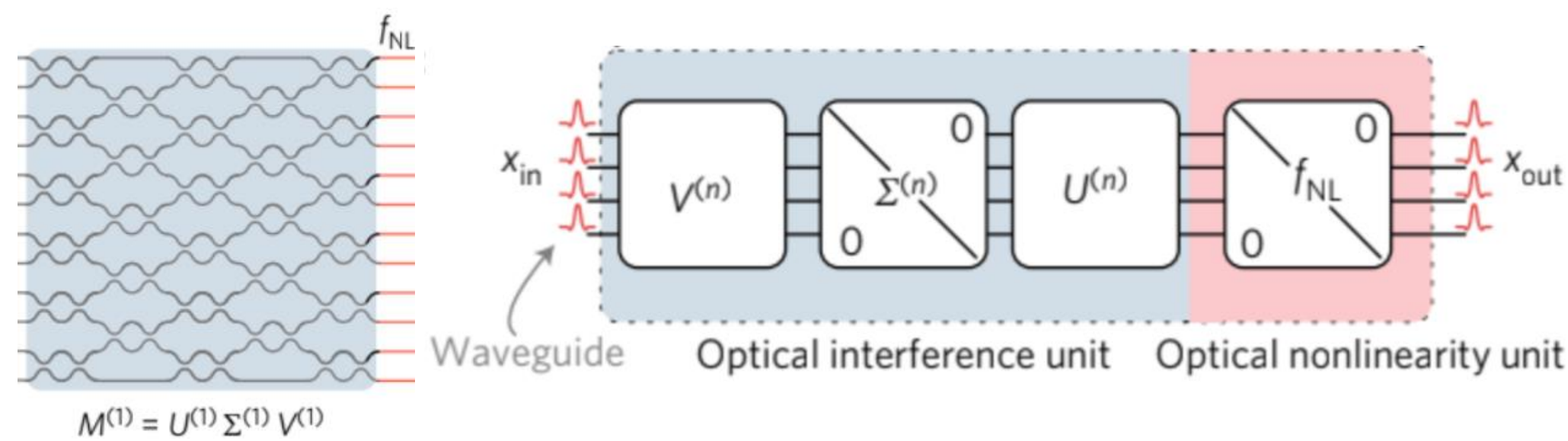## Multi-layer Perceptron Inference

- Input
  - Vector $x$
- Output
  - Vector $y = \sigma(W \cdot x)$
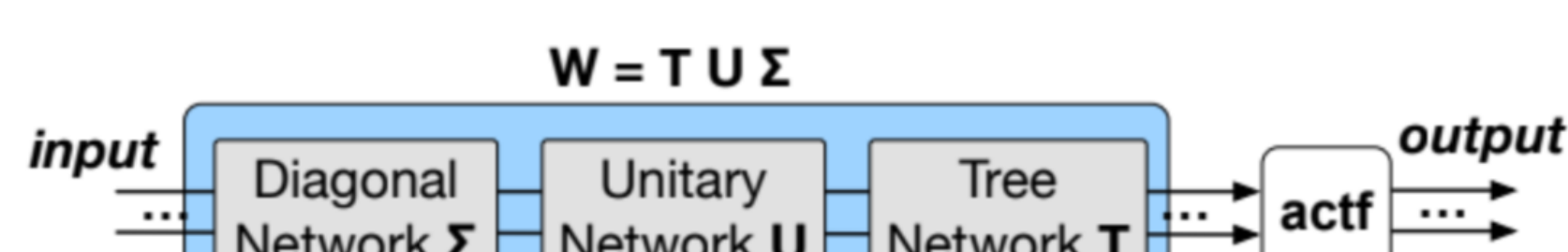- Objective
  - Accuracy



## Previous Work on ONNs



- Photonic microring resonator banks
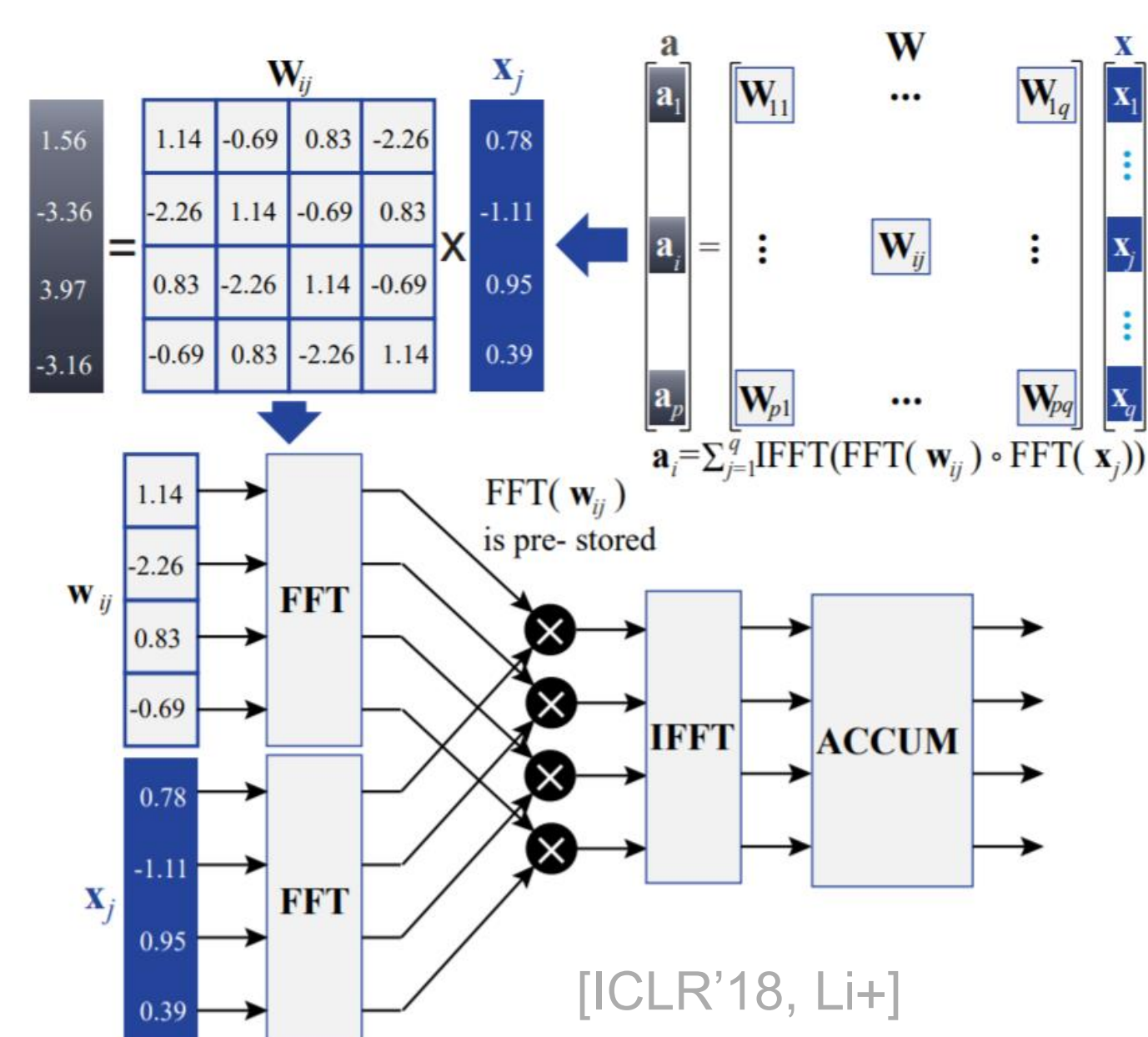
[IEEE SOCC'18, Mehrabian+]



$M^{(1)} = U^{(1)}\Sigma^{(1)}V^{(1)}$

- SVD-based ONNs with MZIs

[Nature'17, Shen+]



$$W = T U \Sigma$$

- $T\Sigma U$-based ONNs with MZIs and sparse tree

[ASPDAC'19, Zhao+]

## Background on Structured Neural Networks

- **ERNN using SNN** [ICLR'18,Li+]
  - Low computational complexity
  - Low storage complexity
  - ~8x parameter reduction
  - <0.5% accuracy degradation

- **Theoretical Proof** [ICML'17,Zhao+]
  - Universal approximation
  - Identical error bound as classical NNs



$a_i = \sum_{j=1}^d \text{IFFT}(\text{FFT}(W_{ij}) * \text{FFT}(x_j))$

FFT($W_{ij}$) is pre-stored

[ICLR'18, Li+]

$$y = Wx \quad\longleftrightarrow\quad y = \mathcal{F}^{-1}(\mathcal{F}(w) \odot \mathcal{F}(x))$$
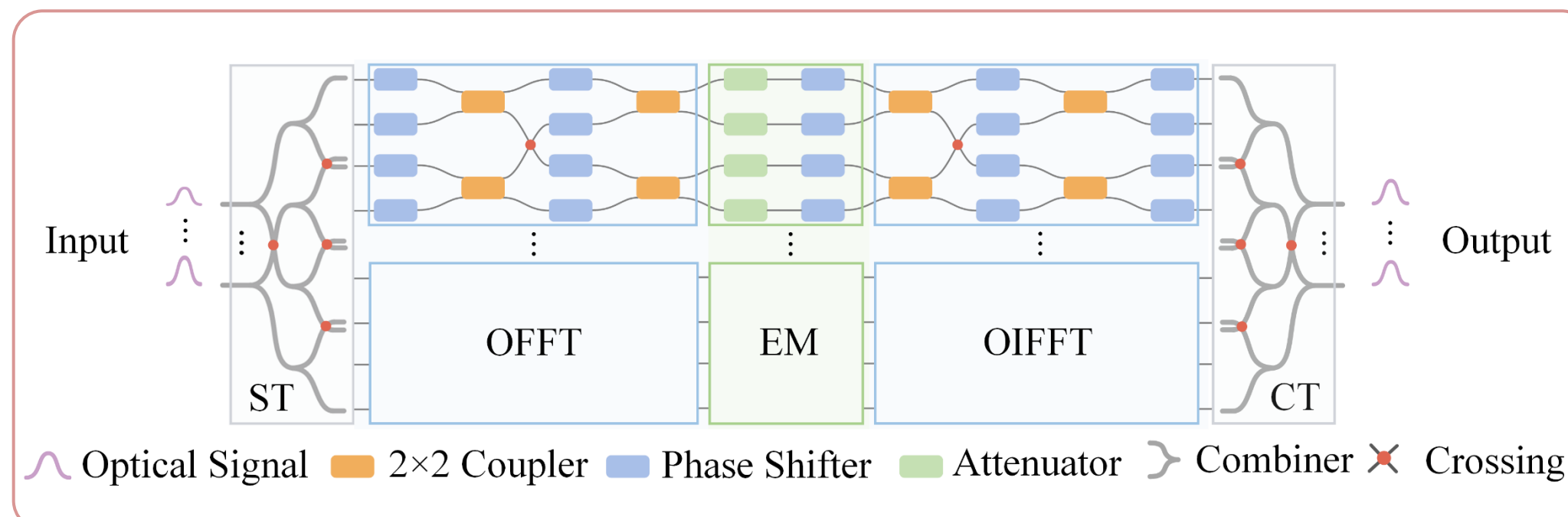
## Advances in Network Pruning

Network slimming with pruning techniques



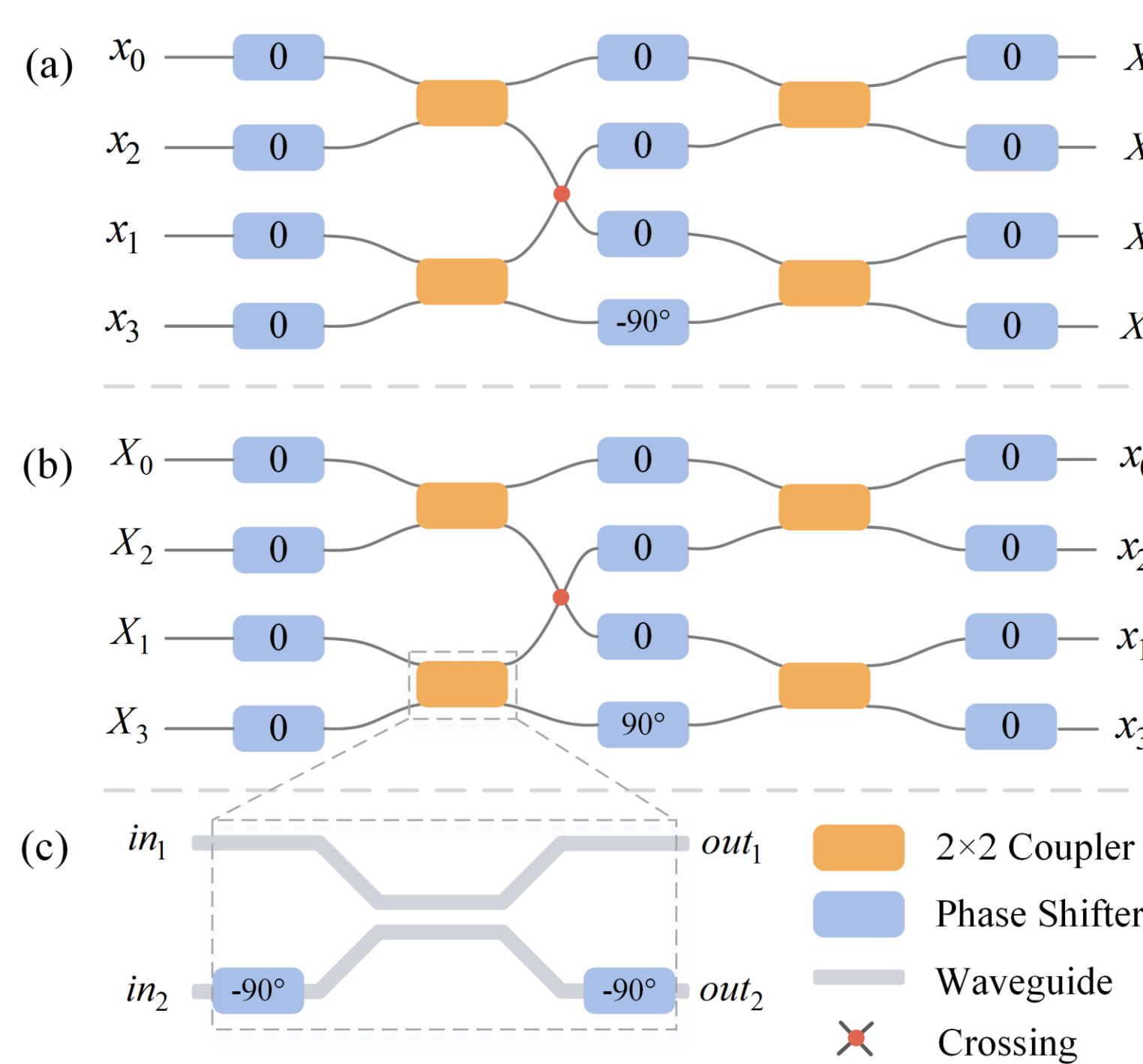**Non-structured pruning**
- Random zero entries
- Irregular

**Structured pruning**
- Zero entries in group
- Regular
- Hardware-friendly

## Proposed ONN Architecture



Input / Output / ST / OFFT / EM / OIFFT / CT

Optical Signal / 2×2 Coupler / Phase Shifter / Attenuator / Combiner / Crossing

## Optical Fast Fourier Transform

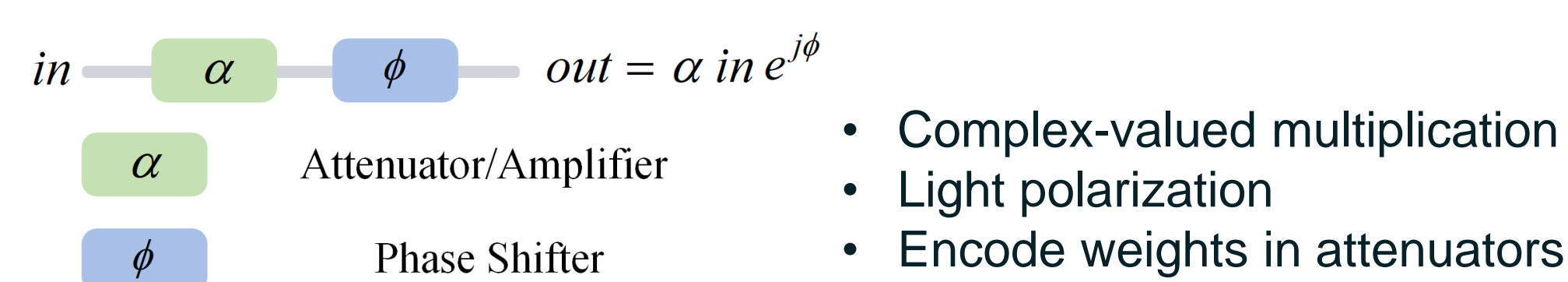$$X_k = \frac{1}{\sqrt{N}} \sum_{n=0}^{N-1} x_n e^{-i\frac{2\pi kn}{N}} \quad k = 0, 1, \cdots, N-1.$$



(a)
(b)
(c)

2×2 Coupler / Phase Shifter / Waveguide / Crossing

$$\binom{out_1}{out_2} = \frac{1}{\sqrt{2}}\binom{in_1 + in_2}{in_1 - in_2}$$
$$= \underbrace{\begin{pmatrix} 1 & 0 \\ 0 & -j \end{pmatrix}}_{\text{output phase shifter}} \underbrace{\frac{1}{\sqrt{2}}\begin{pmatrix} 1 & j \\ j & 1 \end{pmatrix}}_{\text{directional coupler}} \underbrace{\begin{pmatrix} 1 & 0 \\ 0 & -j \end{pmatrix}}_{\text{input phase shifter}} \binom{in_1}{in_2}$$

- 2 × 2 couplers and phase shifters to achieve OFFT

## Element-wise Vector Multiplication



$in \to \alpha \to \phi \to out = \alpha \, in \, e^{j\phi}$

- $\alpha$ Attenuator/Amplifier
- $\phi$ Phase Shifter

- Complex-valued multiplication
- Light polarization
- Encode weights in attenuators

## Splitter/Combiner Tree
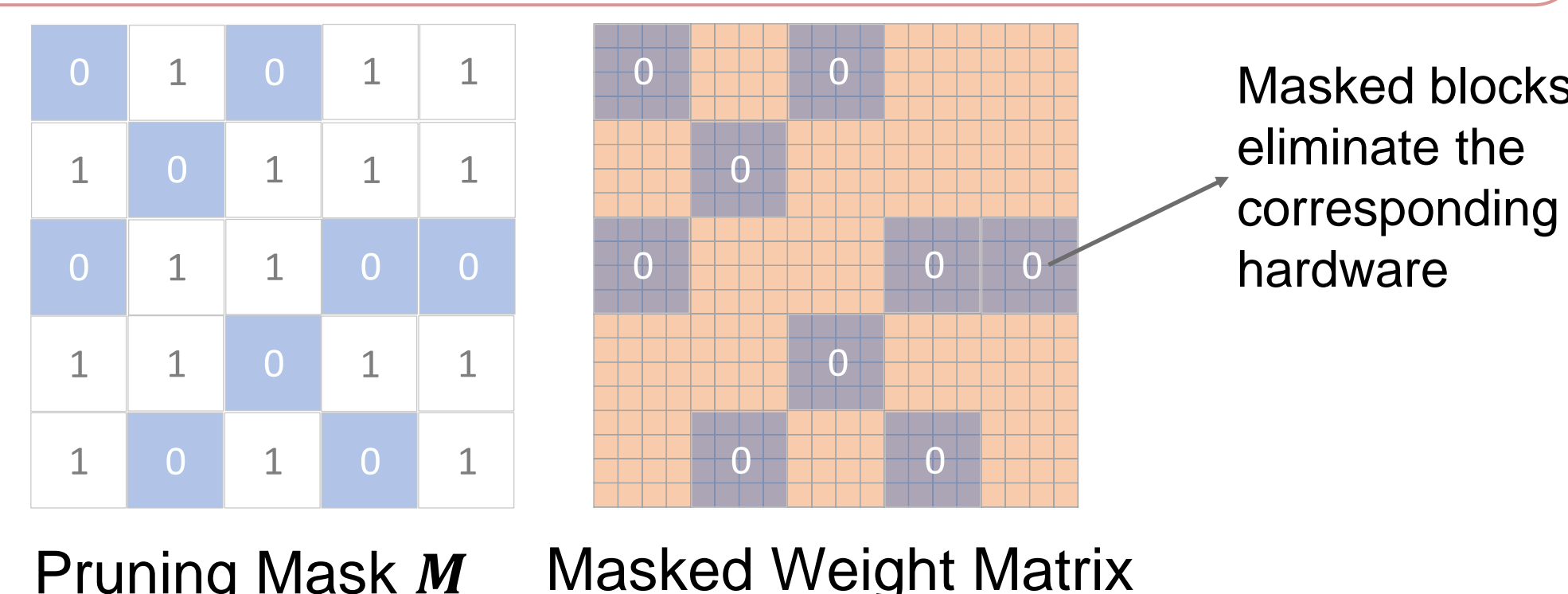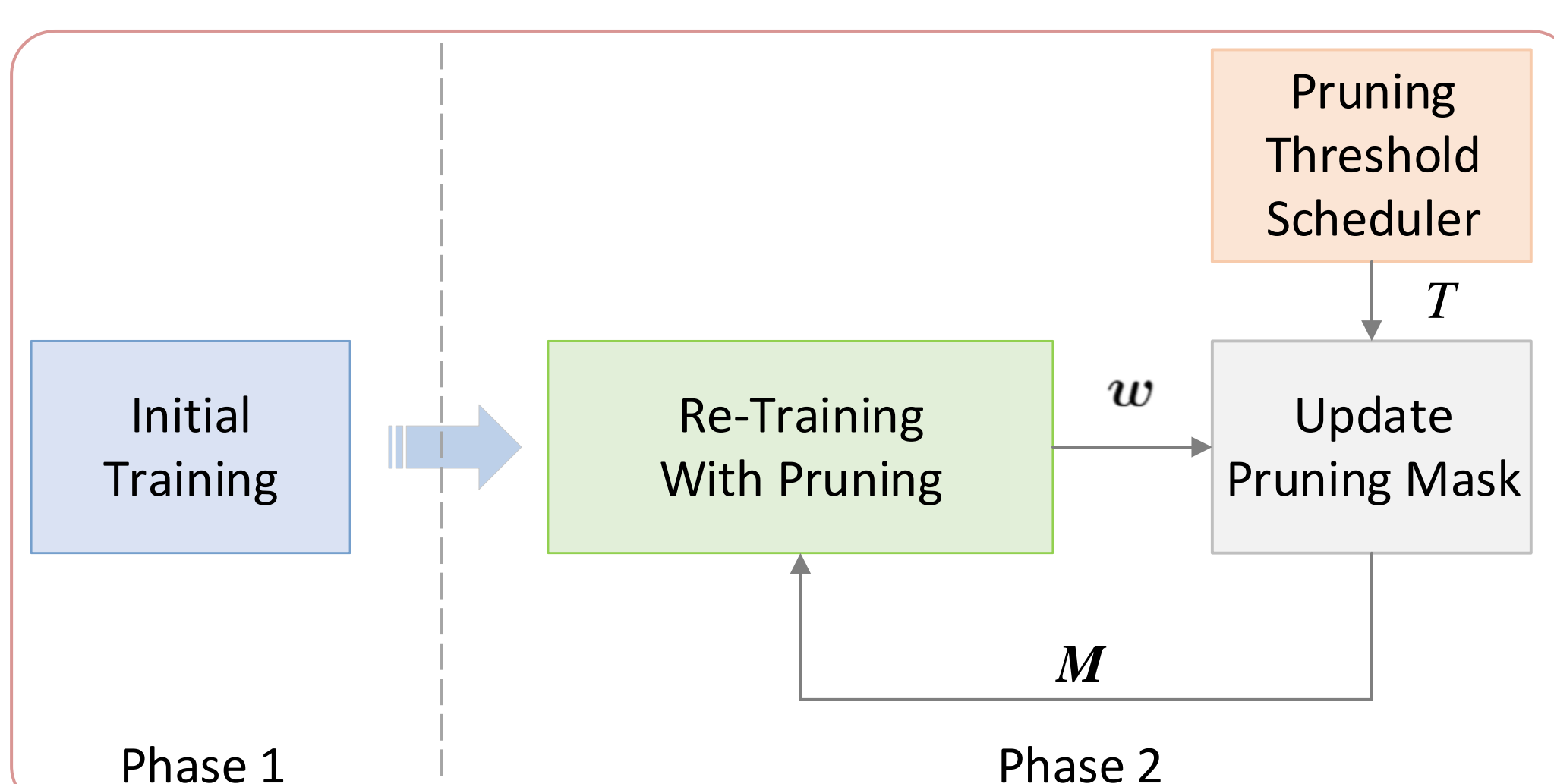


$$y_i = \sum_{j=0}^{q-1} W_{ij} x_j$$

Combiner / Waveguide Crossing

- Fewer waveguide crossings: $k(k-1)(q-1)/2$
- Avoid multi-port combiners

## Two-Phase Training Flow with Structured Pruning



Initial Training / Re-Training With Pruning / Pruning Threshold Scheduler / Update Pruning Mask

Phase 1 / Phase 2

$w$ / $T$ / $M$



Pruning Mask $M$ — Masked Weight Matrix

Masked blocks eliminate the corresponding hardware

## Hardware Utilization Analysis

**SVD-based Architecture ($W \in \mathbf{R}^{m \times n}$)**

$$\#DC_{\text{SVD}} = m(m-1) + n(n-1) + \max(m,n)$$
$$\#PS_{\text{SVD}} = \frac{m(m+1)}{2} + \frac{n(n-1)}{2}$$

**T$\Sigma$U-based Architecture ($W \in \mathbf{R}^{m \times n}$)**

$$\#DC_{\text{T}\Sigma\text{U}} = n(n+1) + \max(m,n)$$
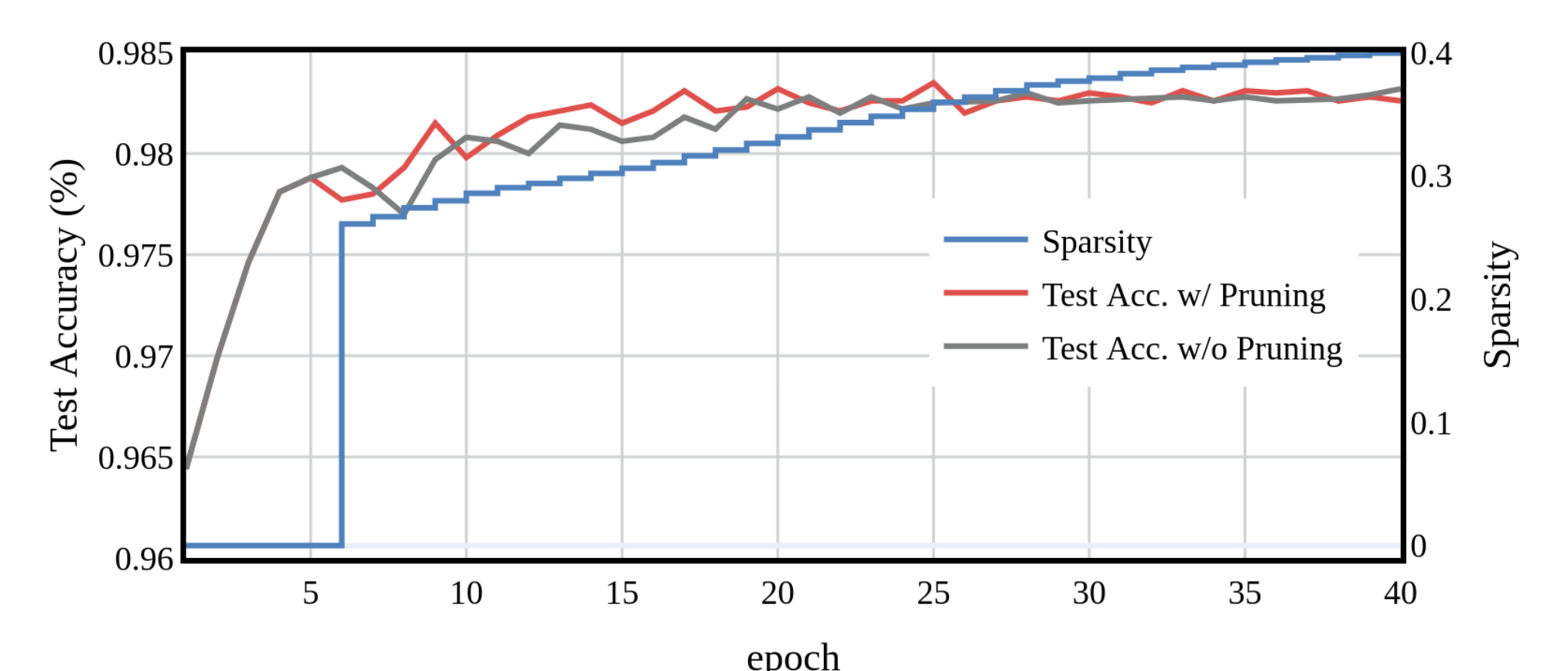$$\#PS_{\text{T}\Sigma\text{U}} = \frac{n(n+1)}{2}$$

**Ours Architecture ($W \in \mathbf{R}^{m \times n}, block = k$)**
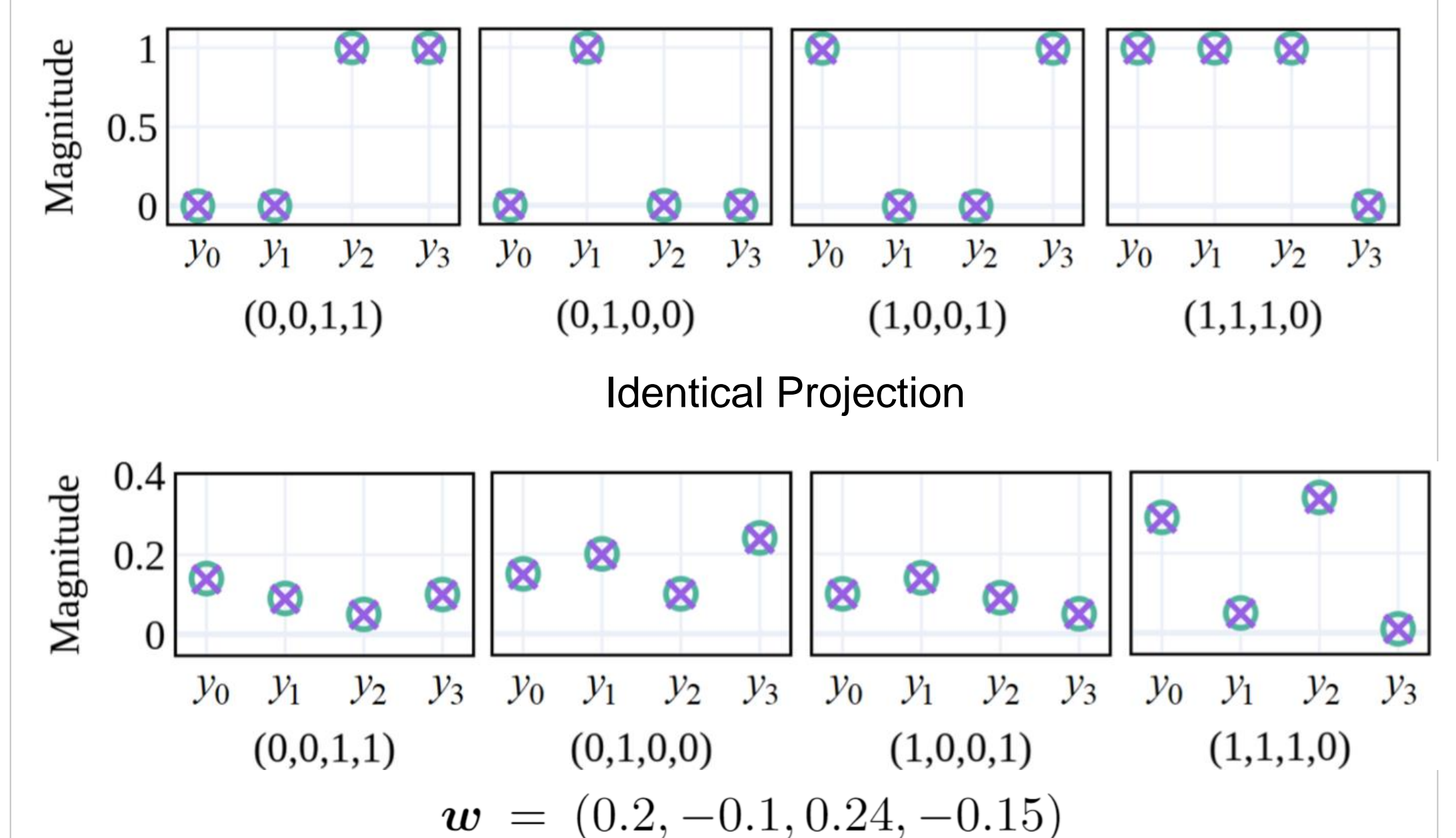
$$\#DC_{\text{Ours}} = \frac{mn}{k}(\log_2 k + 1)$$
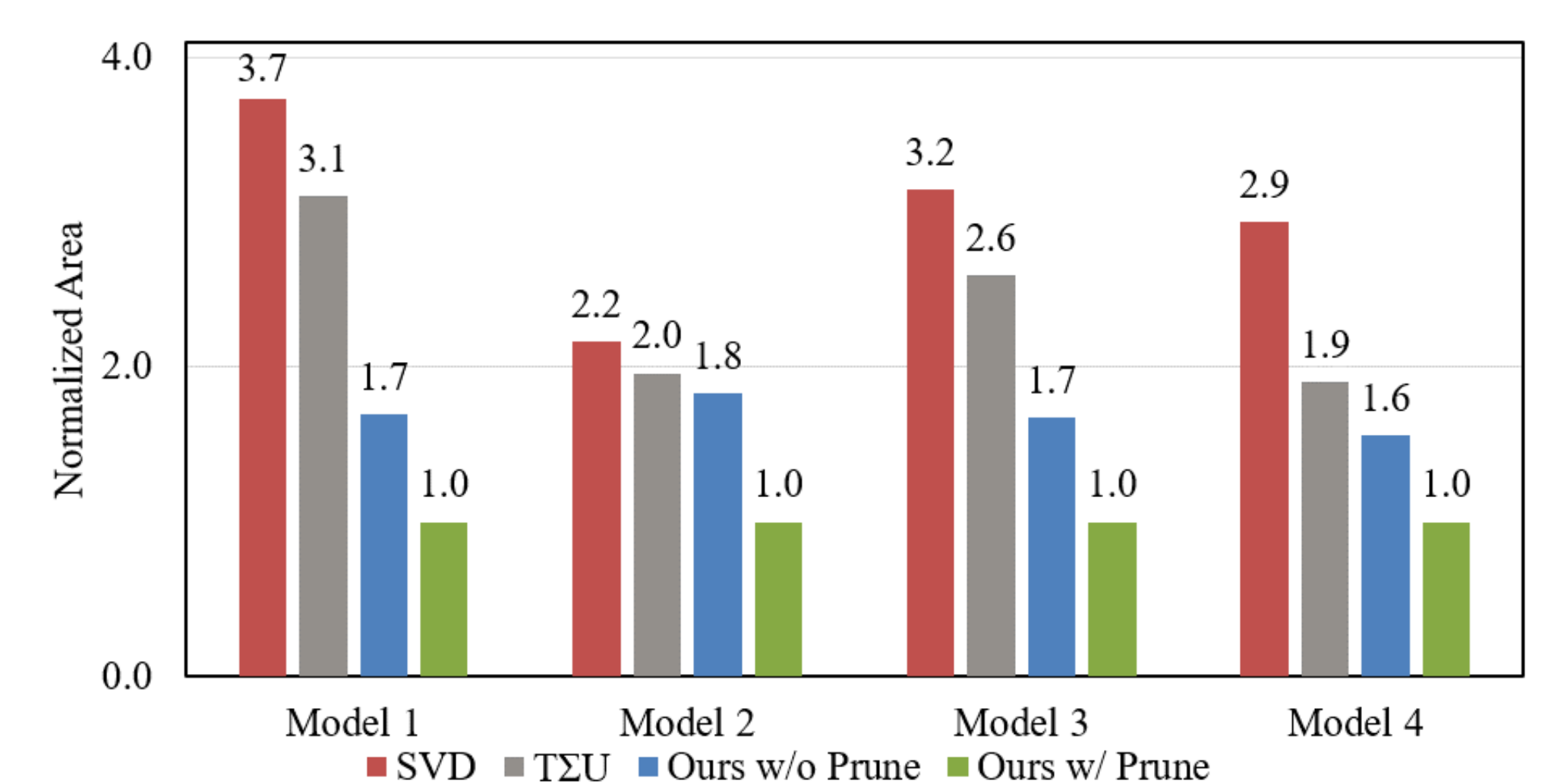$$\#PS_{\text{Ours}} = \frac{mn}{k}(2\log_2 k + 1)$$

## Experimental Results

**Training Curve**



Sparsity / Test Acc. w/ Pruning / Test Acc. w/o Pruning

**Numerical Simulation (Lumerical)**



Identical Projection

(0,0,1,1) (0,1,0,0) (1,0,0,1) (1,1,1,0)

$$w = (0.2, -0.1, 0.24, -0.15)$$

**Area Comparison**



SVD / T$\Sigma$U / Ours w/o Prune / Ours w/ Prune

**Architecture innovation**
- Use circulant matrix representation for better efficiency
- Avoid using MZIs
- Friendly to modern network pruning techniques.
- Fewer parameters

**Software innovation**
- End-to-end training flow to perform structured pruning based on Group Lasso regularization
- Incremental method avoid accuracy degradation

## Conclusion and Future Work

- New architecture to save optical component for better area efficiency
- Enable structured pruning to optical neural networks for network slimming without accuracy degradation
- 2.2~3.7x better area cost than SVD-based architecture
- May extend to OCNN and other compact NNs
- Considering more practical hardware information

Source code: https://github.com/JeremieMelo/fft-onn