

# ELight: Enabling Efficient Photonic In-Memory Neurocomputing with Life Enhancement

Hanqing Zhu, Jiaqi Gu, Chenghao Feng, Mingjie Liu, Zixuan Jiang,  
Ray T. Chen, David Z. Pan

Dept. of Electrical and Computer Engineering  
The University of Texas at Austin

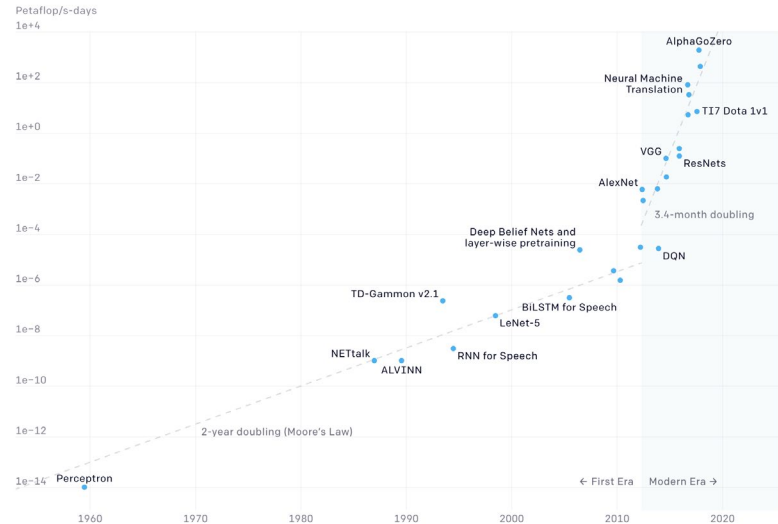
*This work is supported in part by MURI and ONR*

# AI Acceleration and Challenges

- ◆ ML models and dataset keep increasing
  - › Low latency
  - › Low power
  - › High bandwidth



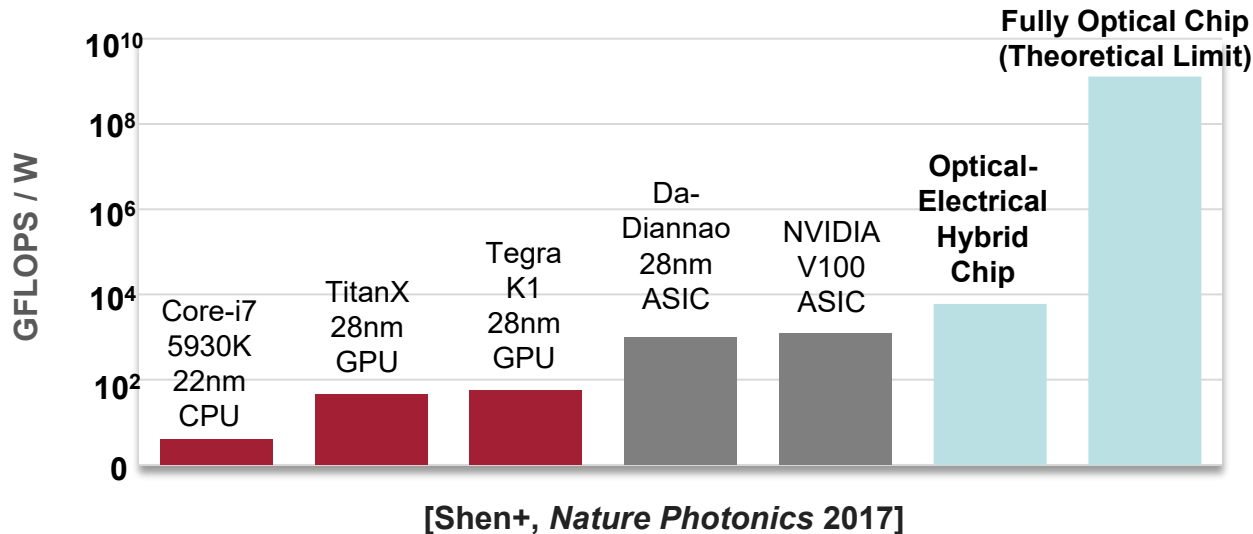
- ◆ AI compute requirement has **5×** the doubling rate of Moore's law



# Photonic AI

- ◆ Use light to continue Moore's law
- ◆ Promising technology for next-generation AI accelerator

## Ultra-high speed & Ultra-low energy



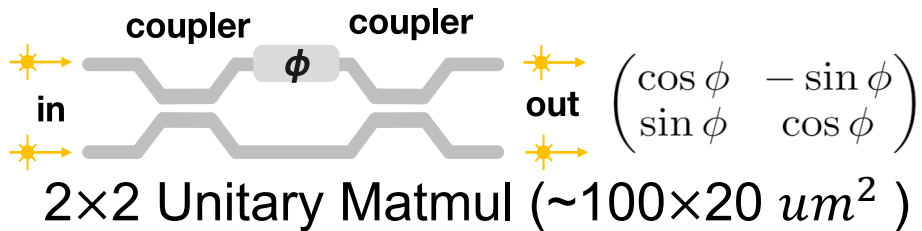
LIGHTMATTER



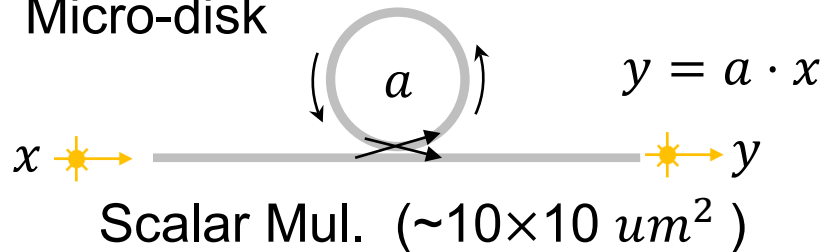
# Optical Computing Basics

## Computing Basics

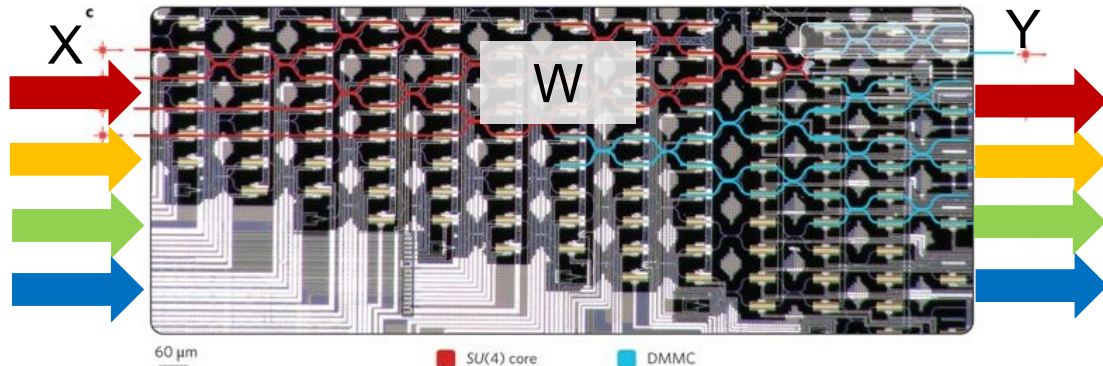
### MZI



### Micro-ring/ Micro-disk

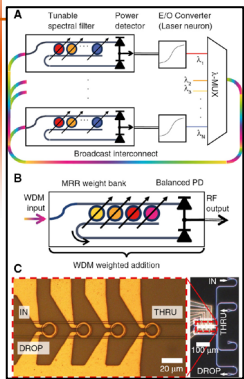


## Optical Neural Network

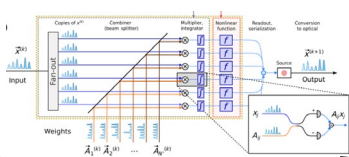


**Photonic tensor unit** for analog GEMM [MIT's Nat. Photonics'17]

# ONN Progress

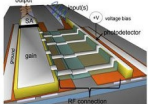


**MRR Neural Network**  
[Brunner+, 2016]  
[Tait+, SciRep 2017]  
Princeton

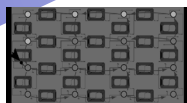


**Holylight and Lightbulb: MRR&PCM**  
[Liu+, Zokaee+ DATE'2019, 2020]  
NTU

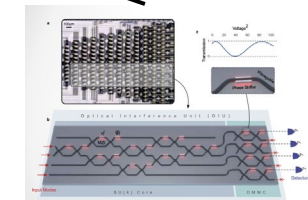
**Quantum ONN**  
[Hamerly+, PhysRev2019]  
MIT



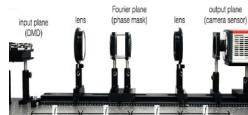
**Optical Spike Neural Network**  
[Tait+, 2016]  
Princeton



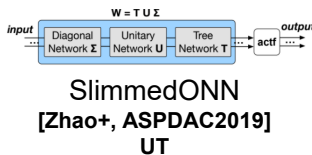
**Optical Reservoir Computing**  
[Vandoorne+, NatureComm 2014]  
Ghent University



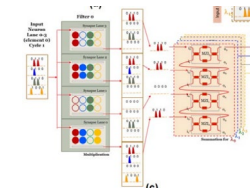
**MZI-based Neural Network**  
[Shen+, Nature Photonics 2017]  
MIT



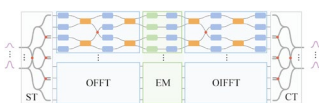
**Free-space ONN**  
[Chang+, SciRep 2018]  
Stanford



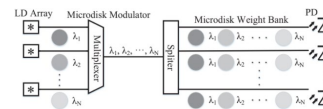
**SlimmedONN**  
[Zhao+, ASPDAC2019]  
UT



**PIXEL, MZI Multiplier**  
[Shieffelt+, HPCA2020]  
Ohio Univ



**FFT-based optical neural network**  
[Gu+, ASPDAC2020]  
UT



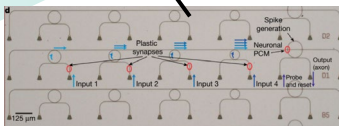
$$P_{out} = x \cdot P_y + \bar{x} \cdot P_y$$

$$\lambda = \lambda_{on} \quad \bar{\lambda} = \lambda_{off}$$

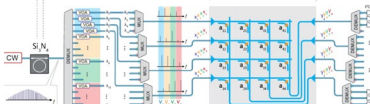
$$x = 0 \quad y = \lambda_{on}$$

$$x = 1 \quad \bar{y} = \lambda_{off}$$

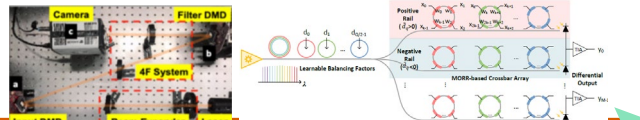
**Holylight and Lightbulb: MRR&PCM**  
[Liu+, Zokaee+ DATE'2019, 2020]  
NTU



**Spiking ONN: PCM**  
[Feldmann+, Nature 2019]  
Munster, Oxford



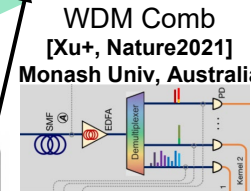
**PCM Xbar**  
[Feldmann+, Nature2021]  
Munster, Oxford



**Free-space ONN**  
[Miscuglio+, Optica2020]  
GWU



**MORR ONN**  
[Gu+, DATE2021]  
UT



**WDM Comb**  
[Xu+, Nature2021]  
Monash Univ, Australia

**PCM Xbar**  
[Miscuglio+, APR2020]  
GWU

2016

2017

2018

2019

2020

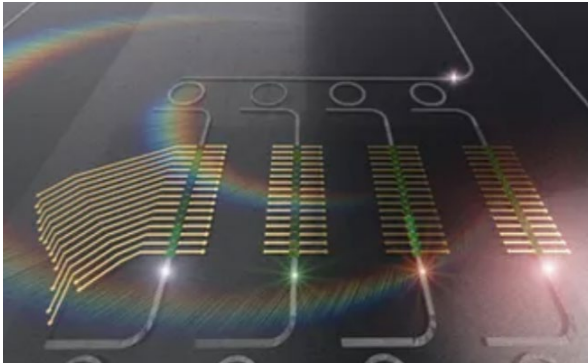
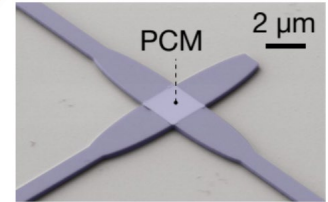
2021

# Computing with Photonic Phase Change Material

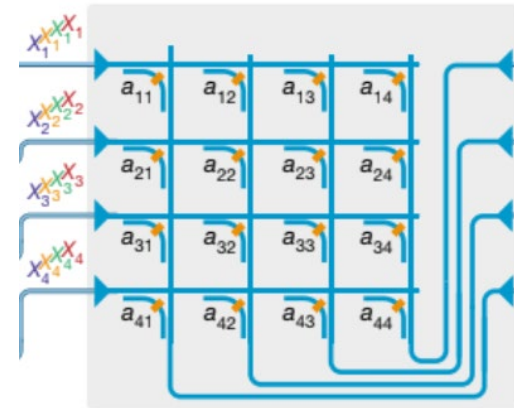
- ◆ Recall photonic devices: MZI, Microring, micro-disks, ...
- ◆ A new device for non-volatile computing
  - › Phase change material (PCM)
  - › **Modulate the light transmission** to achieve multiplication
  - › Store the transmission as a **non-volatile** memory



$$y = w \cdot x$$



[Miscuglio+, APR'20]



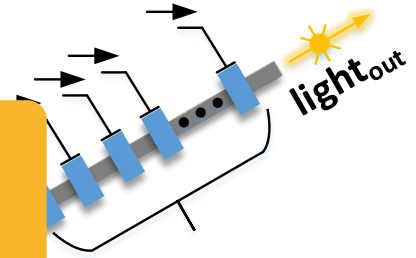
[Feldmann+, Nature'21]

# Barriers Towards Practical Deployment

## ◆ Limited PCM rewrite!

- › Max:  $10^6 - 10^8$  writes
- › Aging

PCM has *limited rewrite times*  
*Short endurance*



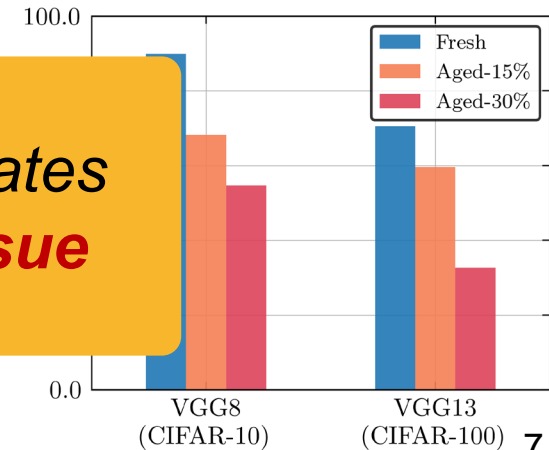
$2^n - 1$  PCM wires

1-bit PCM memory cell  
[Miscuglio +, APR'20]

## ◆ Why we care inference endurance?

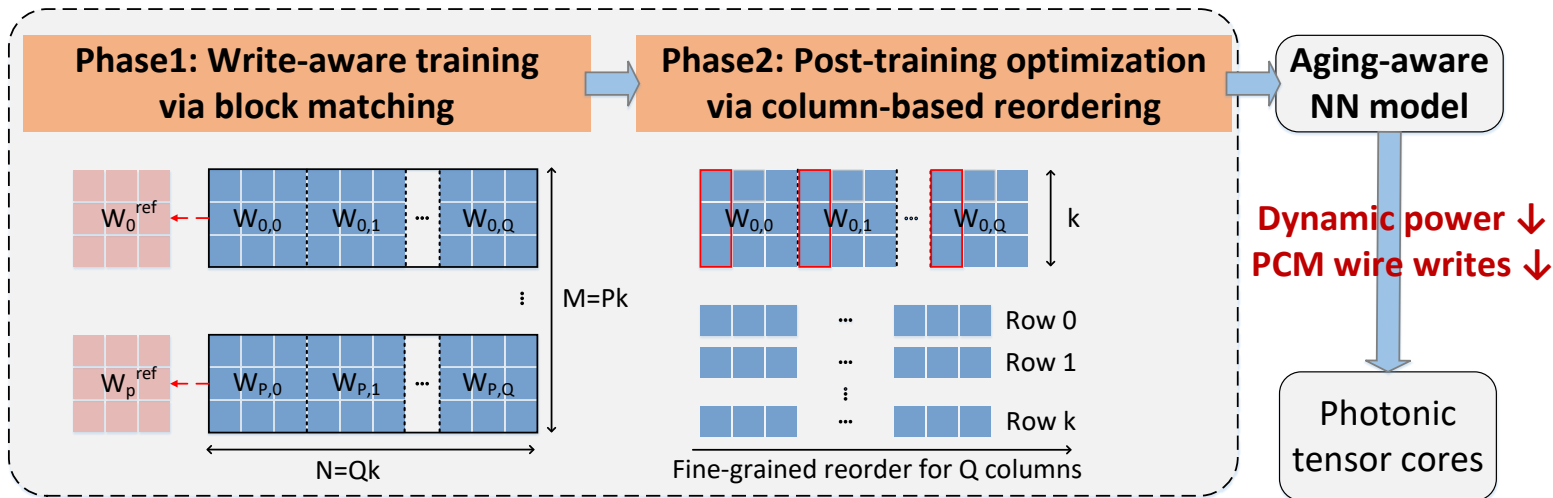
- › Limited
- › Limited
- › Need f

*Frequent reuse further escalates*  
*Dynamic power & Aging issue*



# Our Proposed Aging-aware PCM-ONN: ELight

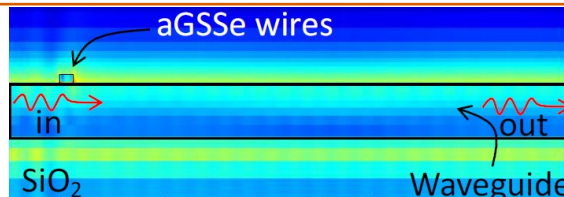
- ◆ A two-phase aging-aware optimization framework
  - › **Minimize PCM write** operations in **inference**
- ◆ Achieved
  - › **> 20** × fewer write operations
  - › **Minimized** # max writes on a single PCM cell
  - › **> 30** × less dynamic energy cost





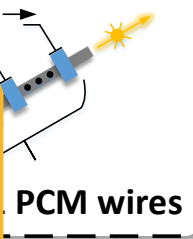
# Distribution-aware Quantization

- ◆ 1 PCM wire:  $y = cx$
- ◆ b-bit PCM cell:  $y = c^b x$
- ◆ Transmission factor (t) follows power-of-c model



Amorphous state (0)  
Crystalline state (1)

b-bit PCM memory cell

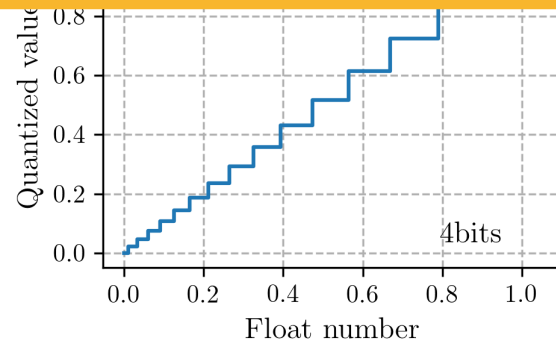
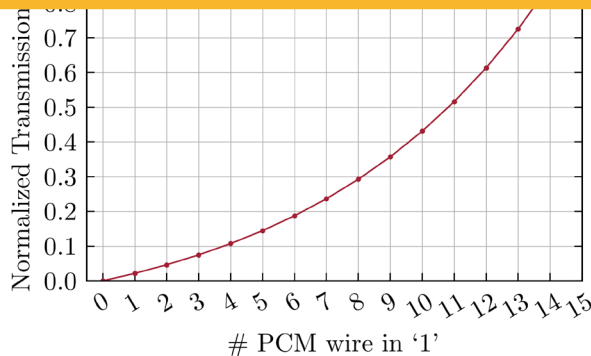


PCM wires

Be Aware of Distribution !

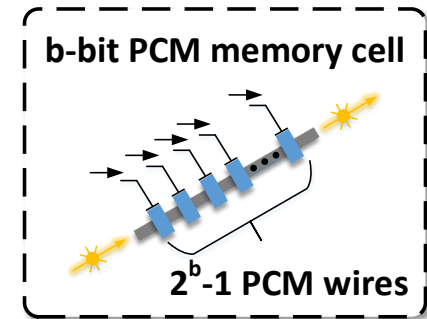
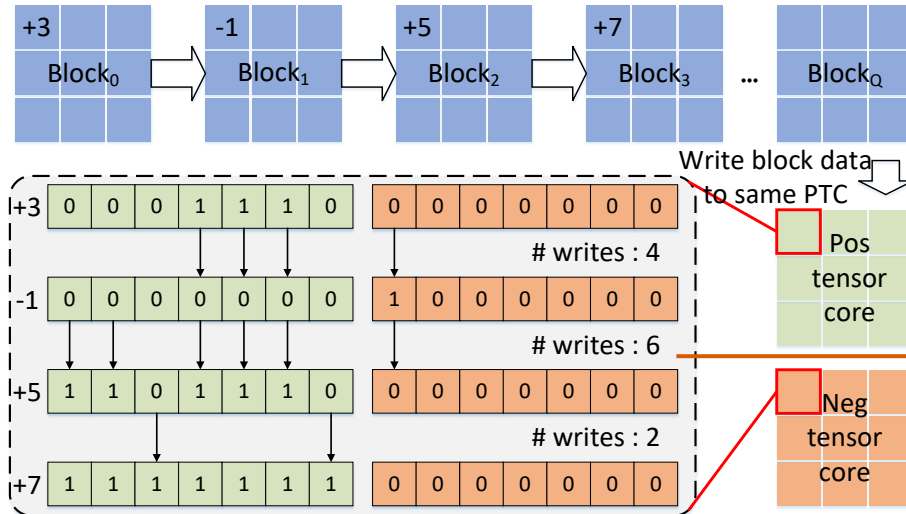
→ **More** quantization levels in **small** weights

→ **boost expressivity**



# Augmented Redundant Write Elimination (ARWE)

- ◆ Block matrix multiplication
- ◆ Assume each PTC is assigned with one row of weight sub-blocks
- ◆ ARWE: Preserve current states at the most
  - › Redundant write elimination scheme [Yang+, ISCAS]
  - › Easy to compare values as weights are known and pre-stored



$$WT(w', w) = |l^+(w') - l^+(w)| + |l^-(w') - l^-(w)|$$

# Write-aware Training: Weight Matching

- ◆ Significant rewrite operations still exists with ARWE!
  - › Deployment of a 5-bit VGG8 model trained on CIFAR10

	Layer 2	Layer 3	Layer 4	Layer 5
<i># total writes</i>	$1.14 \times 10^6$	$4.87 \times 10^6$	$1.66 \times 10^7$	$3.26 \times 10^7$
<i># max writes</i>	294	534	914	1425

- ◆ Pull w

Orchestrate *weight similarity*  
during *training*

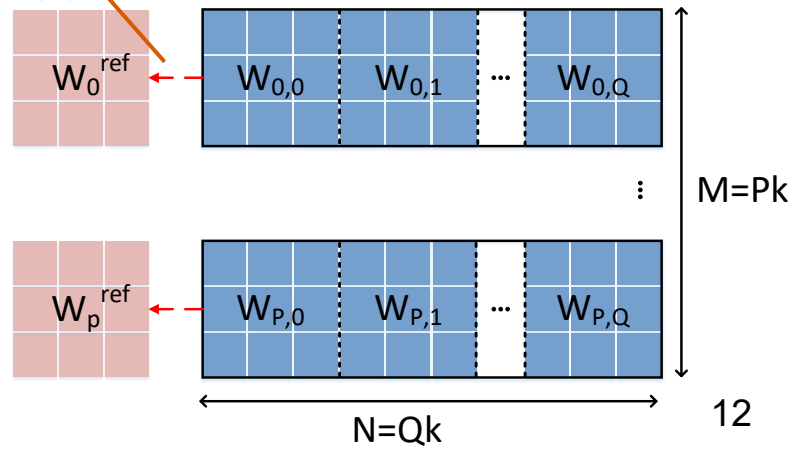


# Write-aware Training: Block Matching

- ◆ Average a group weight blocks into one reference block
- ◆ Compute level difference between two blocks with L2 norm
  - › Penalize large value deviation but Allow slight value deviation
  - › Preserve model expressivity

$$LD(W^{ref}, W) = \sum_i \sum_j^k \|\tilde{l}^+(w_{ij}^{ref}) - \tilde{l}^+(w_{ij})\|^2 + \|\tilde{l}^-(w_{ij}^{ref}) - \tilde{l}^-(w_{ij})\|^2.$$

- ◆ Match weight blocks with the reference block



# Write-aware Training: Optimization issue

- ◆  $\mathcal{L}_{BM}$  is not differentiable

- › Recall the function to get transmission levels

$$l^+(w) = \begin{cases} (2^b - 1) - \text{Clip}(\text{R}(\log_t(s|w| + \delta)), 0, 2^b - 1), & w \geq 0 \\ 0, & w < 0 \end{cases}$$

$$l^-(w) = \begin{cases} 0, & w \geq 0 \\ \text{Clip}(\text{R}(\log_t(s|w| + \delta)), 0, 2^b - 1) - (2^b - 1), & w < 0 \end{cases}$$

STE

- ◆ Not all gradients need to be propagated back

- › Only gradients from physically stored value are valid

Mask out invalid gradients

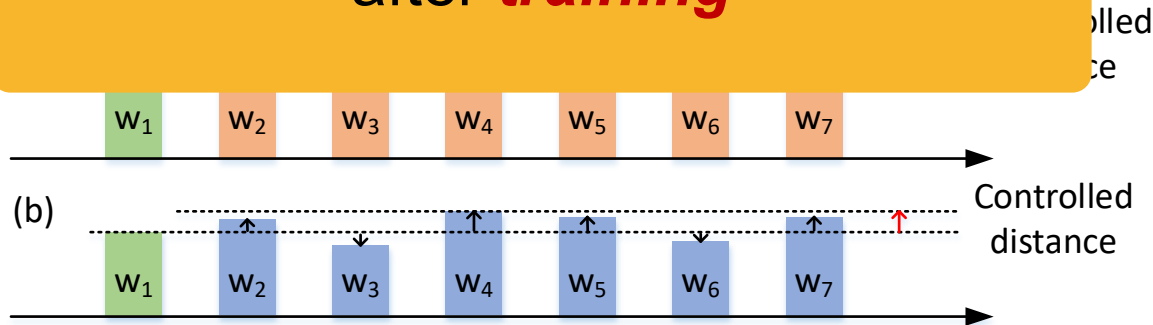
$$\mathcal{L}_{BM} = \sum_l^L \sum_t^G \sum_i^{n_g} \frac{1}{\beta^B} LD(B_i^t, B_{avr}^t)$$

$$LD(B, A) = \sum_i^k \sum_j^k \|\tilde{l}^+(b_{ij}) - \tilde{l}^+(a_{ij})\|^2 + \|\tilde{l}^-(b_{ij}) - \tilde{l}^-(a_{ij})\|^2$$

# Post-training Optimization: Need of Reordering

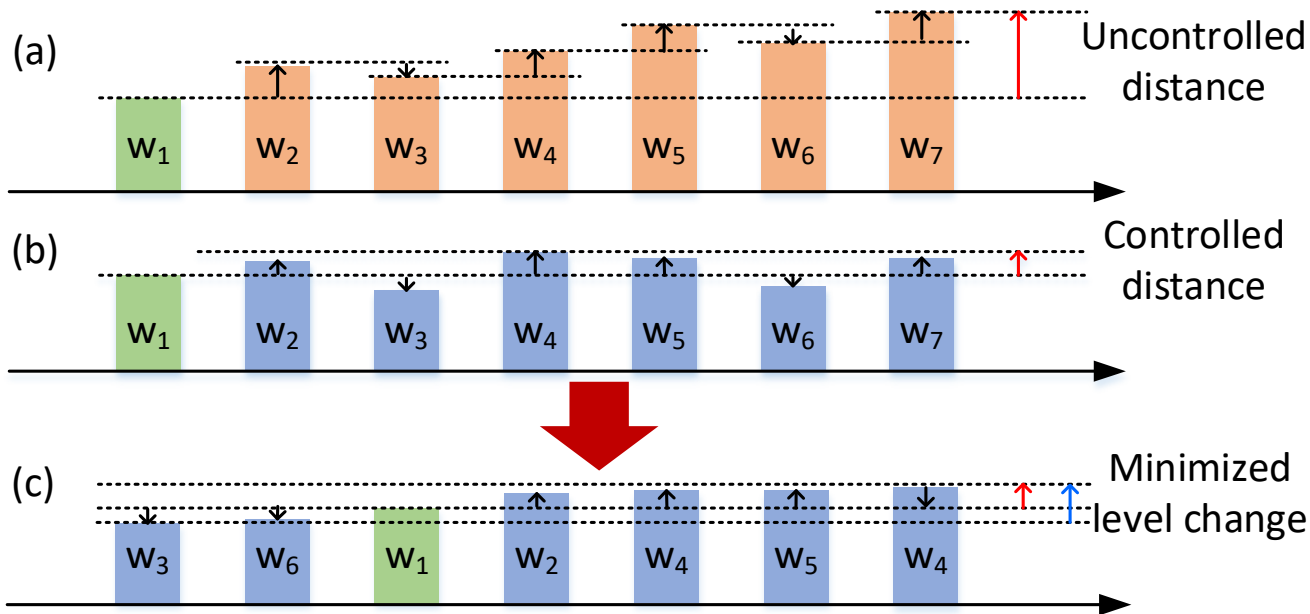
- ◆ We only pull weights close to a reference
- ◆ No order information is introduced during our optimization!
  - › Cannot guarantee best-of-reduction
- ◆ Find an ordering that minimizes the distance

Consider *order information*  
after *training*



# Post-training Optimization: One-shot Reordering

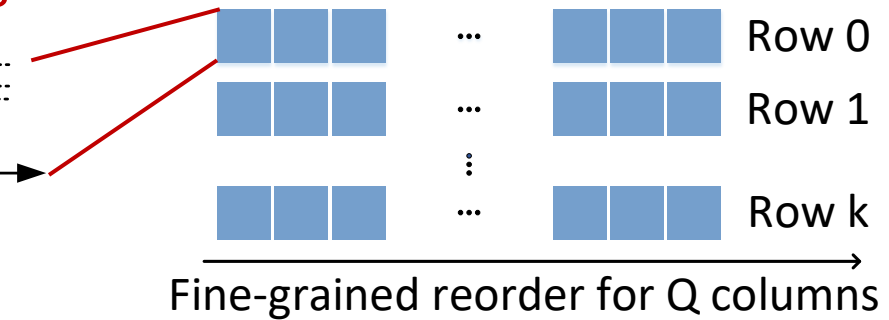
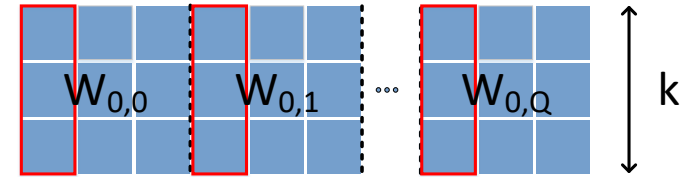
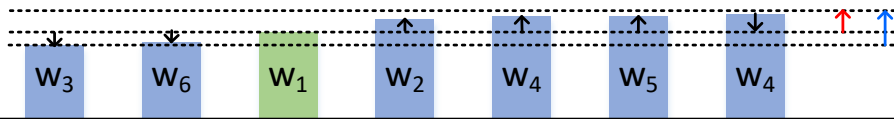
- ◆ Simply reorder the weight sequences
- ◆ Efficient reordering with negligible overhead



# Post-training Optimization: One-shot Reordering

- ◆ Take one more step before real deployment
- ◆ One-shot recording concurrently for different columns
- ◆ No affect on the computation results
  - › General matrix multiplication

Re-order weights to minimize # writes





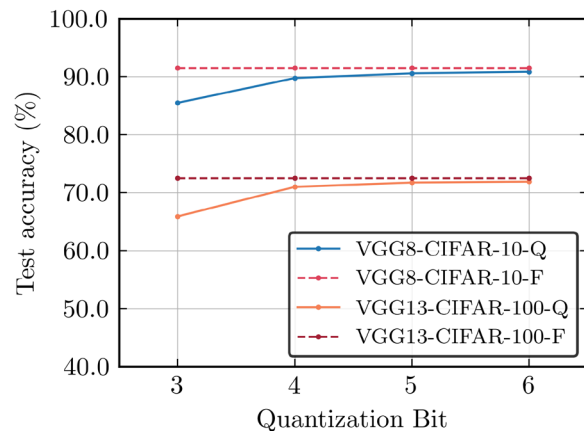
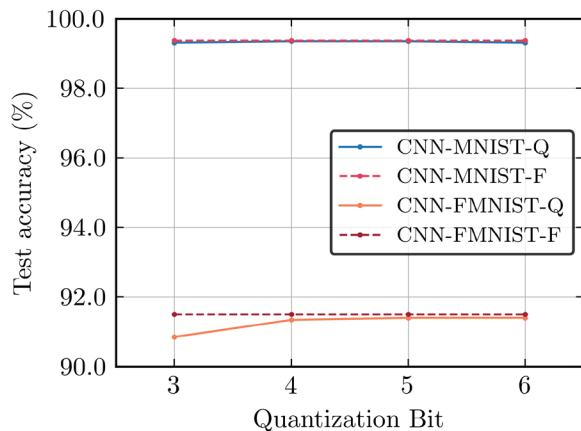
# Experimental Results: Quantization

## ◆ Experiments settings

- › Photonic PCM memory: 3 – 6-bit
- › Photonic tensor core:  $16 \times 16$  and  $64 \times 64$
- › Models: Simple CNN, VGG8 and VGG13
- › Dataset: MNIST, FashionMNIST, CIFAR10 and CIFAR100

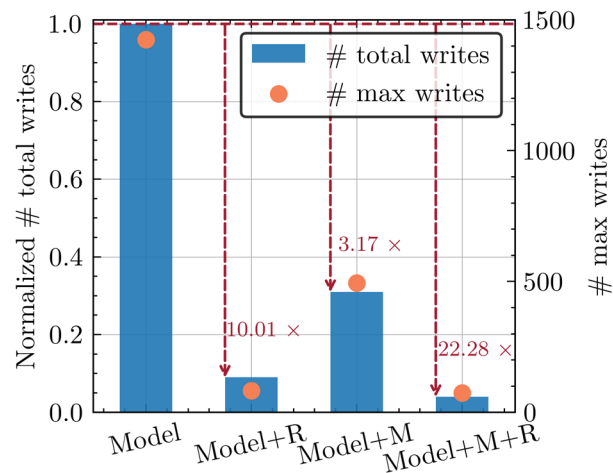
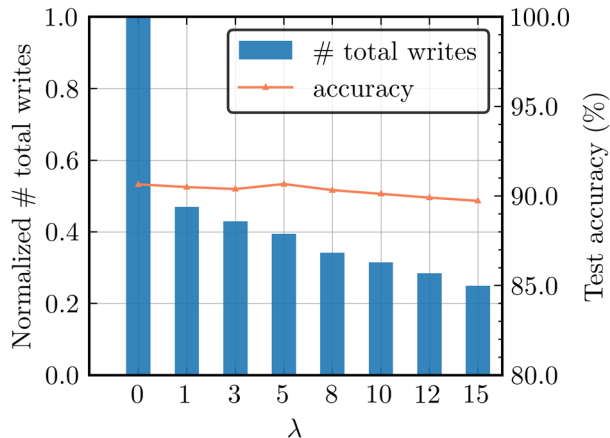
## ◆ Distribution-aware quantization

- › Small accuracy loss with  $> 4$ -bit



# Experimental Results: # total Write

- ◆ 5-bit VGG8 on CIFAR10
- ◆ Write-aware training (M)
  - › Trade-off between accuracy and write elimination
  - › **< 1%** accuracy drop with sweet parameters
- ◆ Post-training reordering (R)
  - › Further cut down redundant writes
  - › Achieve **22.3 ×** reduction on #writes with M



# Experimental Results: Endurance & Energy

- ◆  $> 20\times$  fewer write operations
- ◆ **Minimized** # max writes on a single PCM cell
- ◆  $> 30\times$  less dynamic energy cost

ONN Lifetime  $\uparrow$   
Dynamic energy  $\downarrow$

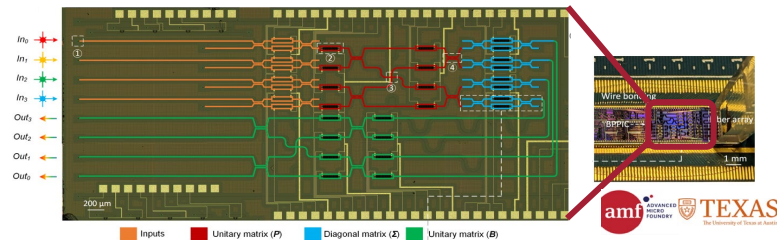
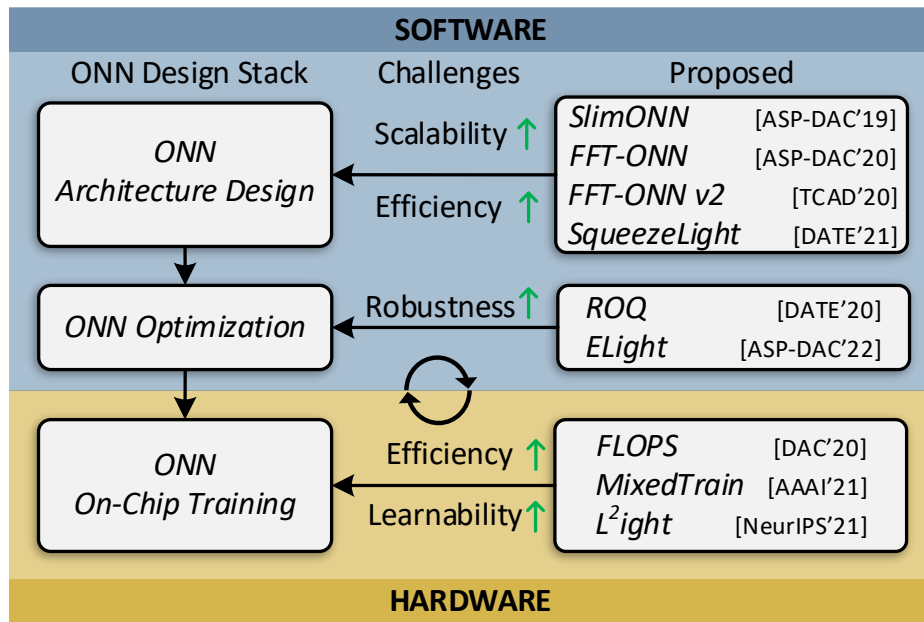
Network	Dataset	Bitwidth	$\lambda$	Acc(%) / AC	# total writes $\downarrow$ ( $\times$ )		Energy cost $\downarrow$ ( $\times$ )		# max writes	
					-	+R	-	+R	-	+R
VGG8	CIFAR-10	3	0	86.71	1	6.52	1	9.27	128	15
			8	86.02/-0.69	22.12	<b>46.11</b>	6.63	<b>69.29</b>	14	<b>7</b>
		4	0	89.75	1	7.84	1	11.31	401	36
			10	89.94/+0.19	3.83	<b>24.45</b>	3.92	<b>35.48</b>	95	<b>19</b>
		5	0	90.56	1	10.01	1	14.35	1425	82
			10	90.12/-0.44	3.17	<b>22.28</b>	3.20	<b>31.17</b>	494	<b>74</b>
		6	0	90.83	1	12.31	1	16.89	4464	180
			5	89.88/-0.95	6.82	<b>26.35</b>	7.15	<b>32.48</b>	1560	<b>146</b>
VGG13	CIFAR-100	4	0	70.99	1	9.66	1	13.84	542	39
			10	70.44/-0.55	3.54	<b>29.25</b>	3.57	<b>42.02</b>	173	<b>33</b>
		5	0	71.73	1	12.06	1	17.29	1771	84
			3	71.95/+0.22	2.19	<b>21.93</b>	2.21	<b>31.41</b>	921	<b>55</b>
		6	0	71.88	1	14.37	1	17.62	4926	182
			3	70.97/-0.91	3.11	<b>22.65</b>	3.19	<b>29.85</b>	3577	<b>156</b>

# Conclusion and Future work

- ◆ The **first** aging-aware optimization framework for Photonic in-memory computing
- ◆ Expressivity: Distribution-aware quantization
- ◆ Lifetime enhancement:  $> 20 \times$  fewer write operations
- ◆ Energy efficiency:  $> 30 \times$  less dynamic energy cost
- ◆ Push forward the real deployment of Photonic in-memory computing
  
- ◆ Future direction
  - › Preserve the accuracy of NN model with aged PCM cells
  - › Counter other non-ideal factors such as device-to-device variations
  - › Investigate the effect of temporal drift

# Our Recent Work & Open-source Framework

- ◆ Build ultra-fast (light-speed), ultra-energy efficient, and highly robust optical neural accelerators with photonic integrated circuits



**Circuit-Architecture-Algorithm Co-Design!**

**PyTorch-ONN Library**



# Thanks!

# Q & A?

