# Photonics for AI

Dr. Yichen Shen
Presentation @ MIT
Apr 26th, 2018

# Nanophotonics

Optical structures (dielectrics) with nanometer-scale features

wavelength of visible light
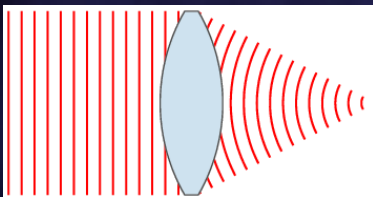


Feature size

m          mm          μm          nm

Macroscopic structures
(Feature size >> Wavelength)

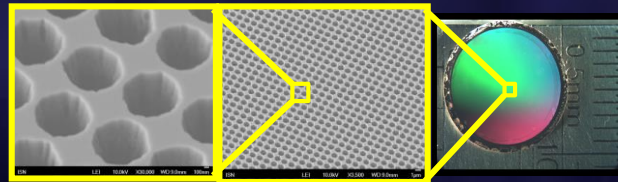Nanophotonic structures
(Feature size <≈ Wavelength)

Metamaterials

Photonic crystals

Ray optics

500 nm

Integrated Photonics

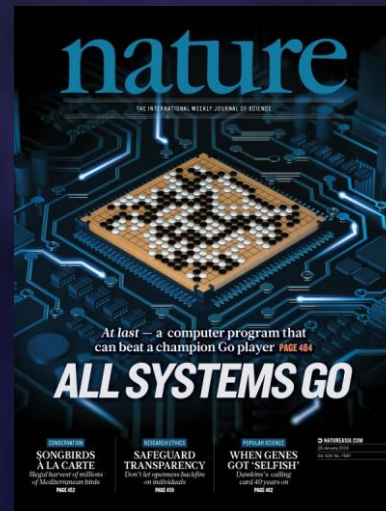Limited ways to
manipulate light

Thanks to improved computation power & fabrication techniques

2

# Artificial Neural Networks (ANN)
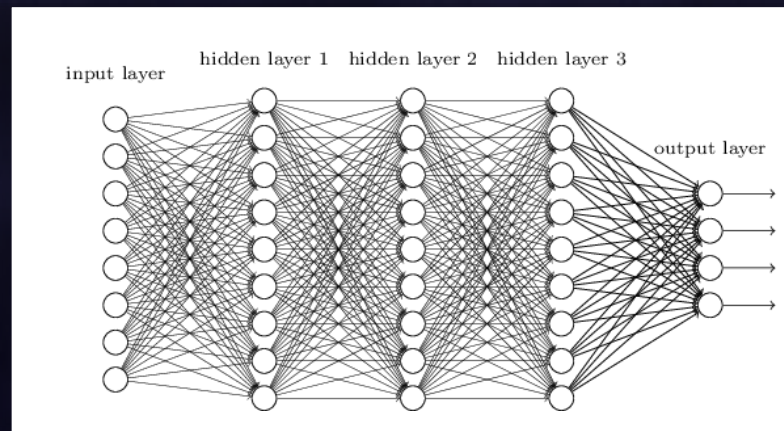
## Breakthroughs in deep learning:

- Natural Language Processing (NLP)
- Game Playing (Go, Atari)
- Autonomous Vehicles
- Control
- Ad Placement
- Researches (drug discovery, material study)
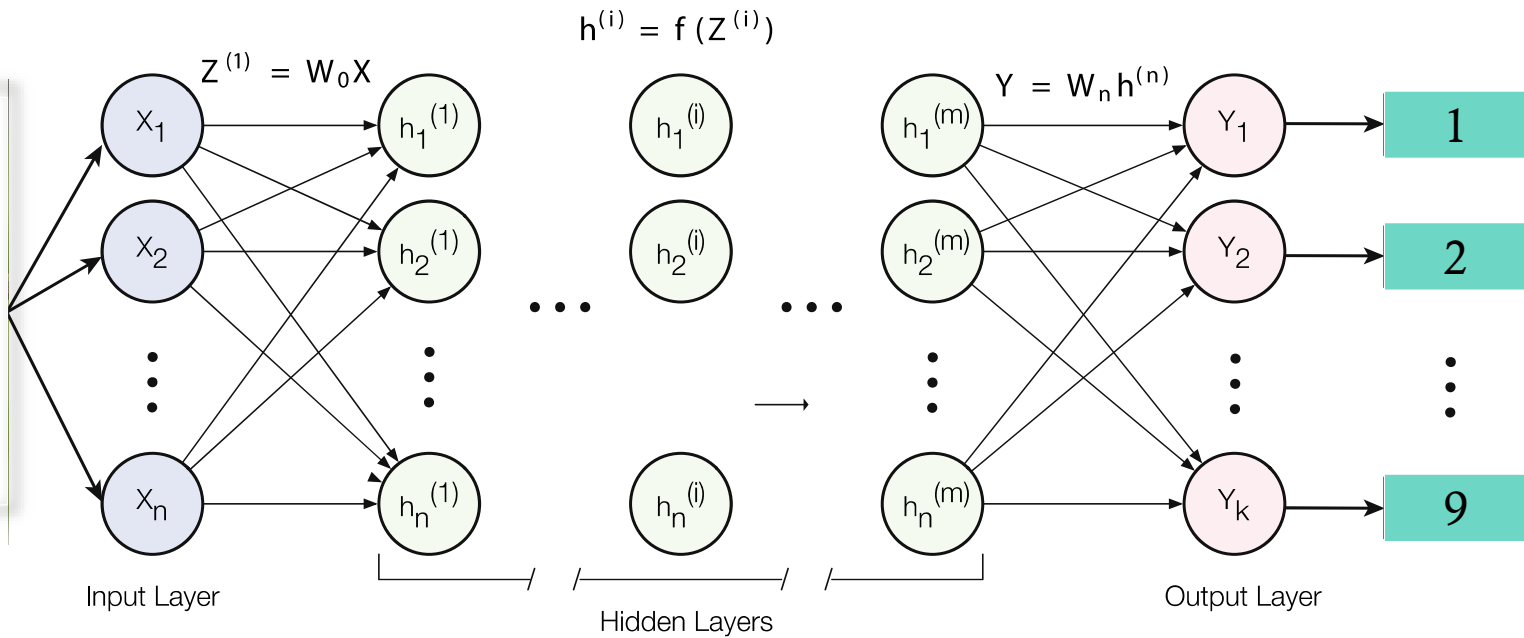- Etc.

# Neuromorphic Computing



Biological Neural Networks
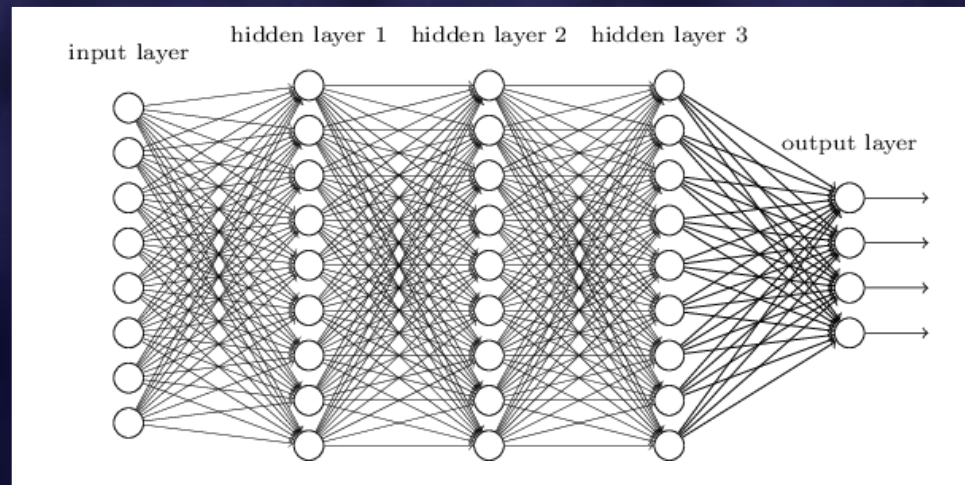


Artificial Neural Networks

# Basic Algorithm of ANN



$$z_j^{(1)} = \sum_{i=1}^{n} w_{ji}^{(1)} x_i \qquad z_j^{(k)} = \sum_{i=1}^{n} w_{ji}^{(k)} h_i^{(k)}$$

Matrix Multiplication:

Nonlinear Activation:

$$h_j^{(1)} = f(z_j^{(1)}) \qquad h_j^{(k)} = f(z_j^{(k)})$$

# Hardware and Data Enable Deep Learning



input layer    hidden layer 1    hidden layer 2    hidden layer 3

output layer

# The Need for Speed

**More Data → Bigger Models → More need for Computation**

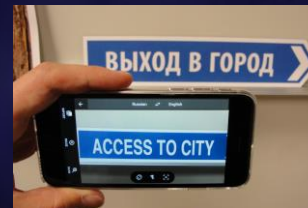**But Moore's Law is no-longer providing more computation…**



The Market:



On clouds:
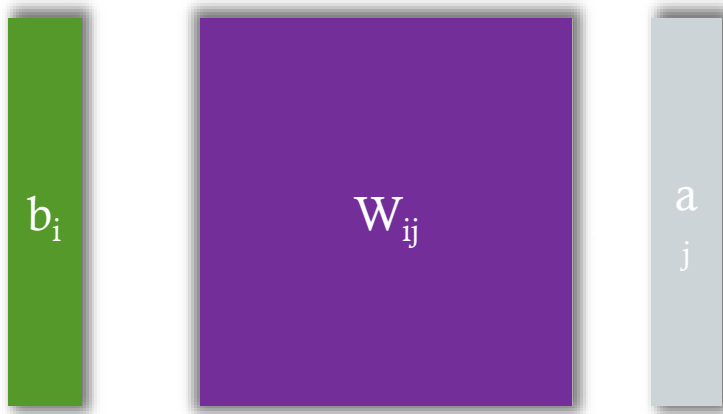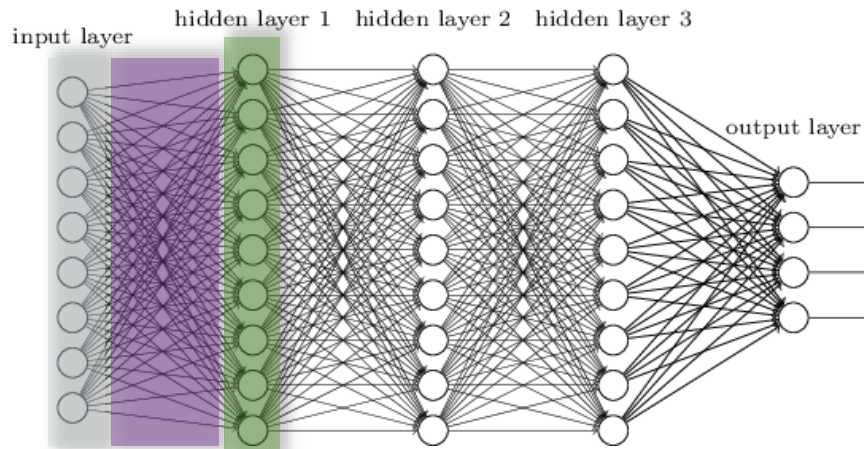Millions of high power AI processors ($10,000 each) in data centers by 2020





On premise:
Billions of compact AI processors needed due to the rise of autonomouse driving, AR and IoT.
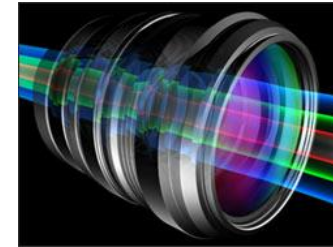
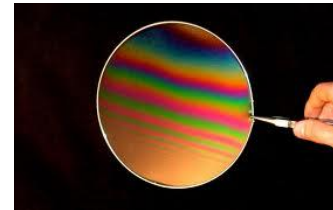Deep Learning with Coherent Nanophotonic Circuits

4/26/2018

# In Deep Learning
# Key Operation is dense M x V
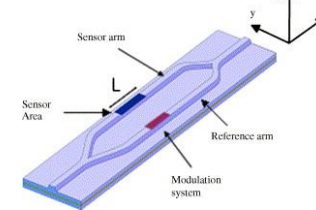


In Optics, Matrix Multiplication is very common & (usually) consumes no energy !



Convolution / FFT



Matrix Multiplication

$b_i$

$W_{ij}$

$a_j$

# ANN does NOT require high resolution

| Category | Method | Weights (# of bits) | Activations (# of bits) | Accuracy Loss vs. 32-bit float (%) |
|---|---|---|---|---|
| Dynamic Fixed Point | w/o fine-tuning | 8 | 10 | 0.4 |
| | w/ fine-tuning | 8 | 8 | 0.6 |
| Reduce weight | Ternary weights Networks (TWN) | 2* | 32 | 3.7 |
| | Trained Ternary Quantization (TTQ) | 2* | 32 | 0.6 |
| | Binary Connect (BC) | 1 | 32 | 19.2 |
| | Binary Weight Net (BWN) | 1* | 32 | 0.8 |
| Reduce weight and activation | Binarized Neural Net (BNN) | 1 | 1 | 29.8 |
| | XNOR-Net | 1* | 1 | 11 |
| Non-Linear | LogNet | 5(conv), 4(fc) | 4 | 3.2 |
| | Weight Sharing | 8(conv), 4(fc) | 16 | 0 |
| * first and last layers are 32-bit float | | | | |

**Sze et al, arXiv:1703.09039 (2017)**
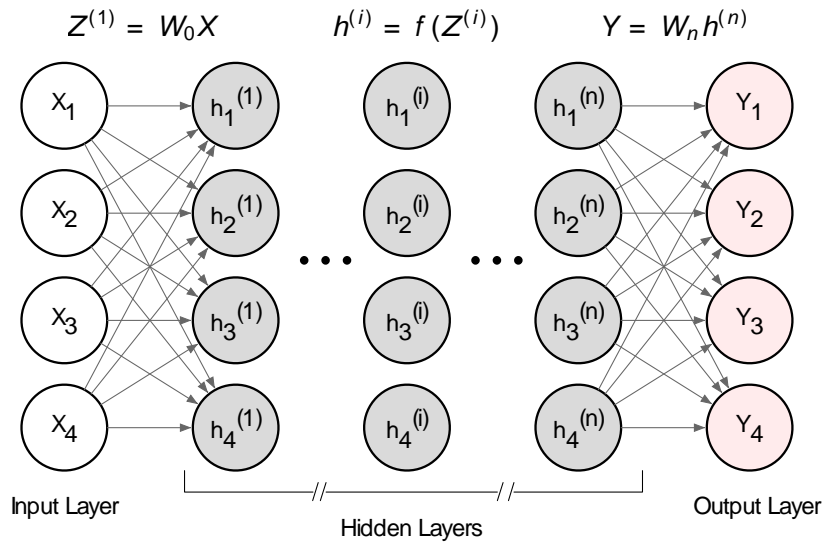
# Deep Learning Inference is "Passive"

Once the Optical Neural Network is trained, no need to update the weights frequently…
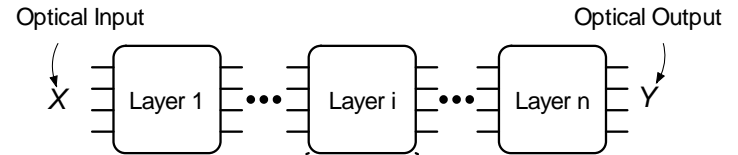
# Deep Learning is very parallelizable

Multiple wavelengths can be used to simultaneously execute batch of data
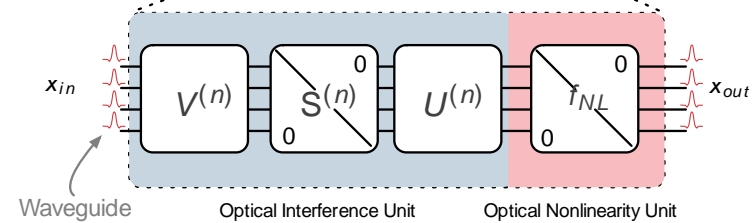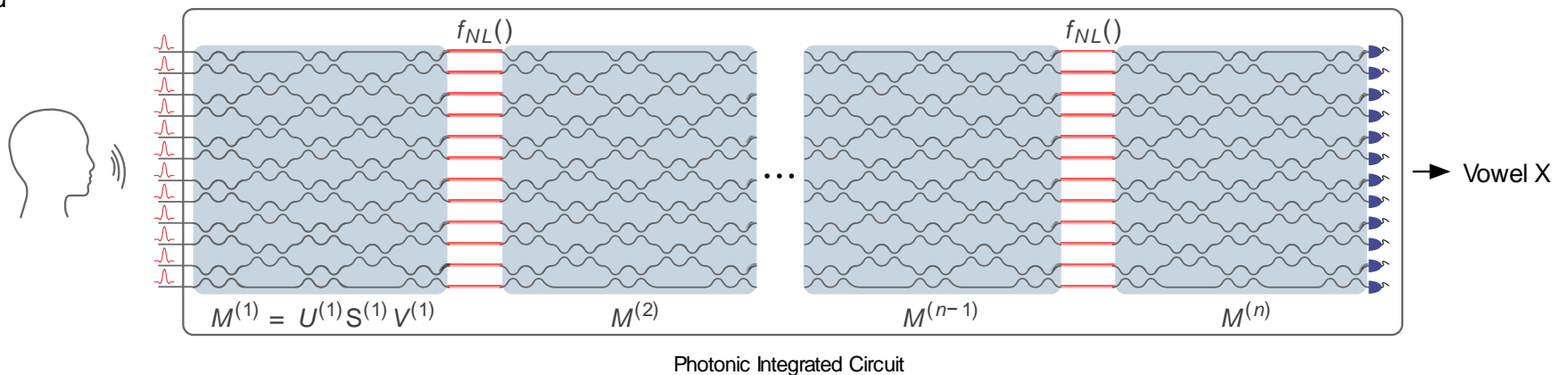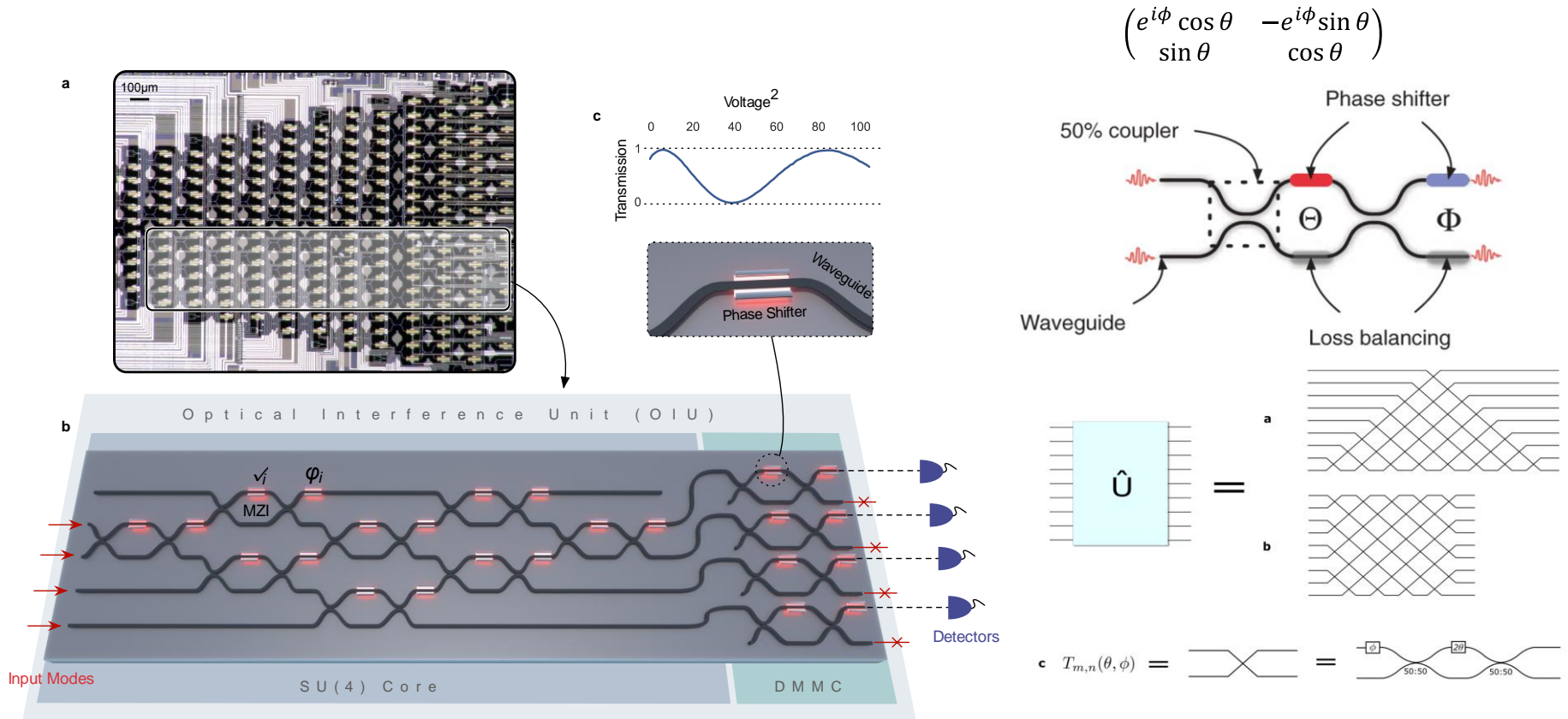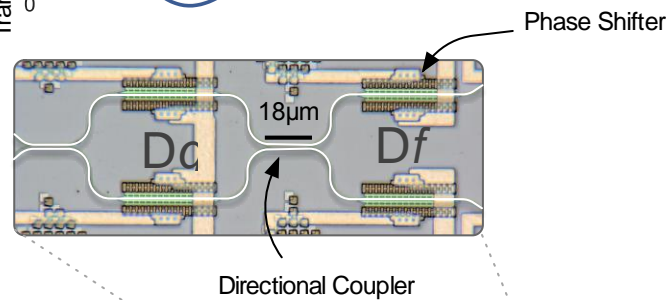
# Coherent Optical Neural Networks (ONN)



a

$Z^{(1)} = W_0 X$     $h^{(i)} = f(Z^{(i)})$     $Y = W_n h^{(n)}$

Input Layer

Hidden Layers

Output Layer

b

Optical Input               Optical Output

$X$   Layer 1   $\cdots$   Layer i   $\cdots$   Layer n   $Y$

c

$x_{in}$

$V^{(n)}$   $S^{(n)}$   $U^{(n)}$   $f_{NL}$   $x_{out}$

Waveguide       Optical Interference Unit       Optical Nonlinearity Unit

d

$f_{NL}()$                   $f_{NL}()$

$\longrightarrow$ Vowel X

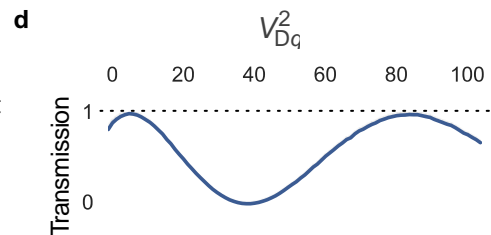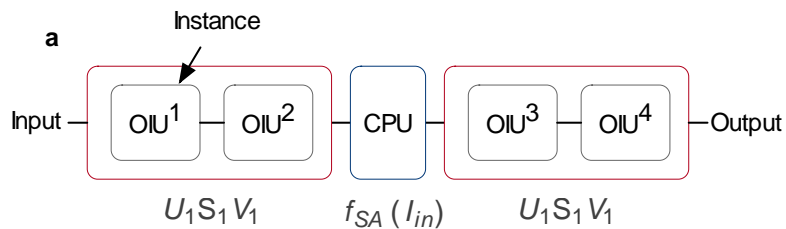$M^{(1)} = U^{(1)} S^{(1)} V^{(1)}$     $M^{(2)}$     $M^{(n-1)}$     $M^{(n)}$
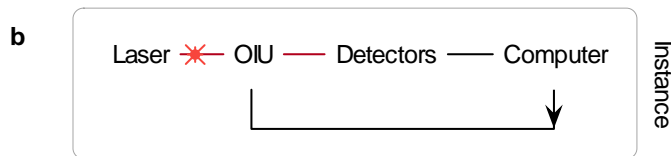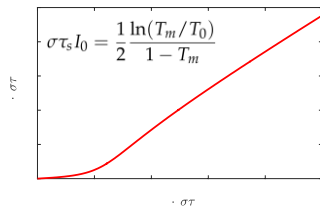
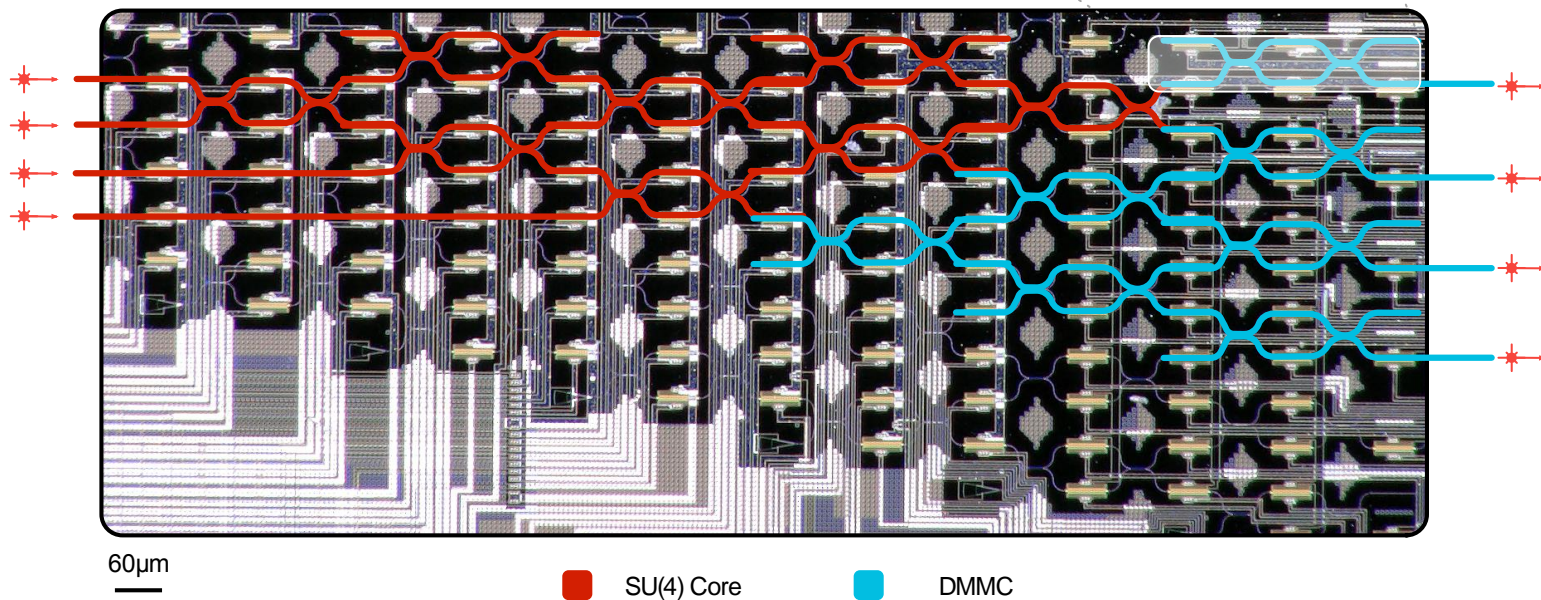Photonic Integrated Circuit

# Programmable Nanophotonic Processors



Y.Shen and N. Harris et al. *"Deep Learning with Coherent Nanophotonic Circuit" Nature Photonics* **11**, 441–446 (2017)
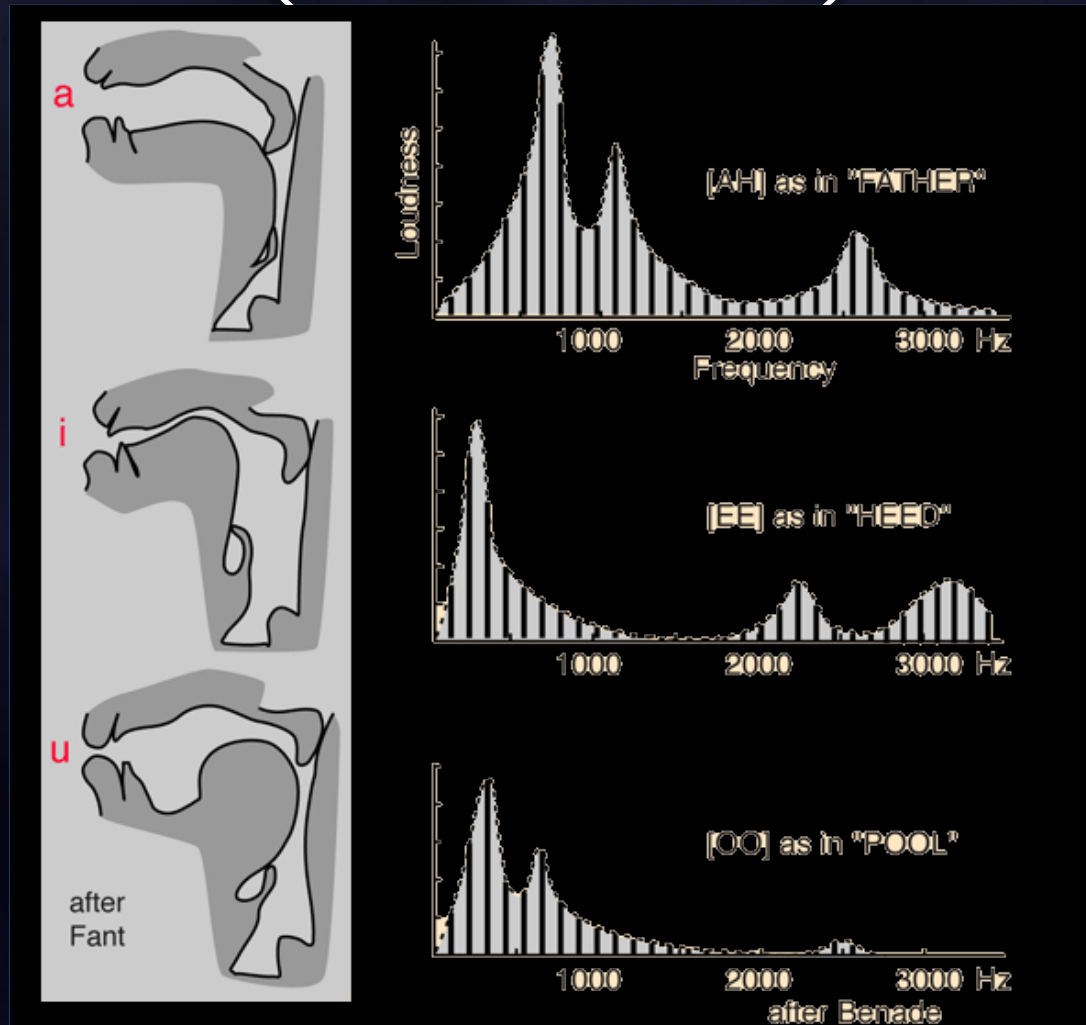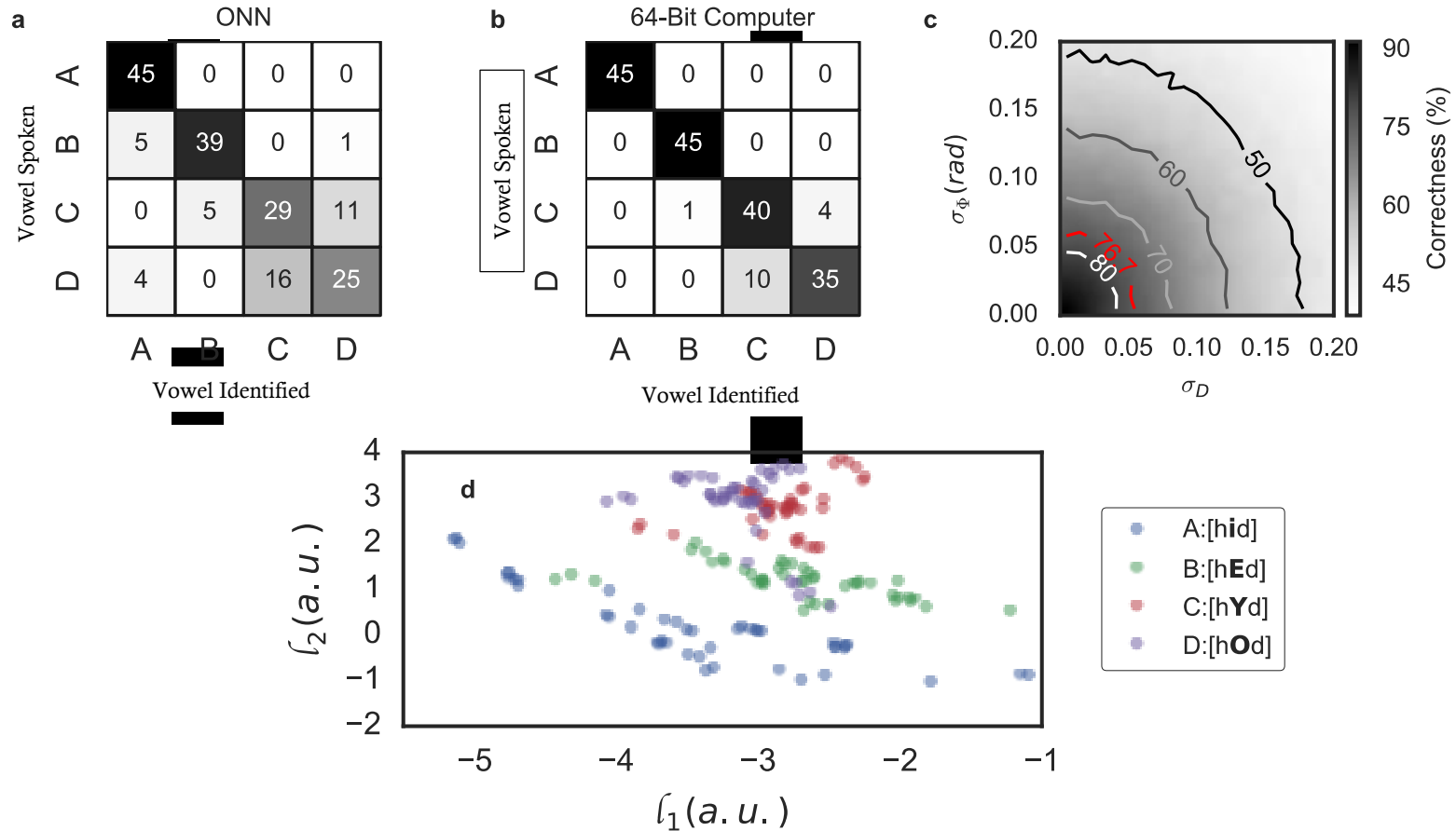
# Simulated Optical Nonlinearity

**a**



Instance

Input — OIU¹ — OIU² — CPU — OIU³ — OIU⁴ — Output

$U_1 S_1 V_1$ ........ $f_{SA}(I_{in})$ ........ $U_1 S_1 V_1$

$$\sigma \tau_s I_0 = \frac{1}{2} \frac{\ln(T_m/T_0)}{1 - T_m}$$

**b**

Laser ✳ OIU — Detectors — Computer

Instance

**d**

$V_{Dq}^2$

0   20   40   60   80   100

Transmission

Phase Shifter

18µm

Dq                    Df

Directional Coupler

**c**

Optical Interference Unit



60µm

■ SU(4) Core      ■ DMMC

Y.Shen and N. Harris et al. *"Deep Learning with Coherent Nanophotonic Circuit" Nature Photonics* **11**, 441–446 (2017)

# Optical Vowel Recognition
## (4d 4 classes)

# Experimental Result



**a** ONN

**b** 64-Bit Computer

**c**

**d**

Simulation Result: 165/180=91.7%
Experiment Result: 138/180=76.7%

# The other side of the Story…

- Immature photonics eco-system (low yield, high cost)

- Large device size

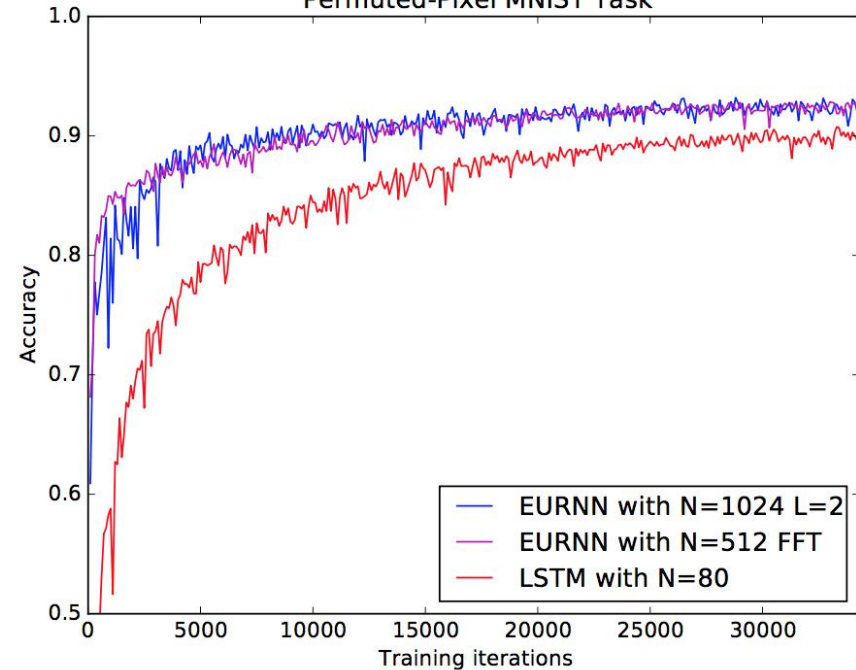- Non-ideal PDK component design (lossy, low resolution, power hungry)

- AD/DA interface

# Software & Hardware

AI algorithms DESIGNED to be run on photonics chip



L. Jing & Y. Shen et al, International Conference for Machine Learning (ICML 2017)
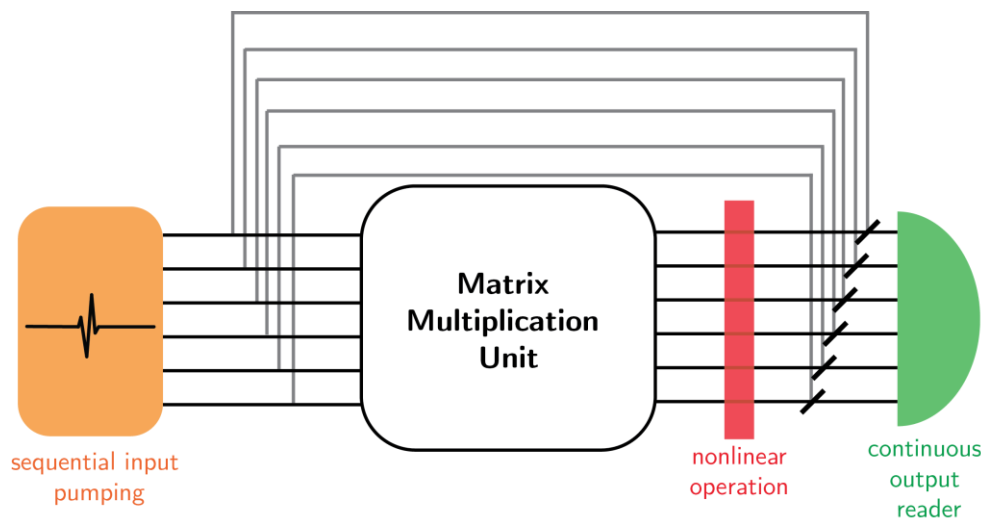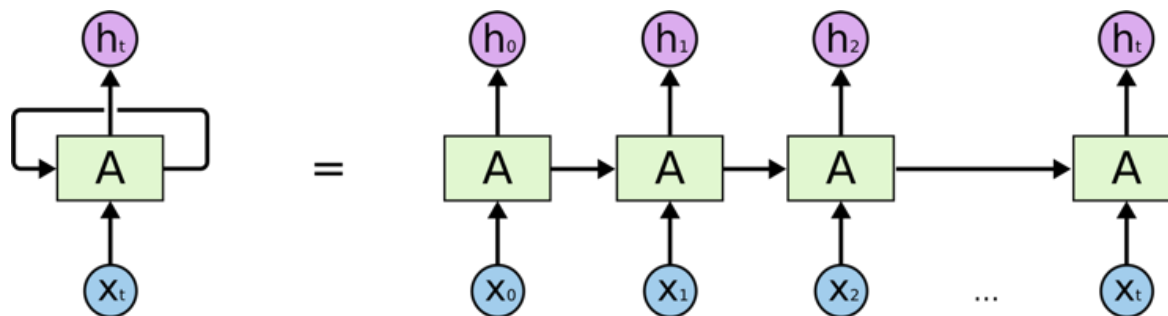
# Fully Connected Neural Networks



Recurrent Neural Networks
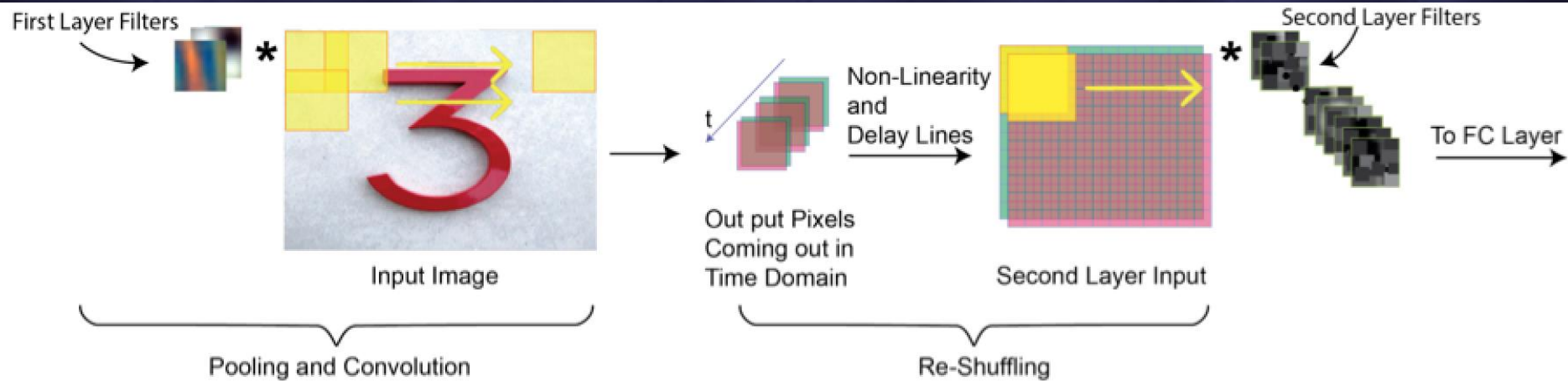
Convolutional Neural Networks

# Recurrent Neural Networks

Commonly used for Speech Recognition and Language Processing
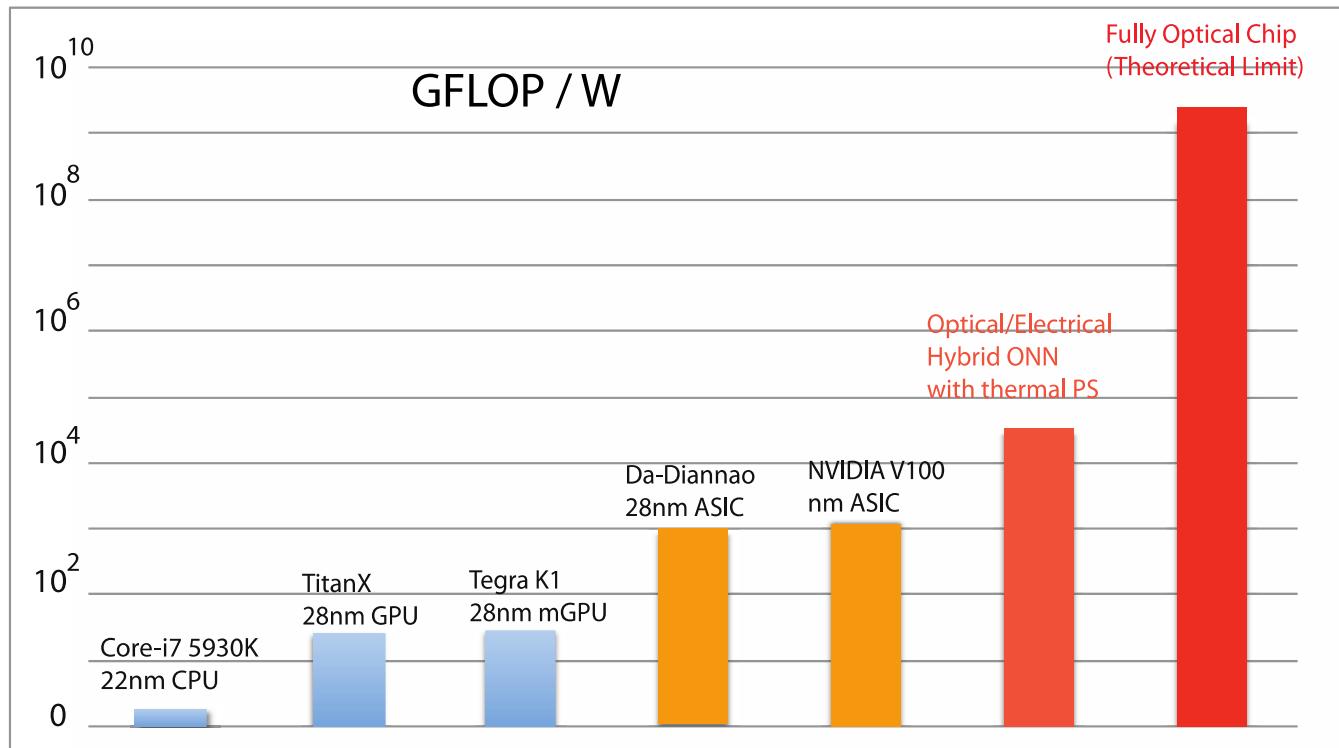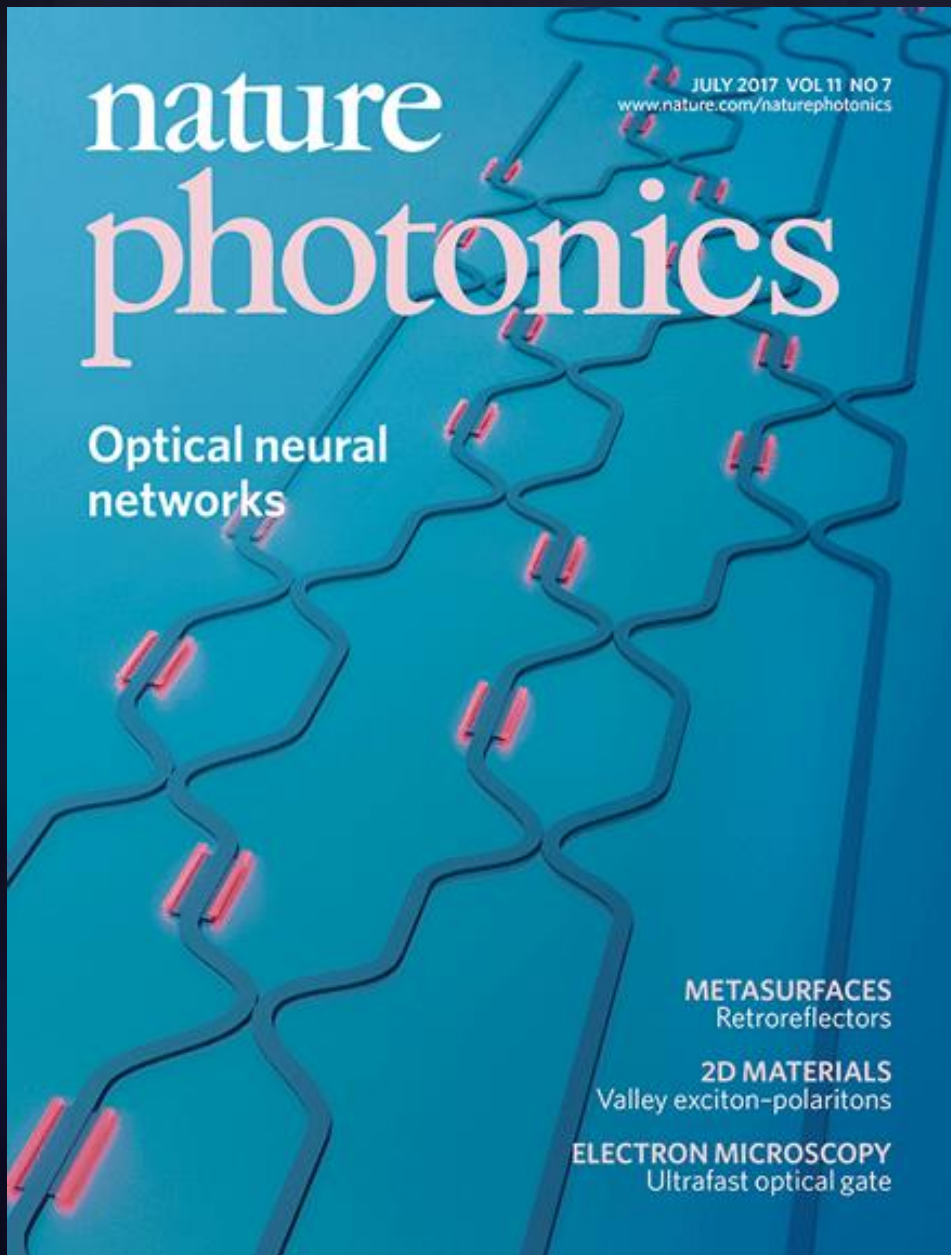
# Convolution Neural Networks



Scott Skirlo and Yichen Shen et al, Manuscript in Preparation

# Speed and Energy Efficiency Comparison with Electrical ANN

|  | NVIDIA TITAN X | ONN (with thermal PS) |
| --- | --- | --- |
| Architecture | Von Neumann | Neuromorphic |
| Power Consumption | 1 kW | 1-2 kW |
| Operation Speed | 10 TFLOP | 10,000 TFLOP |



Y. Shen and N. Harris et al, Nature Photonics 11, 441 (2017)

# Optical AI Computing

## LIGHT MAY BE KEY TO NEW GENERATION OF FAST COMPUTERS

By WILLIAM J. BROAD
Published: October 22, 1985

SINCE its start nearly half a century ago, the computer revolution has advanced by virtue of a simple physical phenomenon: that streams of speeding electrons can start or stop the flow of other streams of electrons. In short, electrons can act as a switch.

FACEBOOK

TWITTER

GOOGLE+

EMAIL

# Some History on Optical Neural Networks

**2005**

> "*The biggest issue with this paper is that it relies on neural networks.*"
>
> Anonymous Reviewer

**2016**

## Springtime for AI: The Rise of Deep Learning

After decades of disappointment, artificial intelligence is finally catching up to its early promise, thanks to a powerful technique called deep learning

———

By Yoshua Bengio on June 1, 2016

**SCIENTIFIC AMERICAN.**