# Ultrafast photonic convolution processing

**Wolfram Pernice**
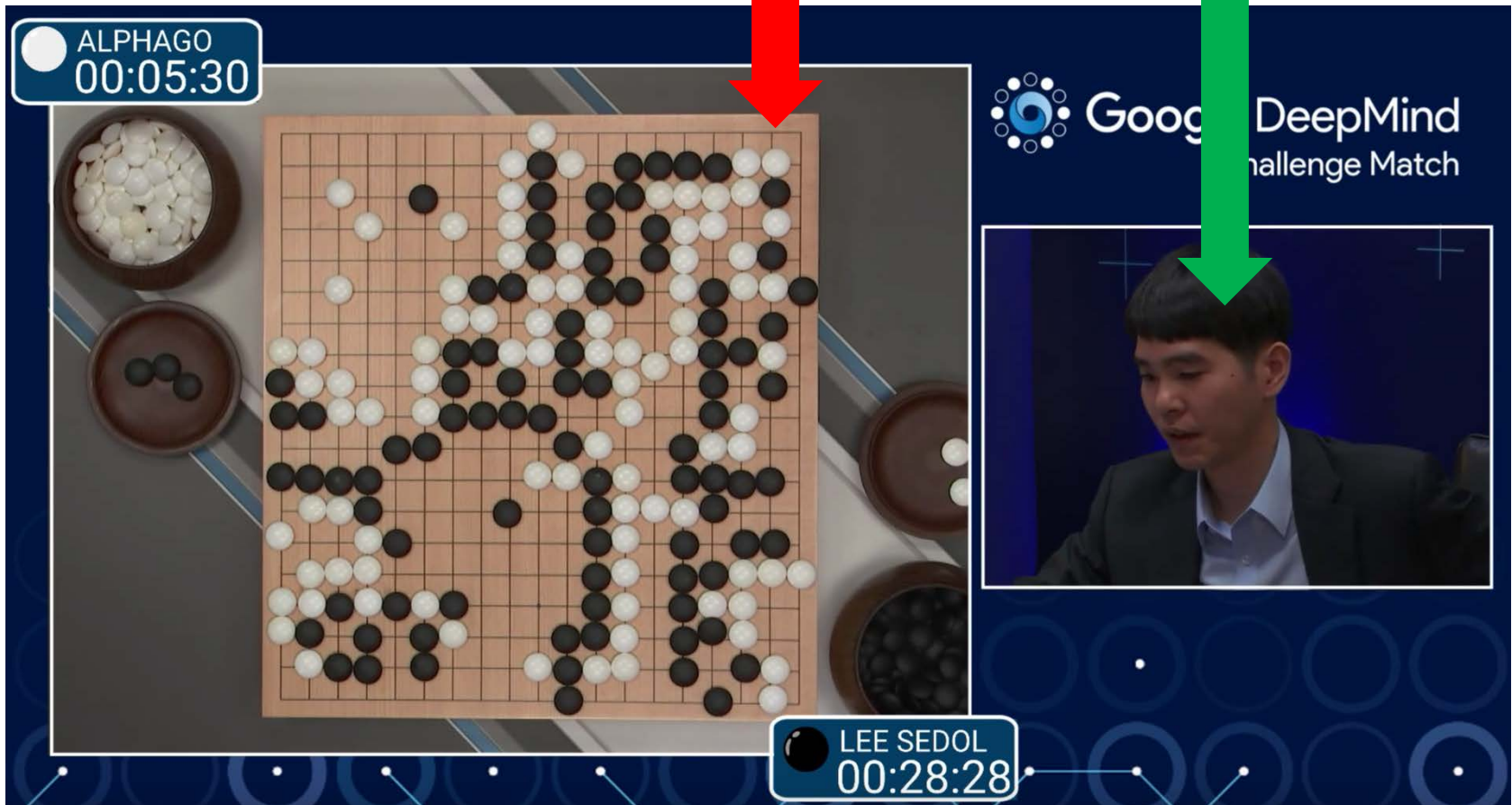http://www.uni-muenster.de/Physik.PI/Pernice/

Universität Münster (WWU), Physikalisches Institut



**Workshop on Optical Computing: current / emerging approaches & applications**
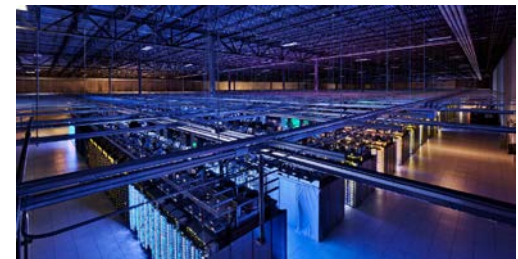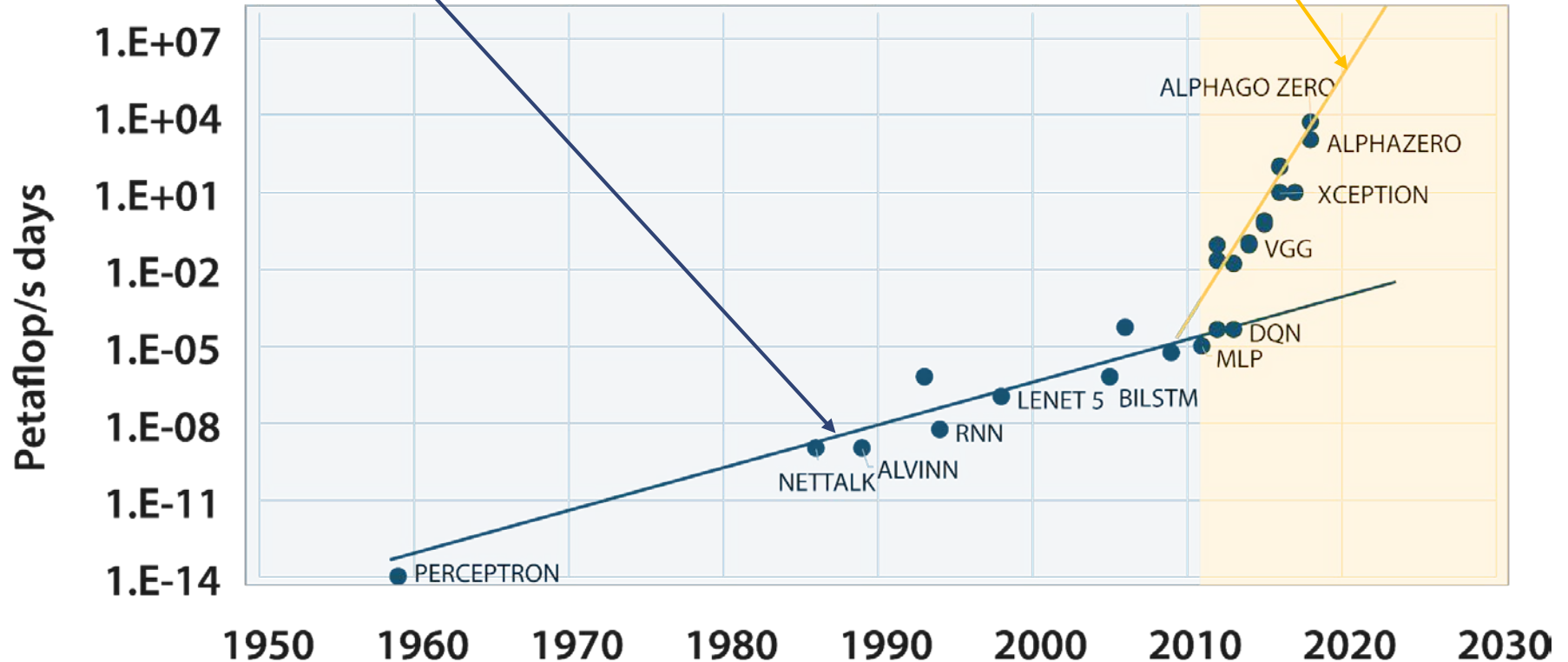
# AlphaGo (2016)

~1,000,000 W

~20 W



- 1202 Central Processing Units (CPUs)
- 176 Graphics Processing Units (GPUs)

# Moore's law – revisited



**Before: Processing power doubles every 2 years**
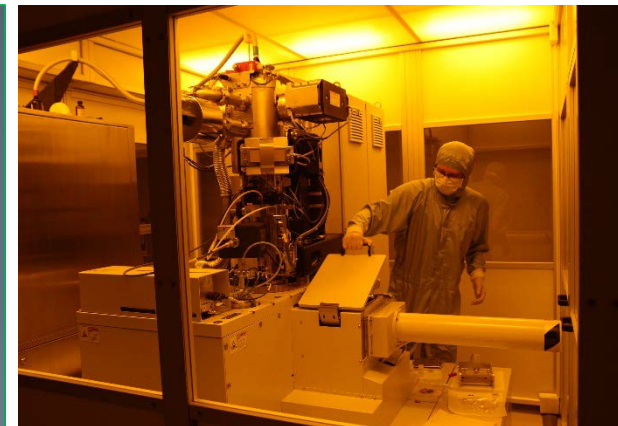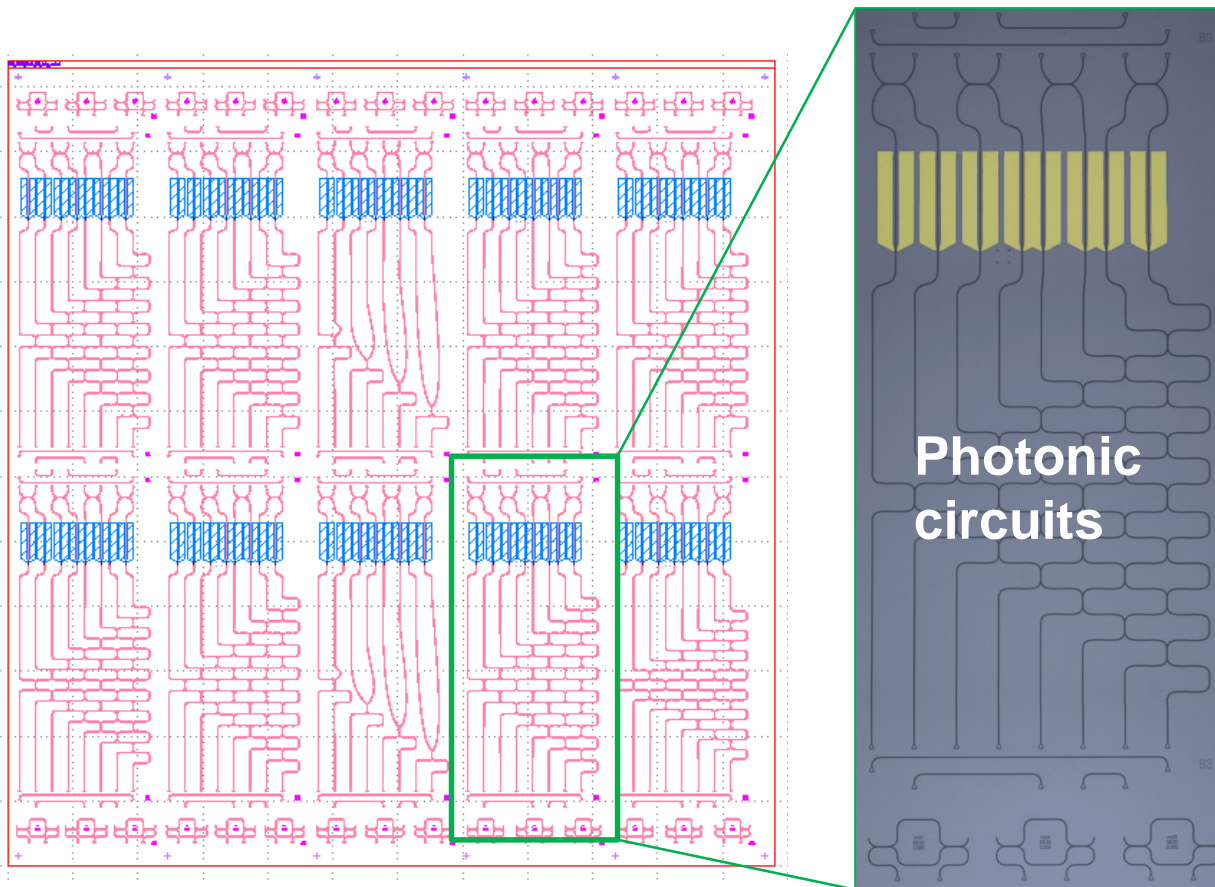
**Now: Processing power doubles every 3.5 months**

# Nanophotonic circuits @ WWU

- Integrated photonic components

- Multiple layers of lithography using alignment
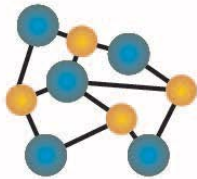
- Photonic CAD with Python framework



Photonic circuits

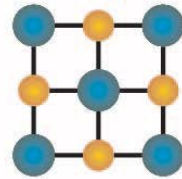*Gehring, et al., OSA Contin. 2, 3091 (2019)*

# Phase-change photonics

- Add active elements to passive waveguides
- Implement synapses and neuron soma with phase change materials (PCMs)
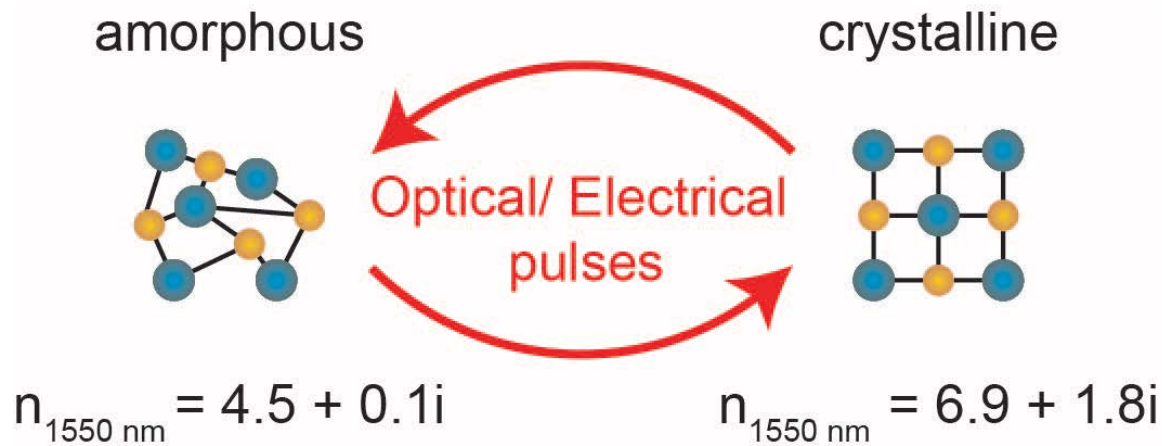


amorphous

$n_{1550\ nm} = 4.5 + 0.1i$

crystalline

$n_{1550\ nm} = 6.9 + 1.8i$

$Ge_2Sb_2Te_5$ (GST)

# Phase-change photonics

- Add active elements to passive waveguides
- Implement synapses and neuron soma with phase change materials (PCMs)
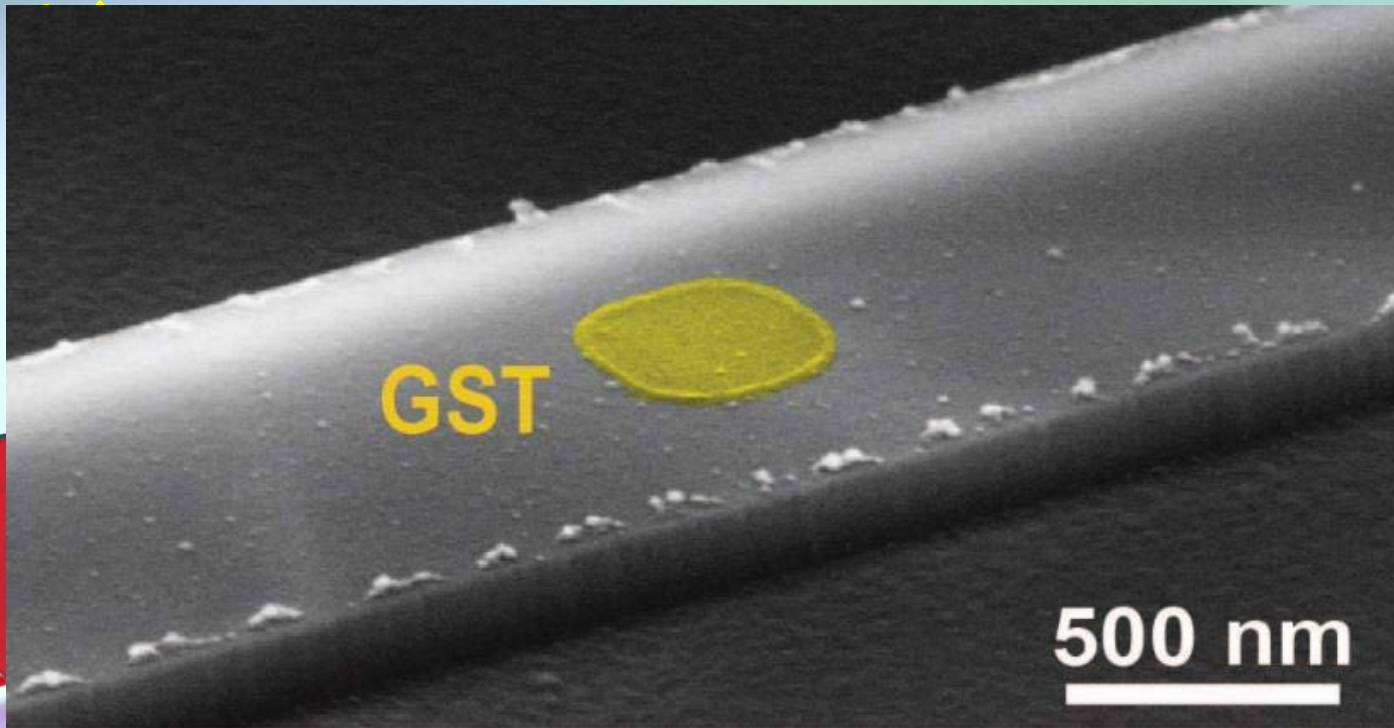- All-optical reconfiguration within sub-nanoseconds



amorphous        crystalline

Optical/ Electrical pulses

$n_{1550\ nm} = 4.5 + 0.1i$       $n_{1550\ nm} = 6.9 + 1.8i$

**$Ge_2Sb_2Te_5$ (GST)**

# PCM nanophotonic devices

- Place PCM in near-field of optical waveguide
- Data is encoded in the amount of transmitted power
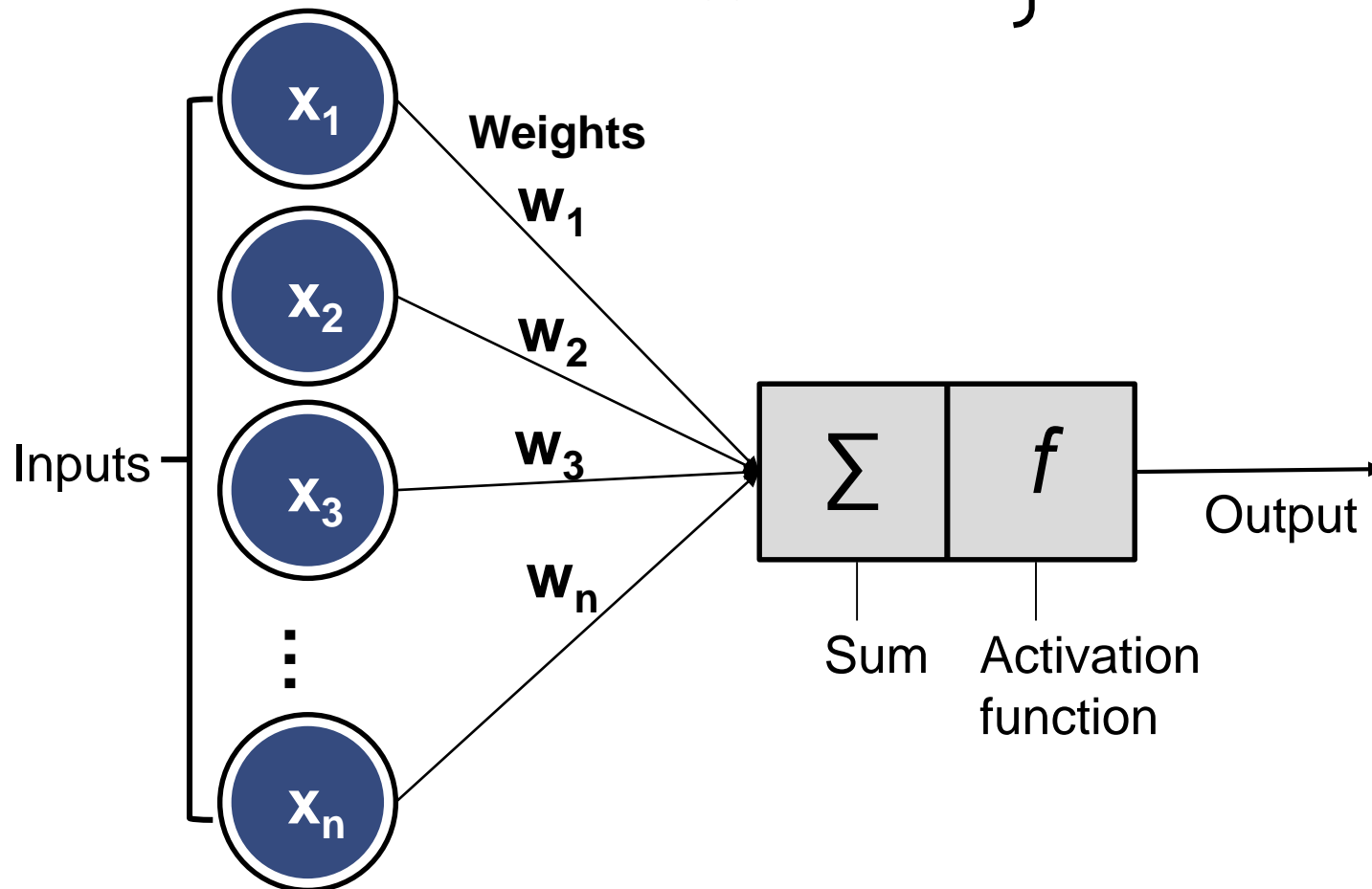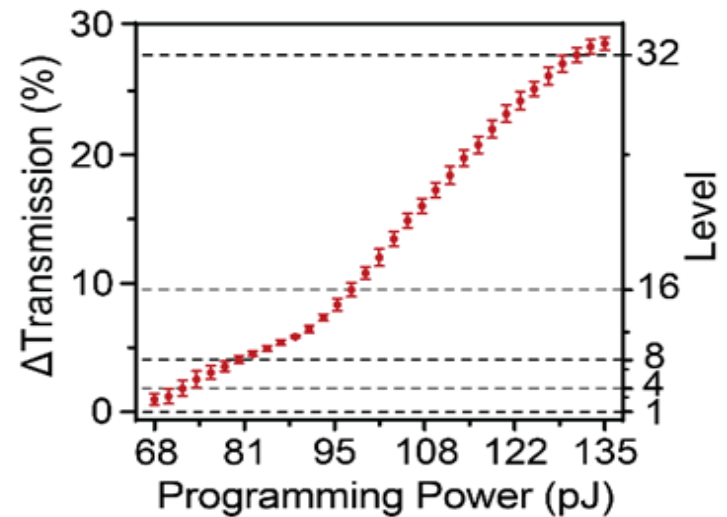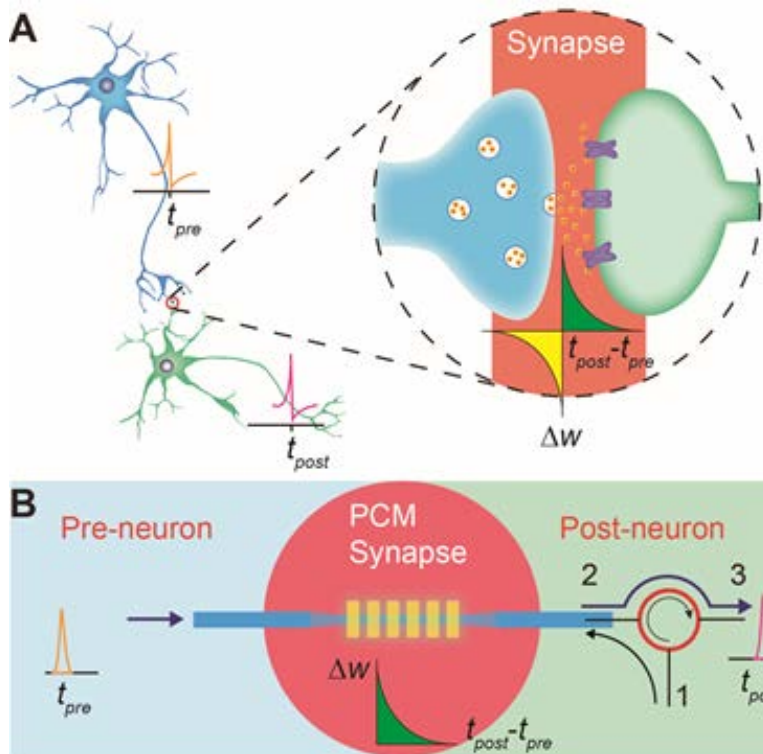
# Photonic neurons

- Use for matrix multiplication:

- Multiplication
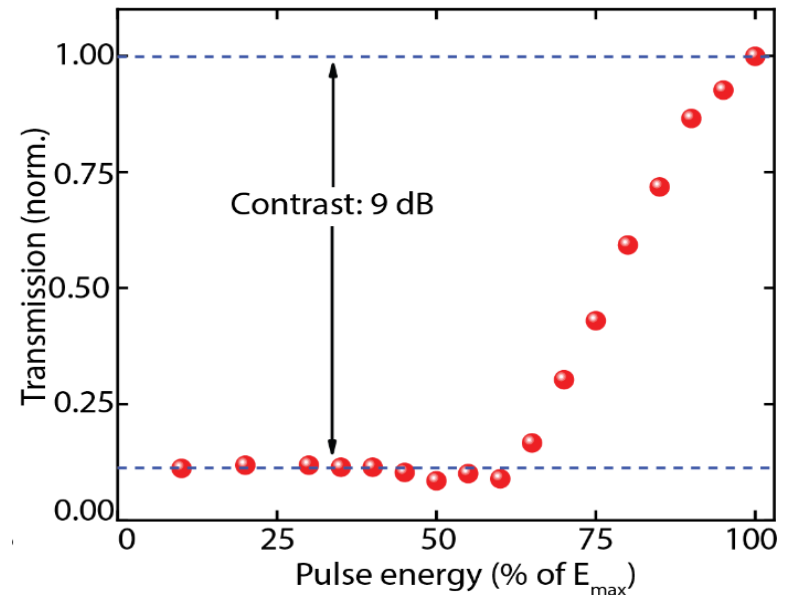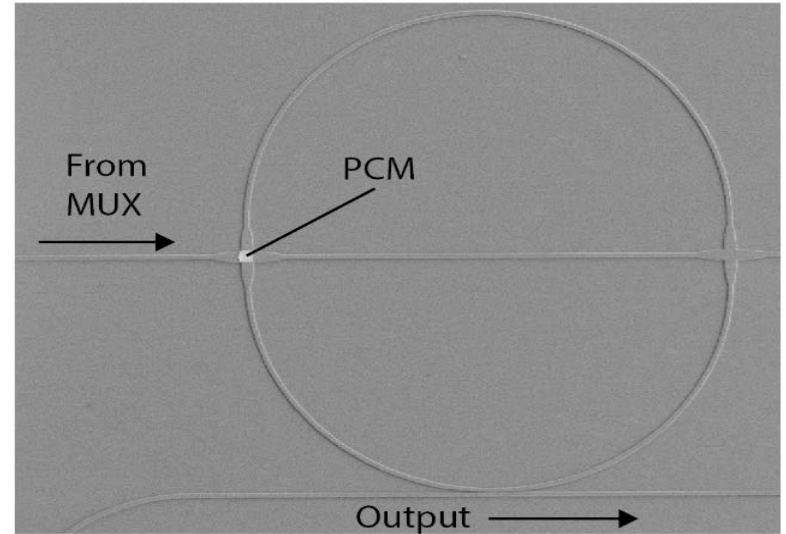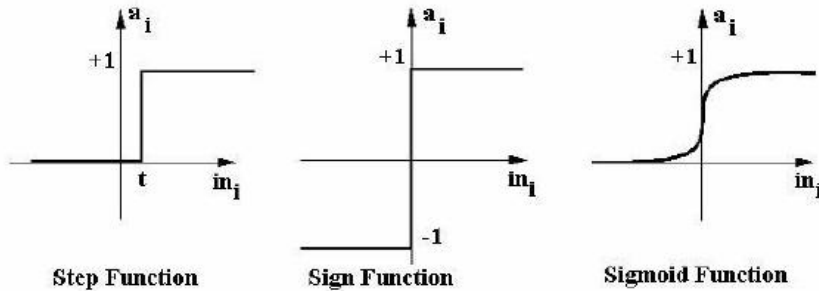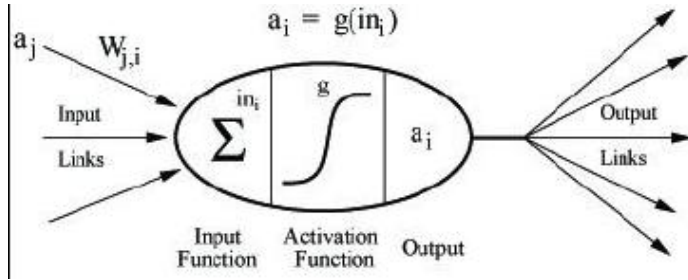- Addition

Multiply-accumulate (MAC)

# A photonic synapse (the weights)



- Partial crystallization allows storage of multiple bits per cell
- Number of levels depends on optical contrast and noise performance
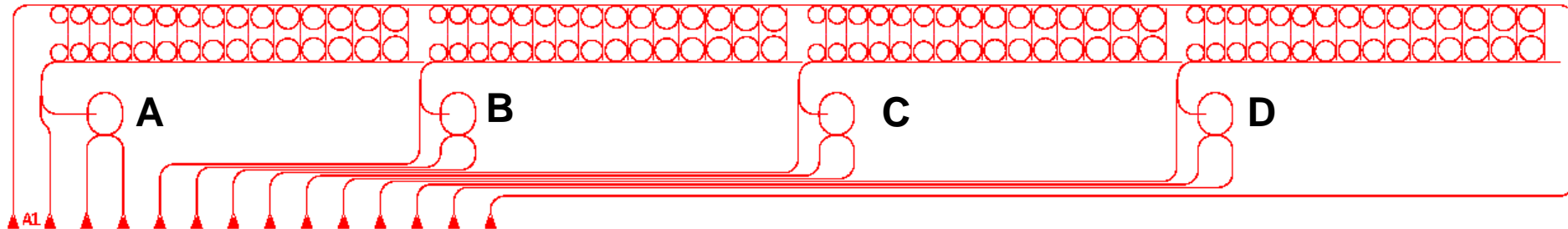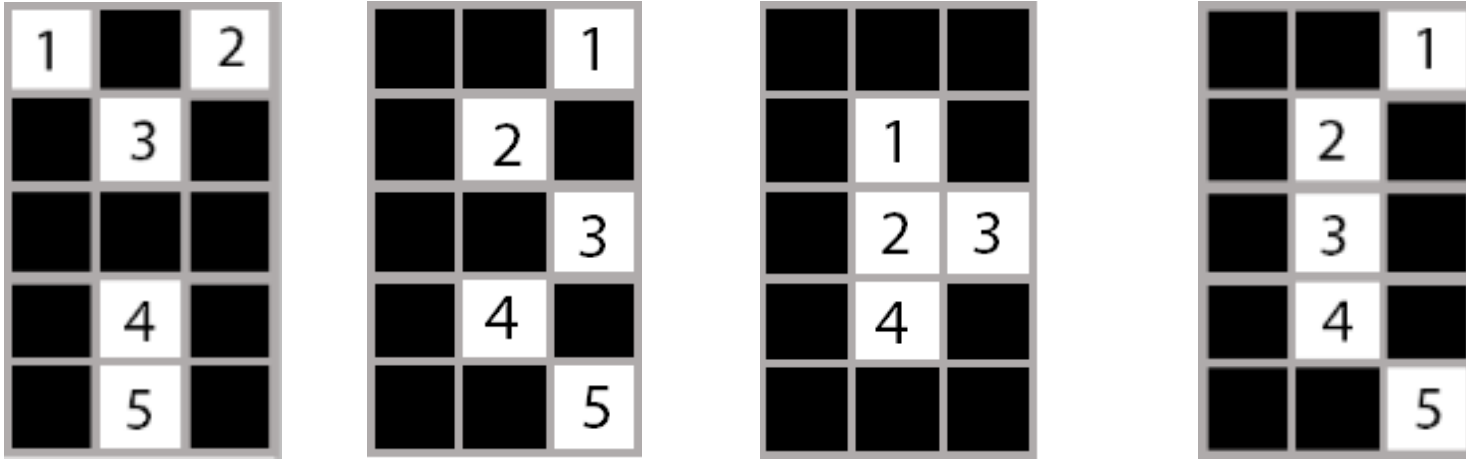- Wheights are stored permanently in crystal state of PCM

Rios et al., Science Advances 5, eaau5759 (2019)
Li, et al., Optica 6, 1, (2018)
Cheng, et al., Science Advances 3 e1700160 (2017)

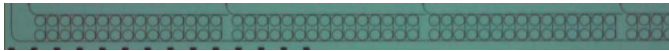# Implementation of threshold function (the „soma")



- Threshold function is provide by ring resonator
- Resonance tuning with embedded PCM element

# A small-scale ANN



- 15 input neurons and 4 output neurons
- Each letter is pixelized into 15 digital elements

*Feldmann et al., Nature 569, 208 (2019)*

# A closer look at the phontonic ANN

$$\begin{bmatrix} a_{11} & \cdots & a_{1N} \\ \vdots & \ddots & \vdots \\ a_{M1} & \cdots & a_{MN} \end{bmatrix} \quad \times \quad \begi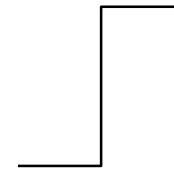n{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_N \end{pmatrix} \quad = \quad \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_M \end{pmatrix} \quad \times$$

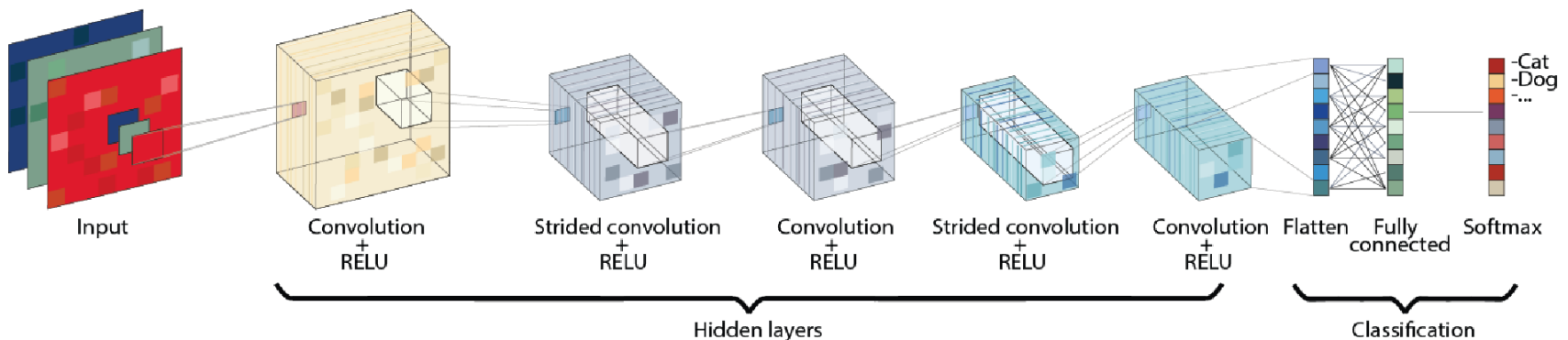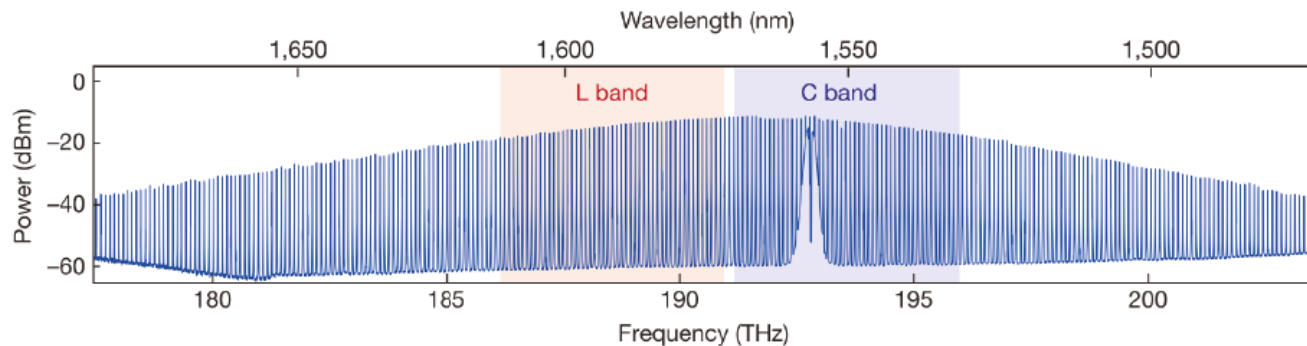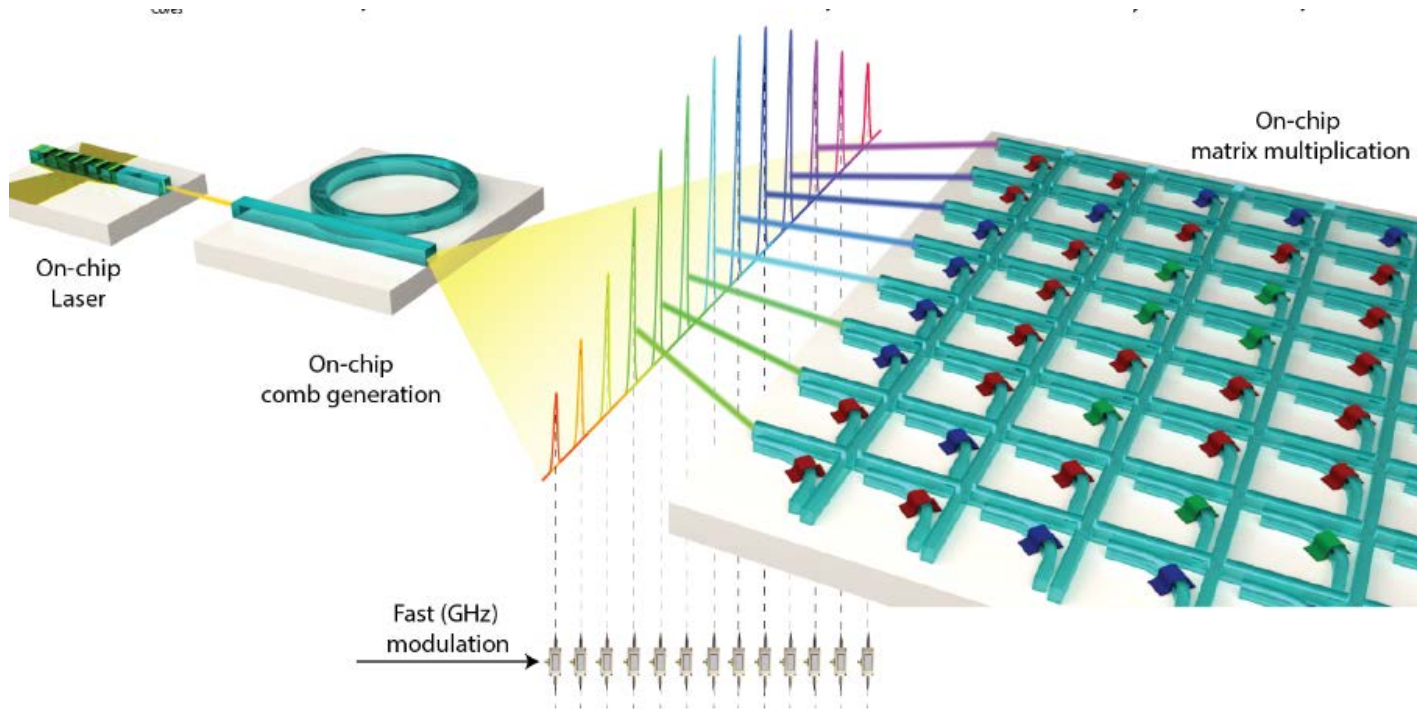**Synaptic weights**
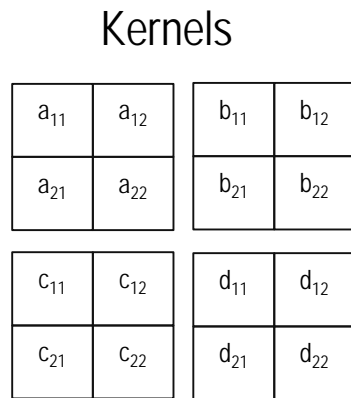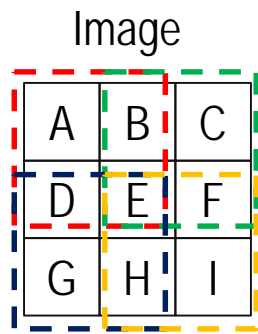
**Input-Vector**

**Rectification**

**Convolutional neural networks**

# Ultrafast convolution processing



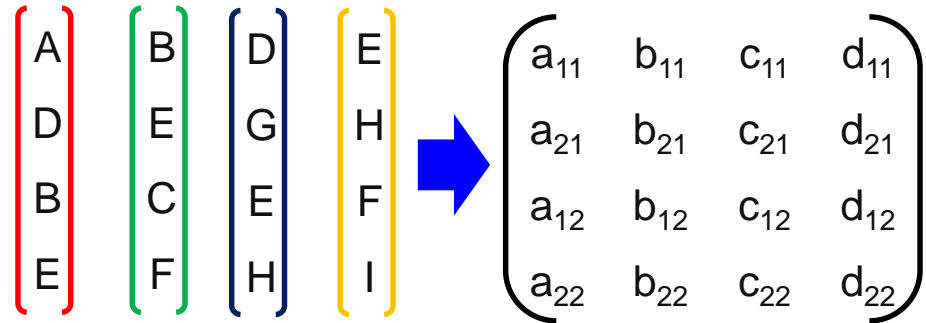**Frequency comb, Kippenberg group (EPFL)**
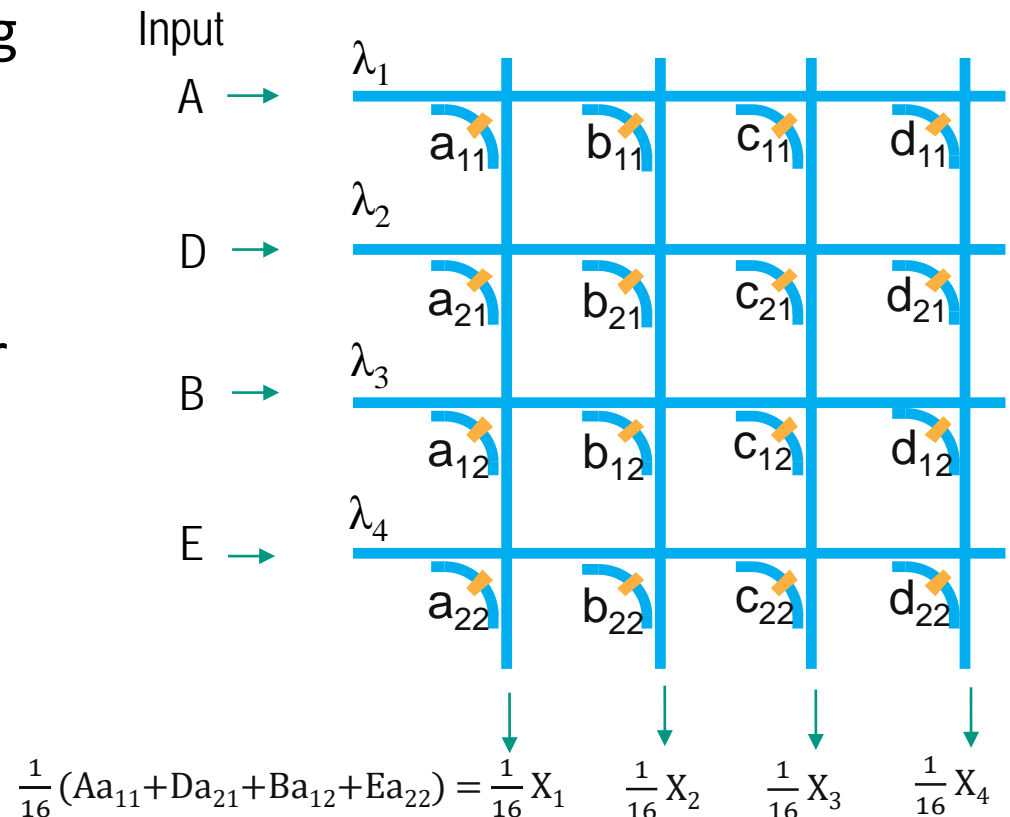
# Example: 3X3 pixel image and 4 kernels of 2X2 dimension

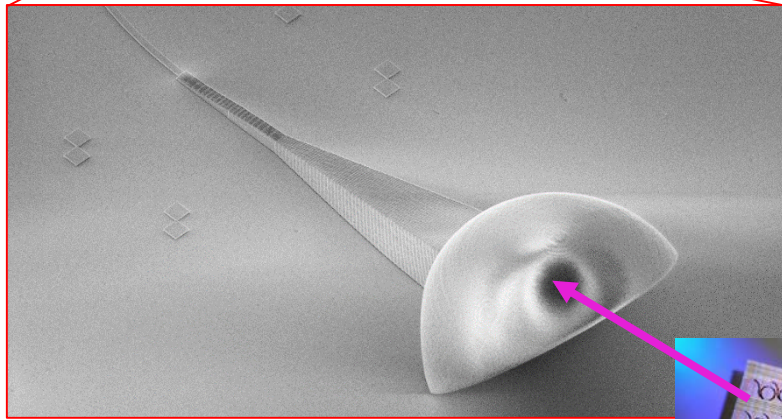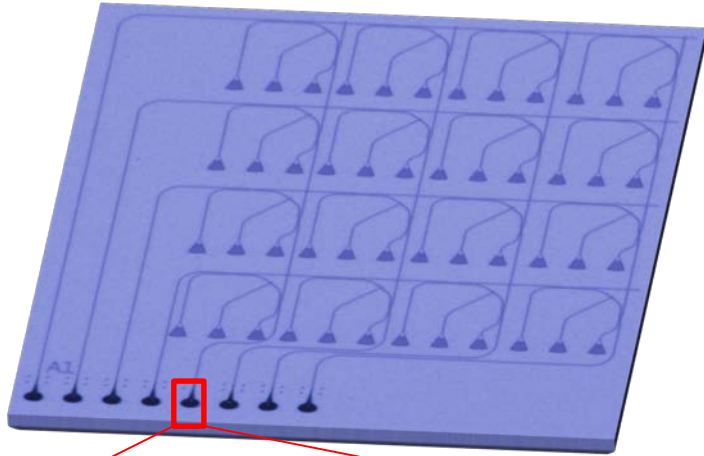# Example: 3X3 pixel image and 4 kernels of 2X2 dimension

- Using waveguide crossing matrix on the right

- Splitting ratios (directional couplers) adjusted for equal power distribution

- Each input on different wavelength (avoid interference)

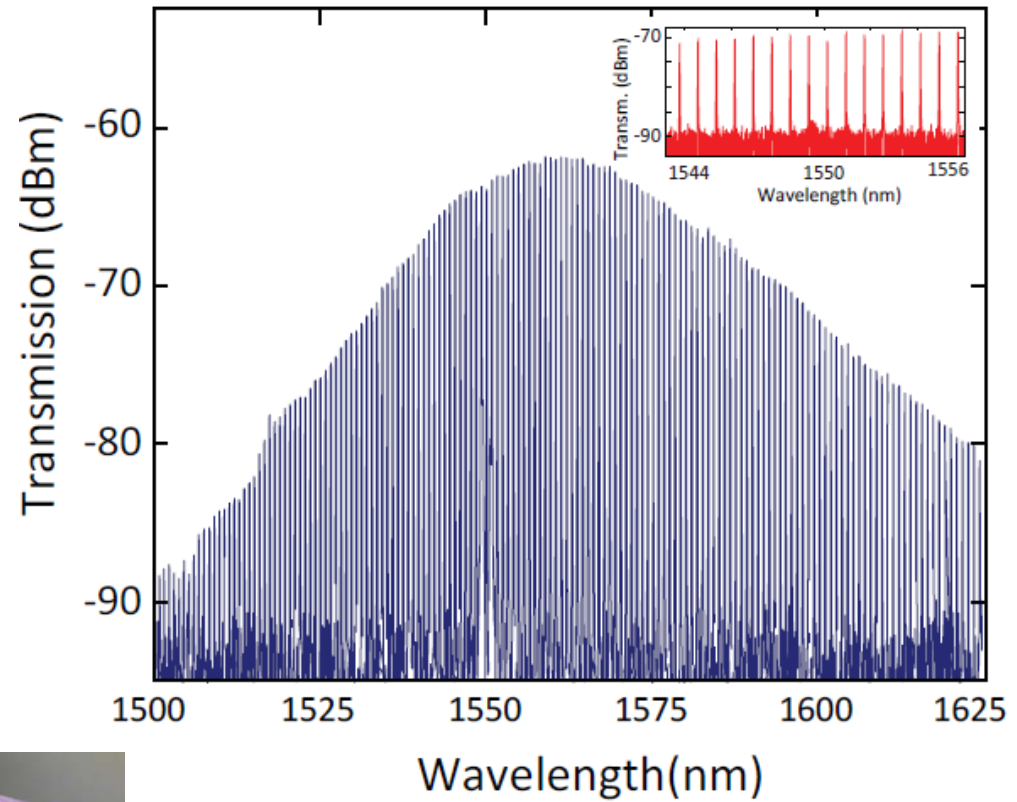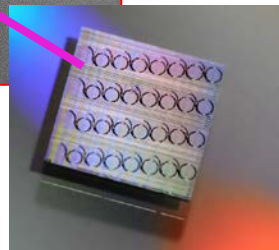- Each column gives convolution output for one of the kernels

Input

$\lambda_1$  A →

$\lambda_2$  D →

$\lambda_3$  B →

$\lambda_4$  E →

$a_{11}$  $b_{11}$  $c_{11}$  $d_{11}$

$a_{21}$  $b_{21}$  $c_{21}$  $d_{21}$

$a_{12}$  $b_{12}$  $c_{12}$  $d_{12}$

$a_{22}$  $b_{22}$  $c_{22}$  $d_{22}$

$$\frac{1}{16}(Aa_{11}+Da_{21}+Ba_{12}+Ea_{22}) = \frac{1}{16}X_1 \qquad \frac{1}{16}X_2 \qquad \frac{1}{16}X_3 \qquad \frac{1}{16}X_4$$

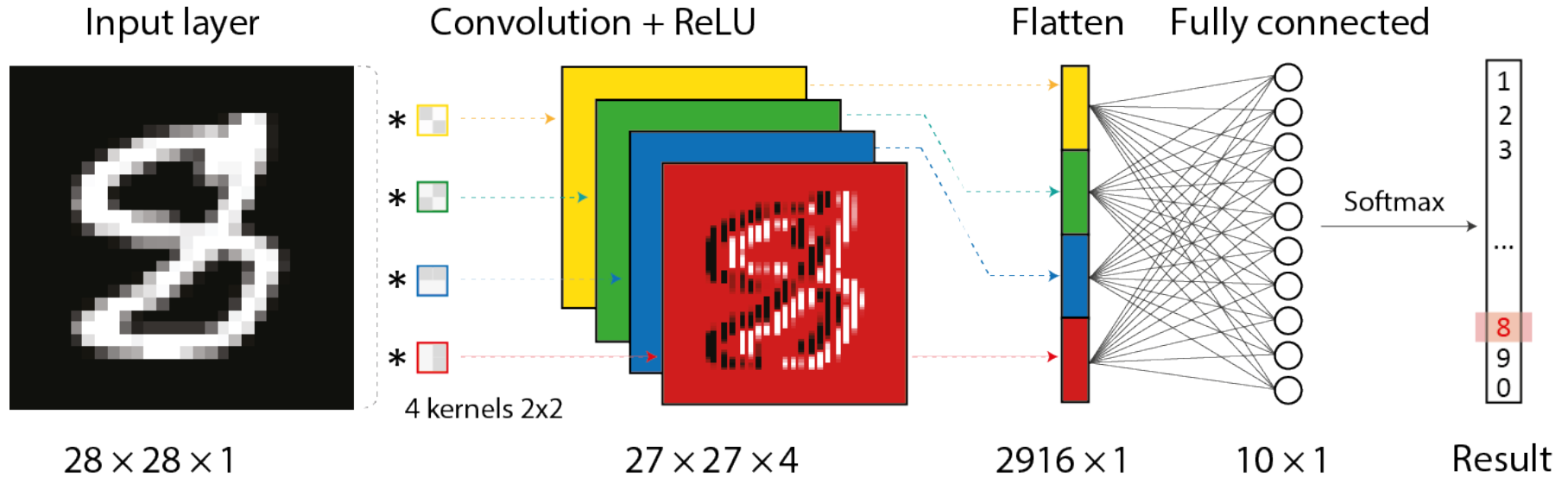Use WDM to perform multiple convolutions at the same time

# Ultrafast convolution processing

**PCM Matrix chip**
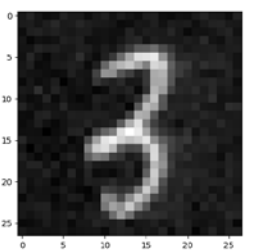


**Comb input, EPFL**



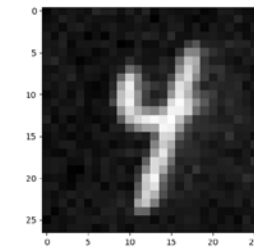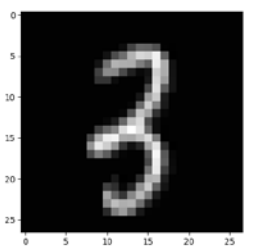*Feldmann et al., Nature 589, 52 (2021)*

# Full digit recognition with photonic NNs



**~95% accurate**

*Feldmann et al., Nature 589, 52 (2021)*

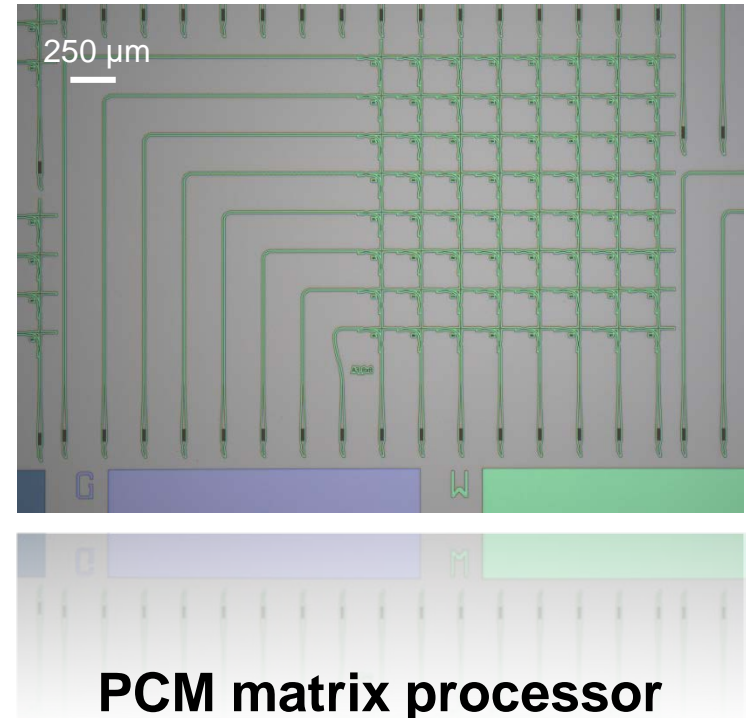# Why this is exciting



**NVIDIA flagsip tensor processor (Tesla V100)**

- 195 GigaFLOPS/core
- 1 FLOP ~ 0.5 multiply-accumulate (MAC) operations



250 µm

**PCM matrix processor**

- **2 TeraMAC/s**

More than 1000 HDTV video streams in parallel

# The people who really do the work:

At WWU:

C. Schuck and team

R. Bratschitsch, J. Kern, P. Tonndorf

J. Feldmann, N. Gruhler, A. Ovvyan, S. Ferrari, W. Hartmann, N. Walter, F. Beutel, M. Stappers, H. Gehring, C. Kaspar, F. Lenzini, T. Grottke, J. Lin, J, Schütte, E. Lomonte, R. Stegmüller, Y Liu

At Oxford:
N. Youngblood
H. Bhaskaran
X. Li

At Exeter:
D. Wright
E. Gemo
S. Garcia-Cuevas
Carrillo

At EPFL:
T. Kippenberg
M. Karpov

At IBM:
A. Sebastian