

MPSoC'25, June 16, 2025

Optical Streaming Processor Utilizing Integrated Photonic Circuits

Guangwei Cong

Photonics-Electronics Integration Research Center
National Institute of Advanced Industrial Science and Technology (AIST), Japan.
E-mail: gw-cong@aist.go.jp

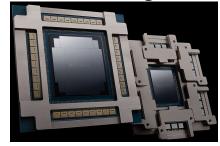
Outline

- Why optical computing is reignited?
- How to do optical computing in photonics?
- Our works: optical streaming processors implementing novel models in photonics

Why optical computing is reignited?

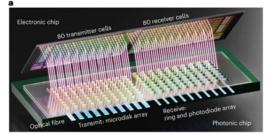
- (1) Power issue and demands for computational power
 - Computational power doubles per 10 months
 - Energy consumption is already substantial, and is increasing rapidly
 - End of Moore's law→ stagnation of digital processor performance
 - End of Dennard law → energy efficiency gains slows
- (2) Advances in optical interconnects and photonics-electronics co-integration
- (3) Physical system computes better in energy efficiency

On-Package



https://nvidianews.nvidia.com/news/nvidia-spectrum-x-co-packaged-optics-networking-switches-ai-factories

On-Chip



https://www.nature.com/articles/s41566-025-01633-0

Physical concepts + machine learning



"They used physics to find patterns in information"

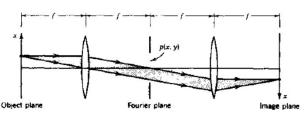
https://www.nobelprize.org/

Intrigue constructing optical physical analog AI system

Optical computing: Hype or hope?

Energy saved or not?

Example: Optical Fourier transformation



I'd to call it lens instead of computing since

Three elements required

- On-demand use input
- Reconfigurable
- Result interpretation
- → Overhead consumes additional powers

Pros and cons

Pros:

Low energy
Low latency
High clock speed
High bandwidth
High parallelism
Network compatibility

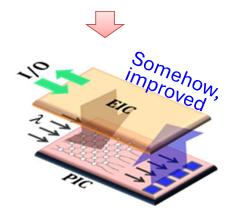
Cons:

Scalability
Programmability
Nonlinearity
Trainability
Stability
Low resolution
Area efficiency

Why it failed?

Lost to VLSI due to

- Scalability
- Low integration
- Programmability
- Not easy to use



What we are expecting

- One-shot
- All-optical
- High-throughput
- Applicable scale
- System-level energy efficiency
- Less digital intervene

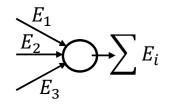
Comments:

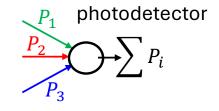
- System evaluation is of necessity
- Concerns of heavy overhead related to digital processors and memory
- It is time to re-evaluate the possibilities of optical computing

How to do computing in photonics?: scalar operations

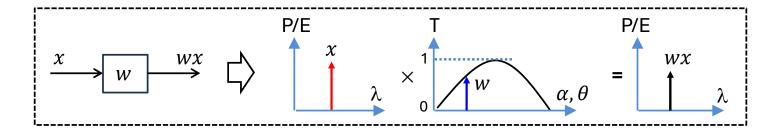
- Addition
 - Electrical field addition (coherent)
 - Optical power addition (incoherent)

interferometer

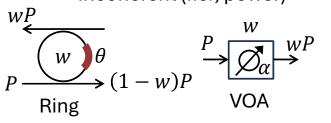




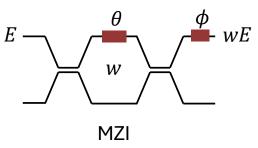
Scalar multiplication

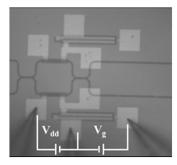


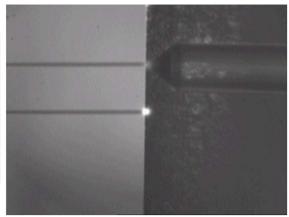
Incoherent (i.e., power)



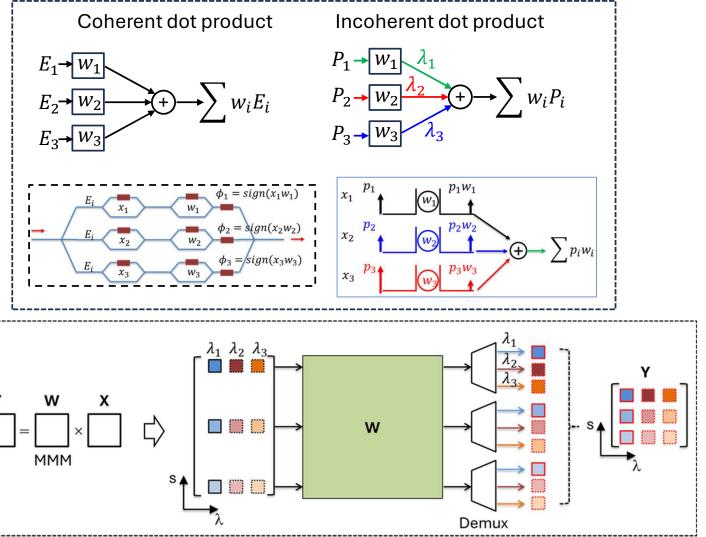
Coherent (i.e., electrical field)

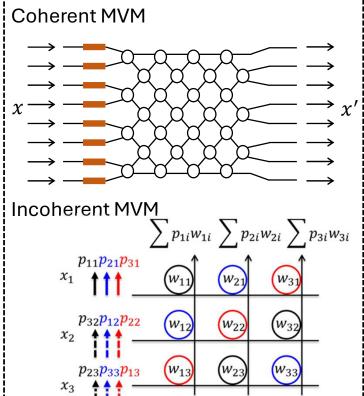






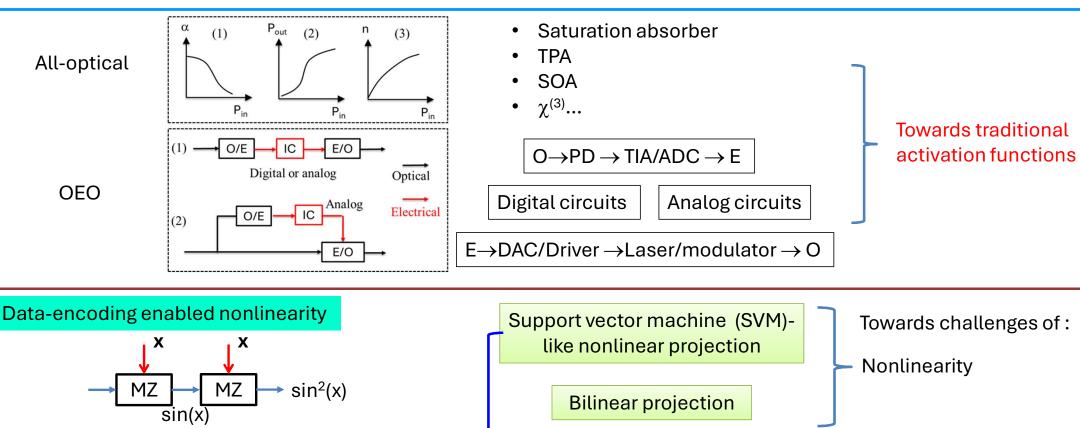
Multiplexing for linear operations





Massive parallelism by sharing hardware, offering scalable and energy efficient linear operations.

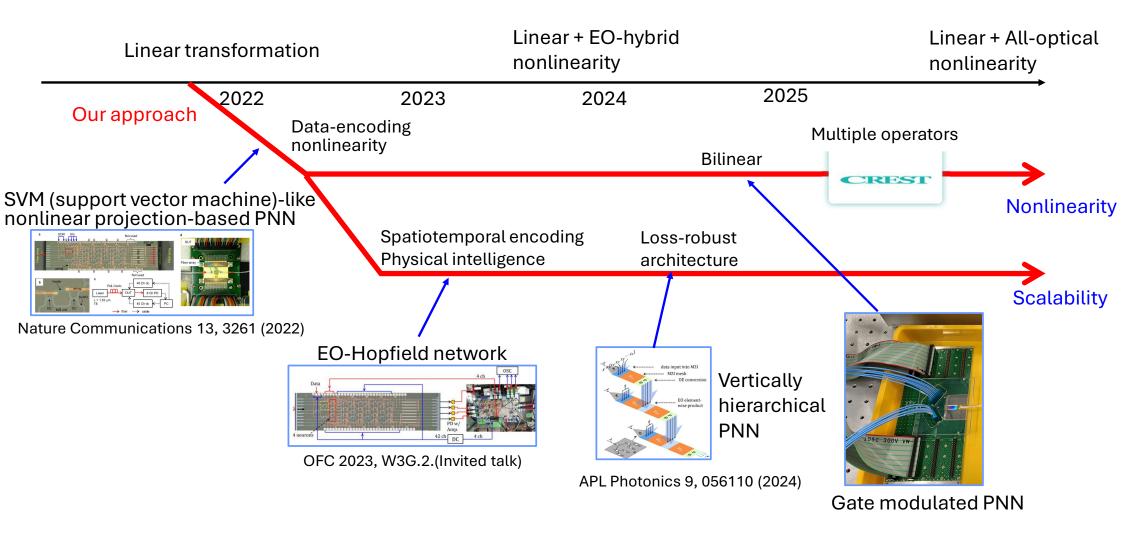
How nonlinear functions can be achieved in photonics?



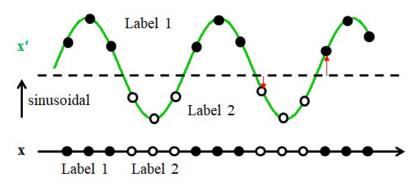
- We are focusing on wide nonlinearities to construct new models instead of replacing traditional ones
- Despite of same technology, new possibility and effects are achievable
- Offering new merits in solving some remaining $_{AIS}$ challenges

7

Our works: optical streaming processors implementing novel models



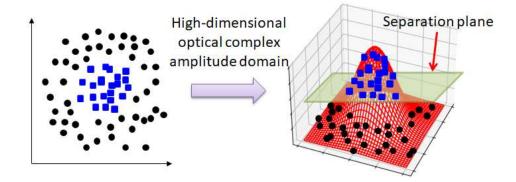
(1) Nonlinear projection-based PNN

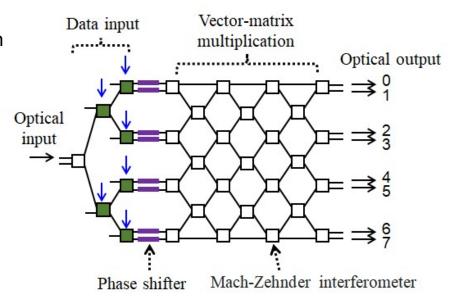


- For multilayer perceptron (MLP), it is using y = wf(...(wf(wx + b) + b)) to approximate a nonlinear function.
- We are thinking reversely: constructing nonlinear projection function firstly, that may enable easy separation.

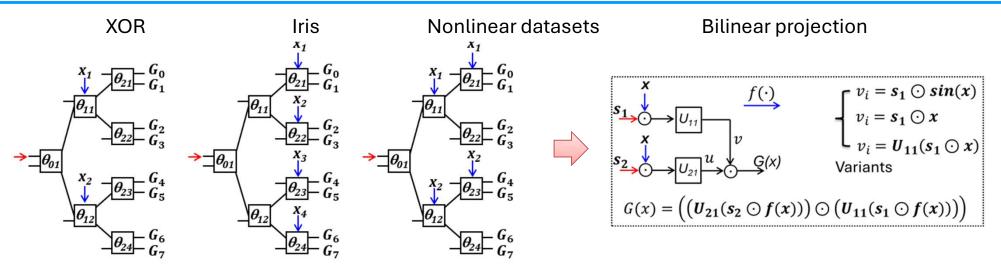
Nonlinear projection + linear separation

- Leveraging phase-amplitude (EO) nonlinearity of MZI
- Linear separation by VMM afterwards
- Maximum-power port position indicates the result
- Only passive circuits

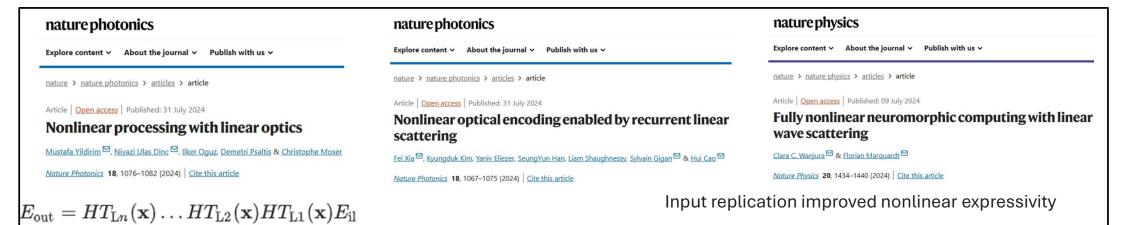




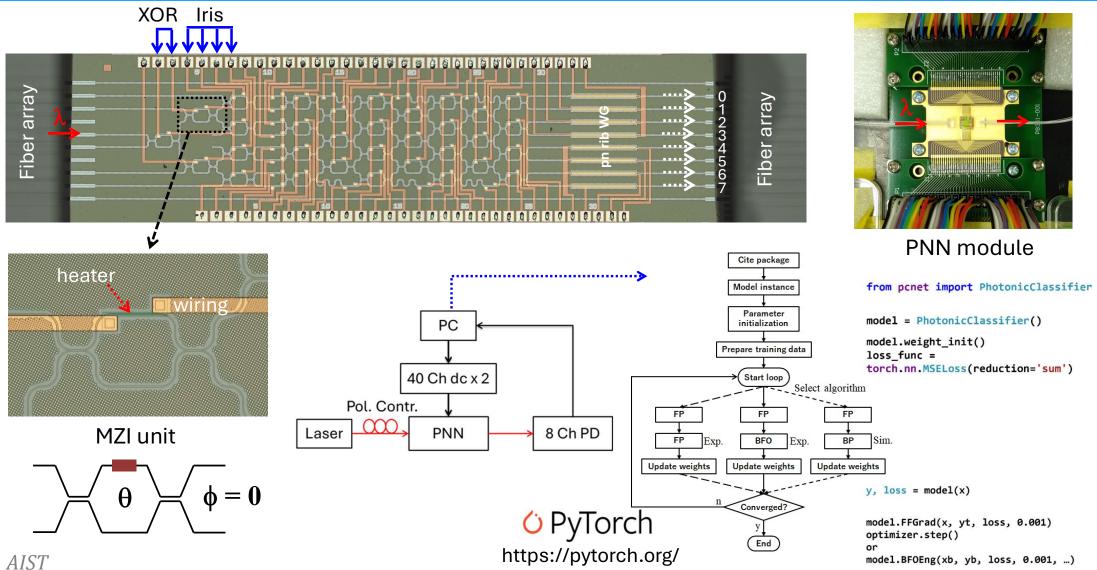
Encoding the data into MZI network



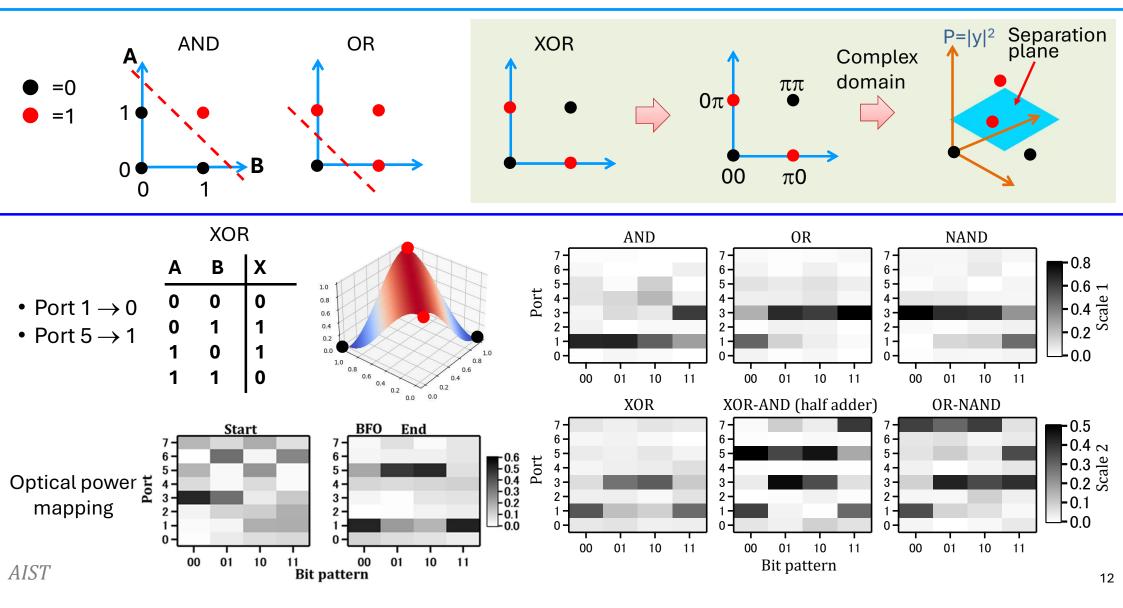
- Encoding the data into MZI networks
- · Leveraging EO nonlinearity associated with data encoding



 ΔIST



Classification experiment: Boolean logic

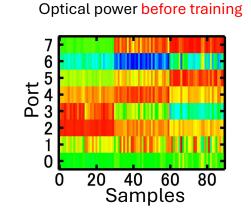


Classification experiment: Iris dataset

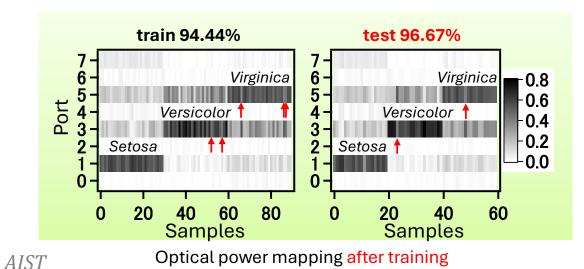
A DE

https://archive.ics.uci.edu/ml/datasets/iris

	Species	Petal width	Petal length	Sepal width	Sepal length	No.
THE S	Setosa -	0.2	1.4	3.5	5.1	1
S 4. A	Setosa	0.2	1.4	3	4.9	2
Setos	Setosa	0.2	1.3	3.2	4.7	3
						-
	Versicolor	1.4	4.7	3.2	7	51
	Versicolor	1.5	4.5	3.2	6.4	52
	Versicolor	1.5	4.9	3.1	6.9	53
Versico	-			3 1		:
	Virginica	2.5	6	3.3	6.3	101
	Virginica	1.9	5.1	2.7	5.8	102
1	Virginica	2.1	5.9	3	7.1	103
			(2)	[8]	8	9

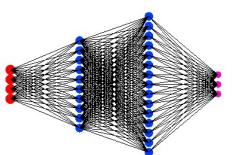


- Data normalization
- Port assignments
 - Port 1 → Setosa
 - Port 3 → Versicolor
 - Port 5 → Virginica



Two examples of digital NN

https://python.atelierkobato.com/variety/



- $4\times10\times15\times3$
- ReLU
- 235 weights
- ~97%

Conferences > 2018 International Conference...

A Model of Deep Neural Network for Iris Classification With Different Activation Functions

Publisher: IEEE

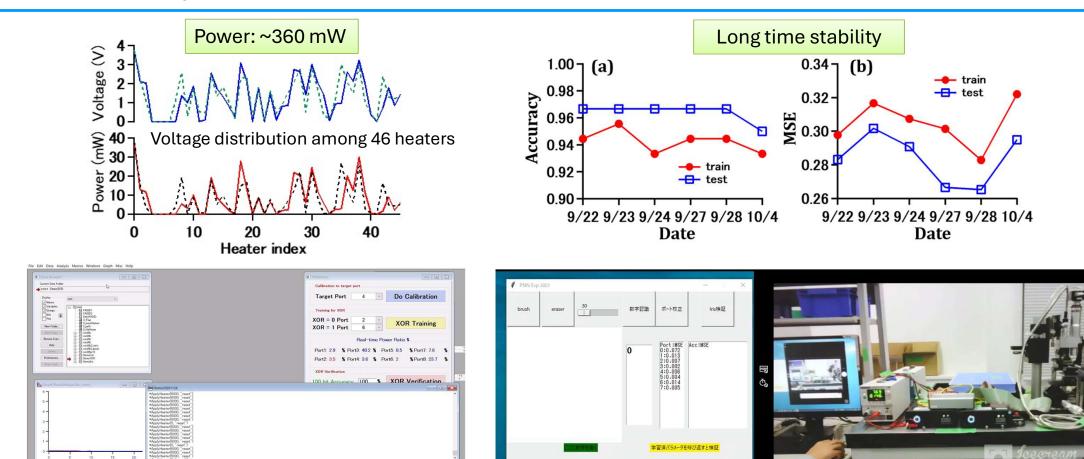
Cite This

A PDF

DOI: 10.1109/IDAP.2018.8620866

1st Hidden Layer	2 nd Hidden Layer	Epoch Number	Accuracy Rate	
	,	100	86%	
D.I.	G''1	200	86%	
Relu	Sigmoid	300	86%	
		400	86%	
		100	90%	
n. 1	m 1	200	90%	86-90%
Relu	Tanh	300	90%	
		400	90%	
		100	90%	
	n	200	90%	
Relu	Relu	300	90%	
		400	90%	

Details of experimental results



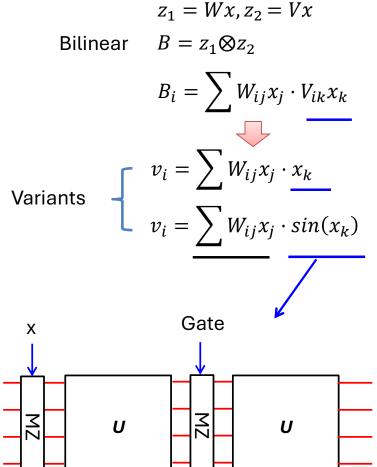
XOR video using BFO algorithm

✓ Electronics for IO only

- ✓ Computing is done just by optical propagation
- ✓ Latency < 100 ps

Demo video

(2) Bilinear projection in gate modulated PNN



arXiv:2002.05202v1,12 Feb 2020 GLU Variants Improve Transformer

Noam Shazeer Google noam@google.com

February 14, 2020

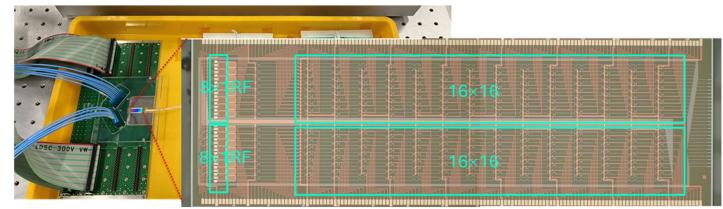
These variants have shown improvements over traditional activation functions when used in Transformer models.

Abstract

Gated Linear Units [Dauphin et al., 2016] consist of the component-wise product of two linear projections, one of which is first passed through a sigmoid function. Variations on GLU are possible, using different nonlinear (or even linear) functions in place of sigmoid. We test these variants in the feedforward sublayers of the Transformer [Vaswani et al., 2017] sequence-to-sequence model, and find that some of them yield quality improvements over the typically-used ReLU or GELU activations.

$$\mathrm{GLU}(x, W, V, b, c) = \sigma(xW + b) \otimes (xV + c)$$

$$Bilinear(x, W, V, b, c) = (xW + b) \otimes (xV + c)$$



Acknowledgments







JST CREST (JPMJCR24R3) (2024 newly launched)





AIST SCR Platform







Thank you for kind attention!