



# Photonic AI Accelerators

**Nikos Pleros**

Department of Informatics, Aristotle University of Thessaloniki, Greece  
Center for Interdisciplinary Research and Innovation, AUTH, Greece

# AI computing and energy consumption



## MEGATRON-TURING NLG

### Using DeepSpeed and Megatron to Train Megatron-Turing NLG 530B, A Large-Scale Generative Language Model

Shaden Smith<sup>§,†</sup>, Mostofa Patwary<sup>§,‡</sup>, Brandon Norick<sup>†</sup>, Patrick LeGresley<sup>‡</sup>, Samyam Rajbhandari<sup>†</sup>, Jared Casper<sup>‡</sup>, Zhun Liu<sup>†</sup>, Shrimai Prabhumoye<sup>‡</sup>, George Zerveas<sup>\*,†</sup>, Vijay Korthikanti<sup>‡</sup>, Elton Zhang<sup>†</sup>, Rewon Child<sup>‡</sup>, Reza Yazdani Aminabadi<sup>†</sup>, Julie Bernauer<sup>‡</sup>, Xia Song<sup>†</sup>, Mohammad Shoeybi<sup>‡</sup>, Yuxiong He<sup>†</sup>, Michael Houston<sup>‡</sup>, Saurabh Tiwary<sup>†</sup>, and Bryan



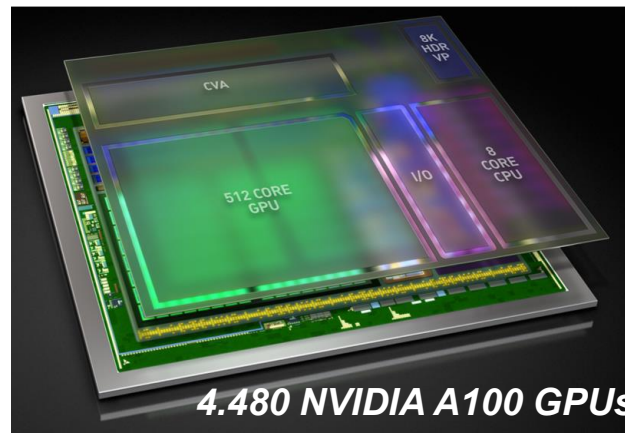
Catanzaro<sup>‡</sup>

arXiv, 2022

**530B parameters &  $10^{24}$  FLOPs compute power**



consumes energy that  
**powers a small city**  
(16k people) **for 1 year**



**4.480 NVIDIA A100 GPUs**

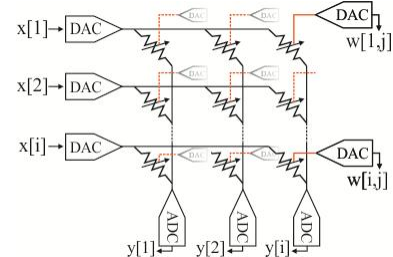
**...for 11 days**

**400MWh**

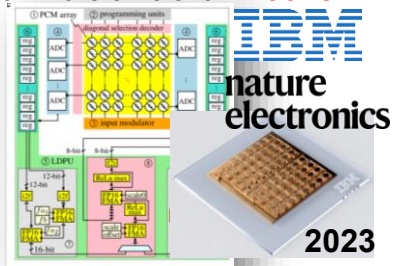


# Escape through analog photonic MVMs

## Neuromorphic electronics

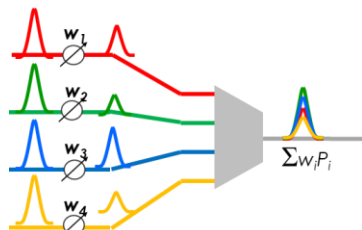


>5x more efficient: ~200 fJ/MAC



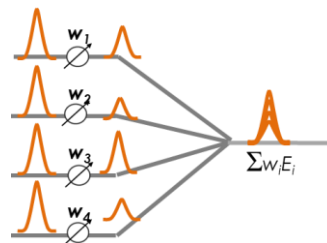
## wavelength

## Incoherent photonic accelerators



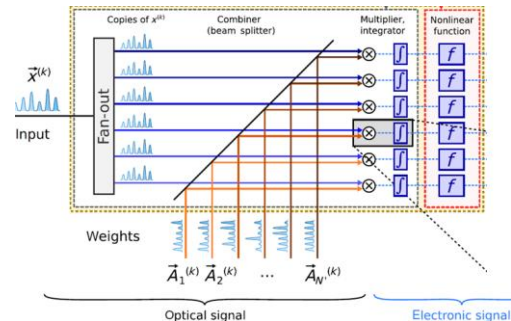
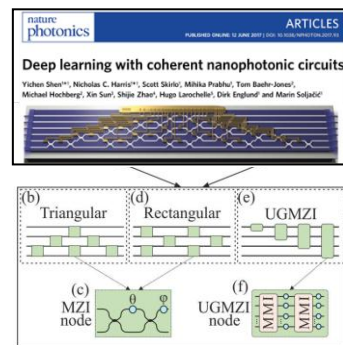
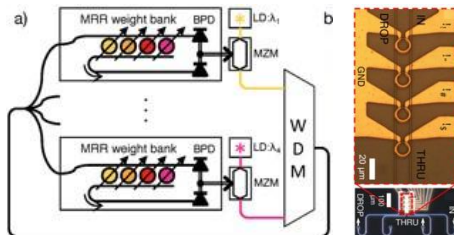
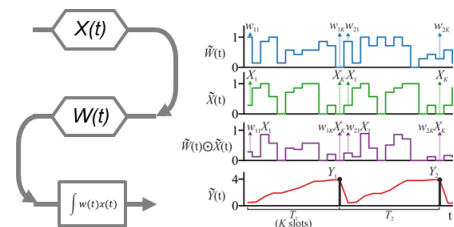
## space

## Coherent photonic accelerators

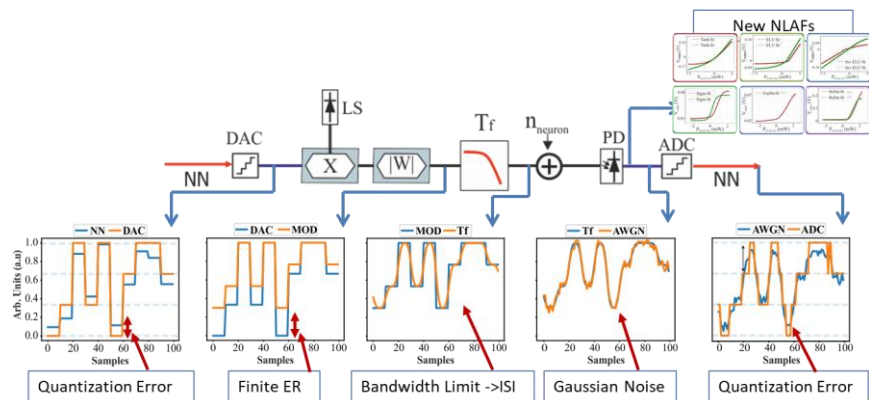


## time

## Time-multiplexed photonic accelerators



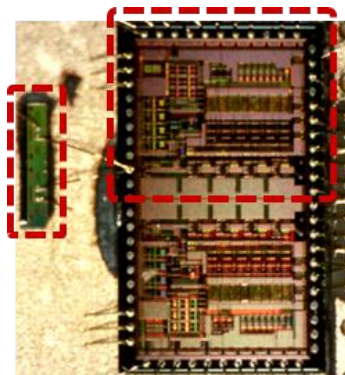
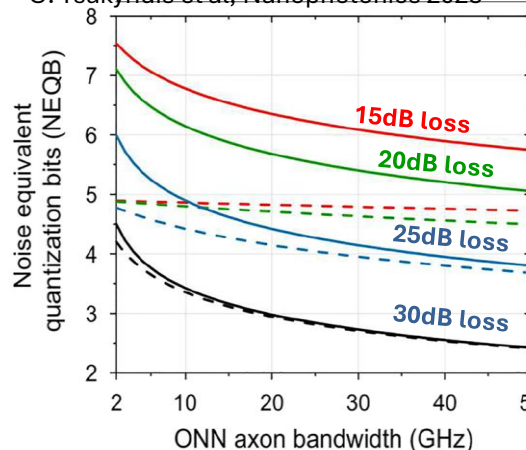
# Lessons Learned



A. Tsakyridis et al, "Photonic Neural Network and Optics-informed Deep Learning Fundamentals", (Invited Tutorial) APL Photonics 2024

✓ **Optics-Informed Deep Learning:**  
embeds ER, BR, noise, BW etc

G. Tsakyridis et al, Nanophotonics 2023



C. Pappas et al, JSTQE 2023  
S. Kovaivos et al, JLT 2024

✓ High-speed ADCs reduce bit-resolution and energy efficiency

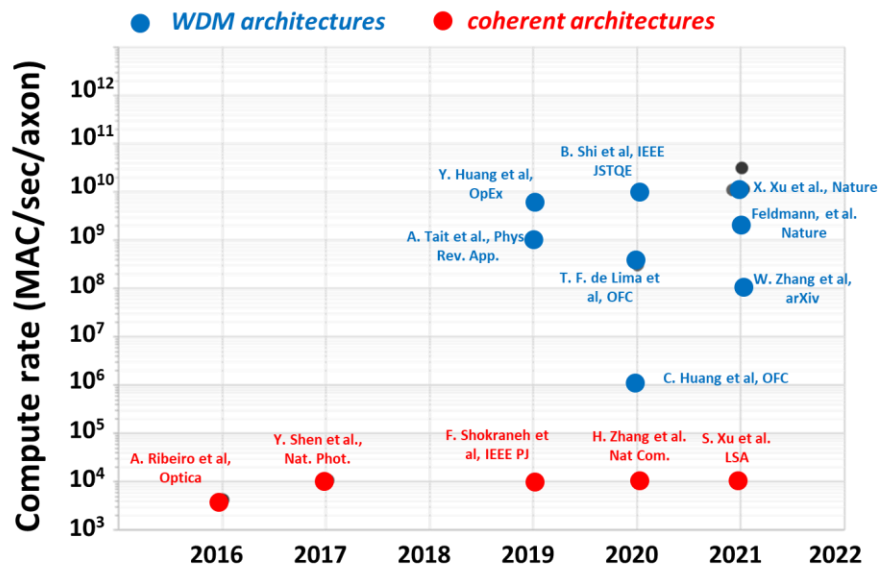
**Lesson #1:**

**Hardware/Software co-design**

**Lesson #2:**

**<GHz ADC via TDM PNNs**

# Lessons Learned



✓ **Incoherent architectures** support 10's GHz clock rates, but:

- require one wavelength per axon: **no WDM capability**

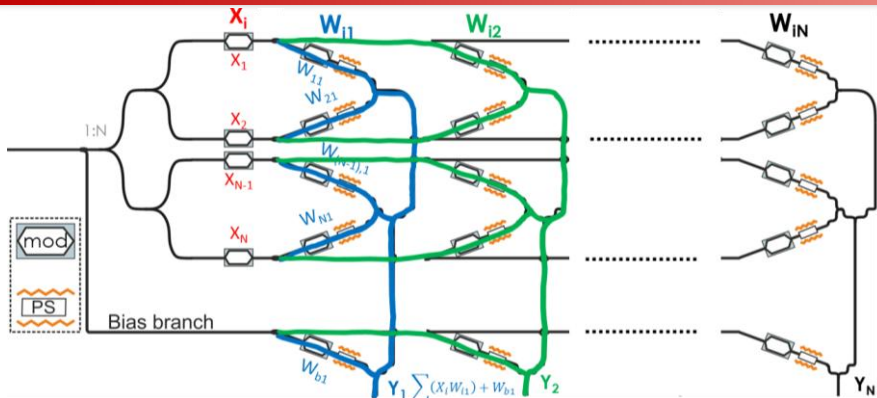
✓ **SVD coherent architectures** require a single wavelength, but:

- hardly exceed **GHz range**
- slow weight update – **small NN sizes**
- **degrade fidelity**

**Lesson #3: need for a new MVM architecture** (that supports TDM)



# The photonic Xbar processor

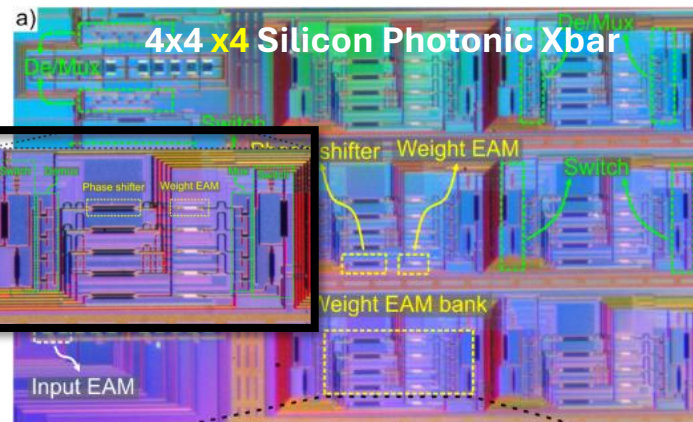
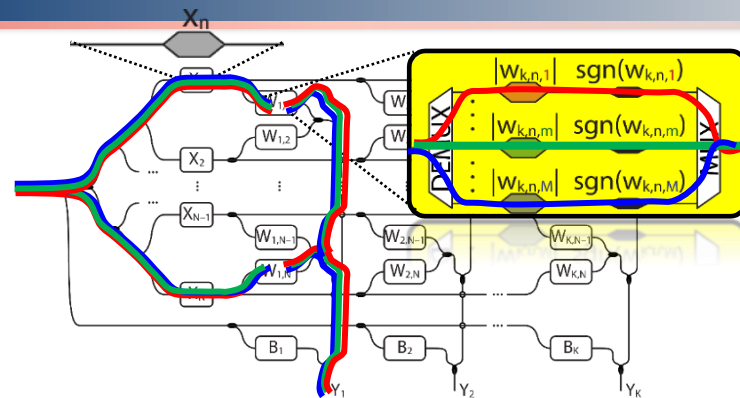
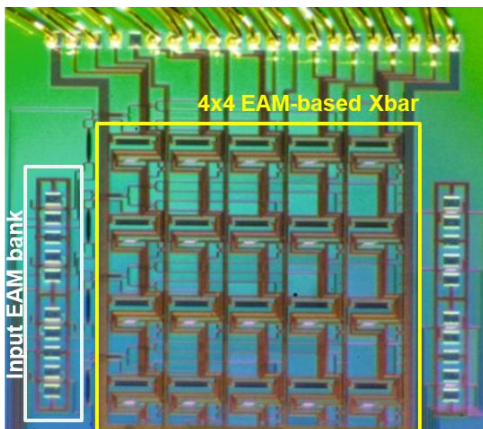
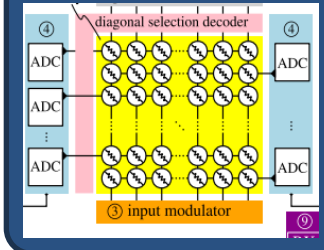


G. Giamougiannis et al, JLT, 2023

G. Giamougiannis et al, JSTQE 2023

M. Moralis-Pegios et al, Nat. Comms, 2024

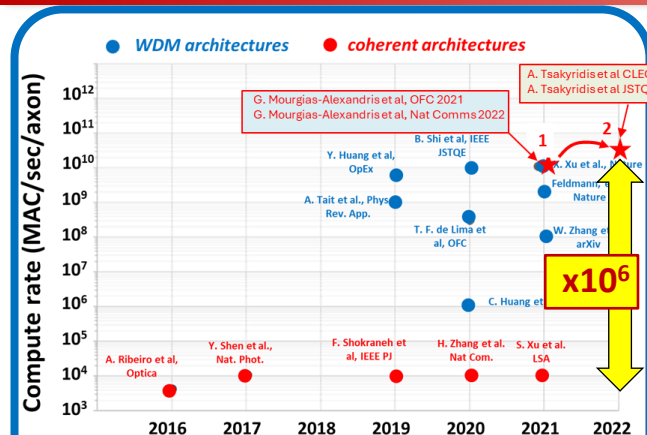
the equivalent  
to the analog  
electronic Xbar



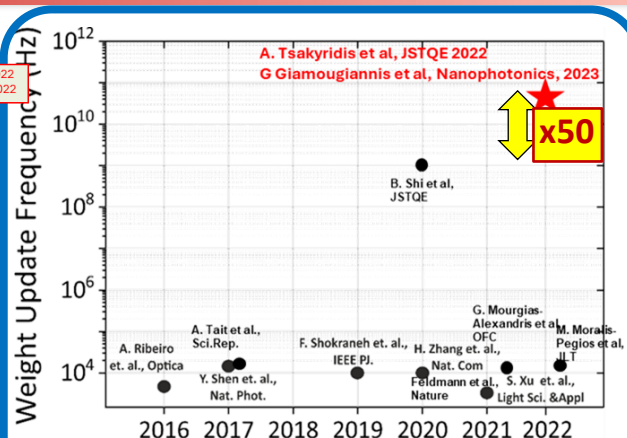
S. Kovaivos et al, JLT 2025

A. Totovic et. al., Neuromorph. Comput. Eng. (2022)

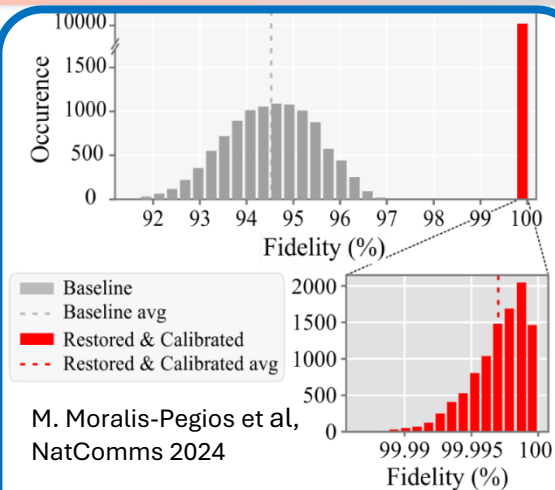
# The photonic Xbar processor



- ✓ Record-high **50GHz compute rate**
- ✓  **$\times 10^6$  improvement** (now 60GHz, Z. Lin et al, NatComms 2024)

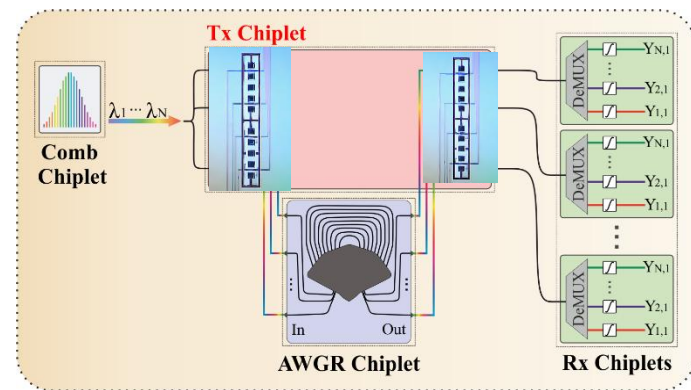
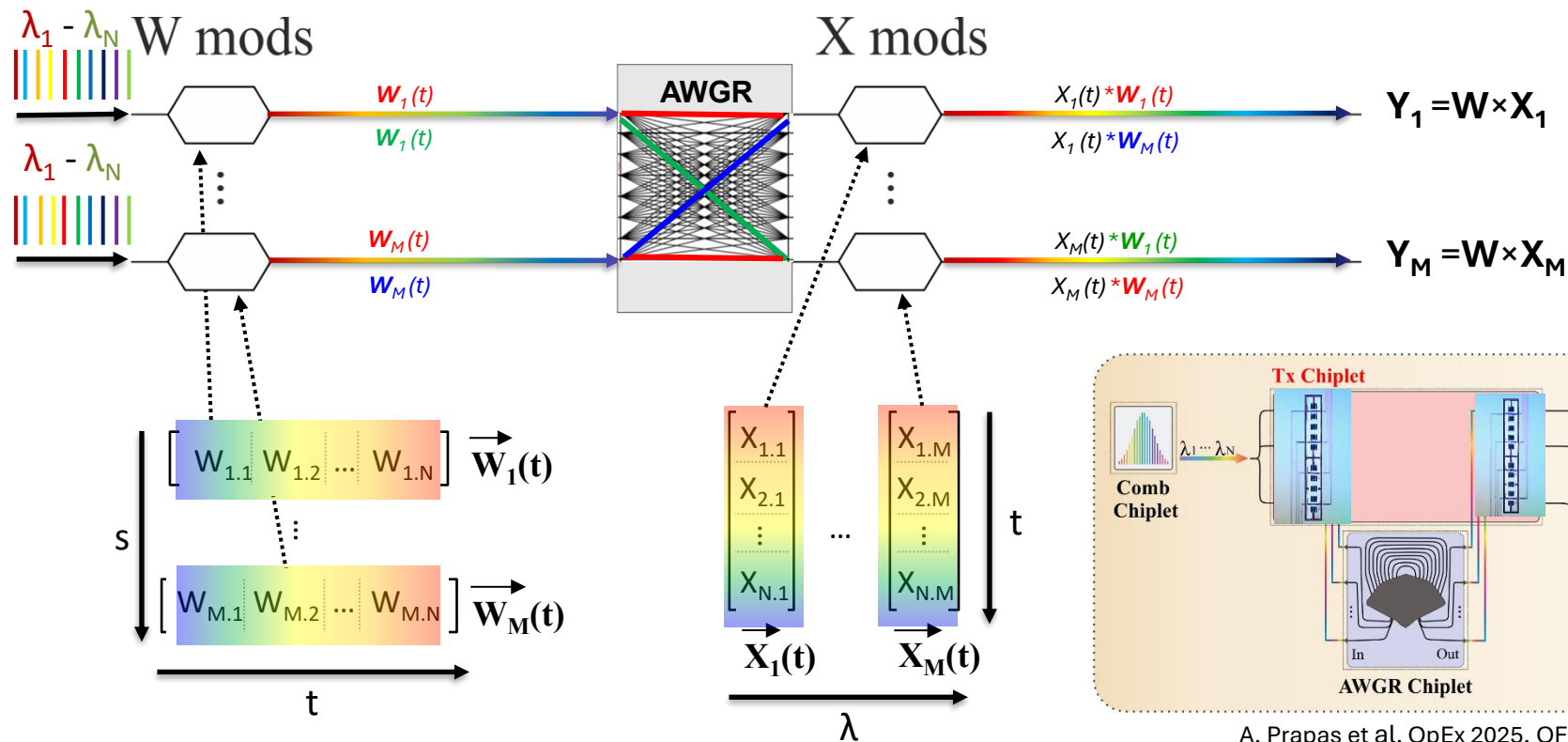


- ✓ **50Gs/sec photonic GeMM** via 50x faster weight update rates
- ✓ **Supports TDM**, i.e. NN size  $\gg$  PNN



- ✓ Fidelity restoration &..
- ✓ **Record at 99.997 %**

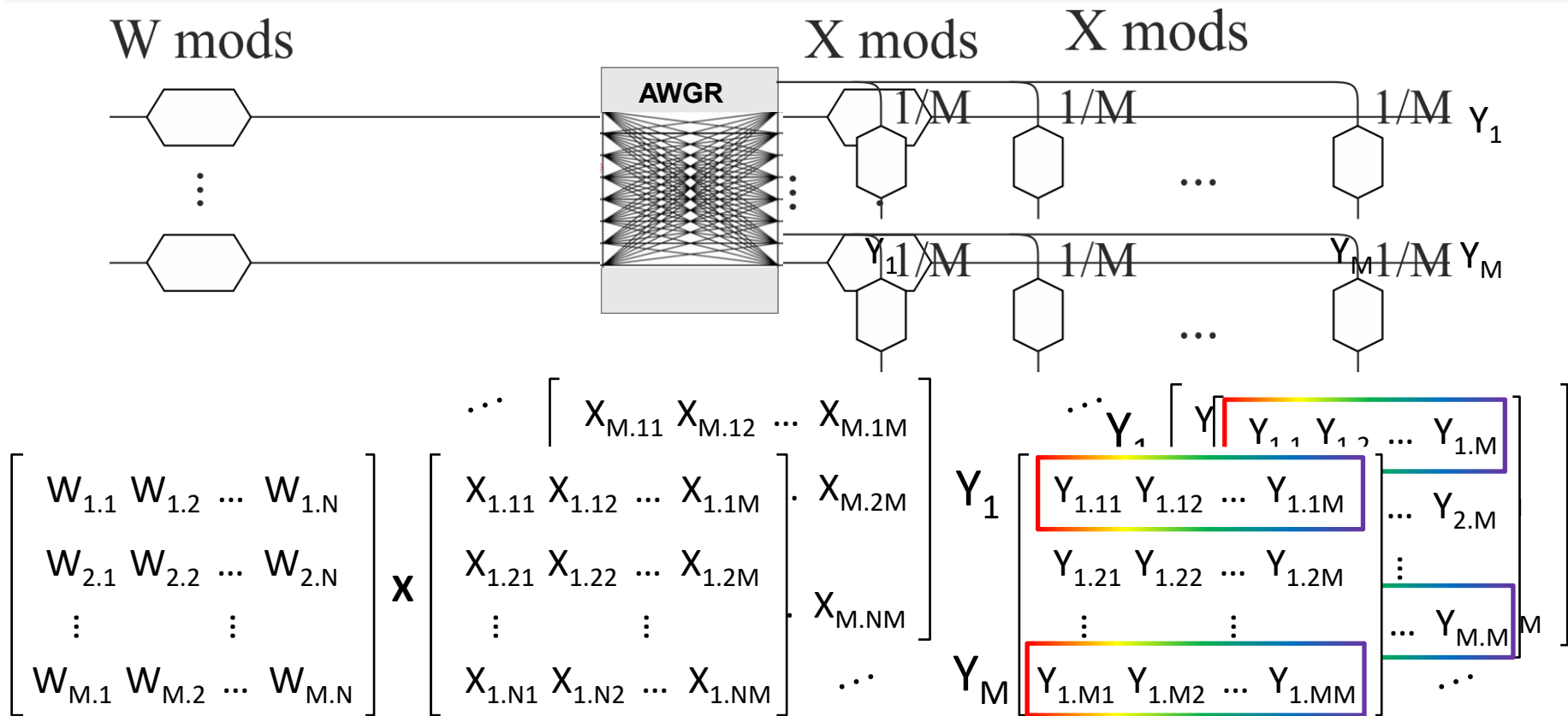
# ...and a recent AWGR-based TSWDM design



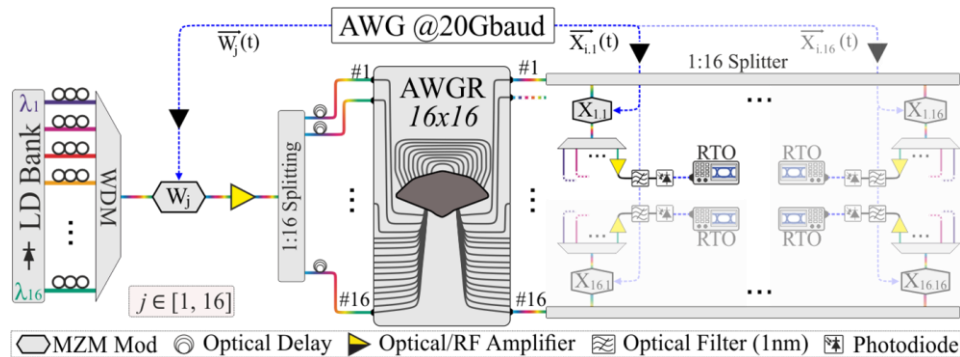
A. Prapas et al, OpEx 2025, OFC 2025



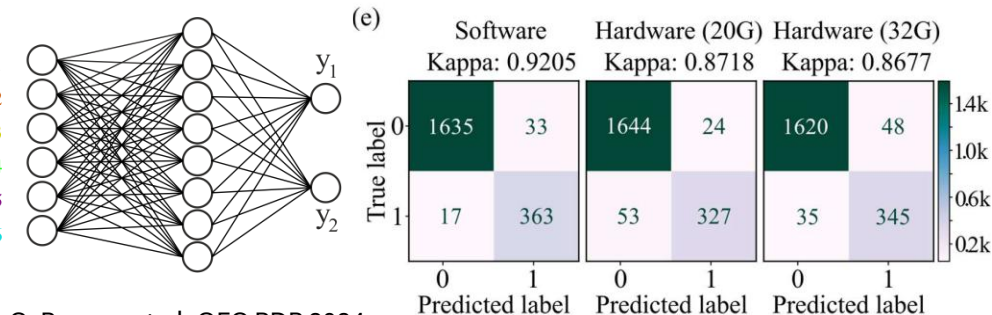
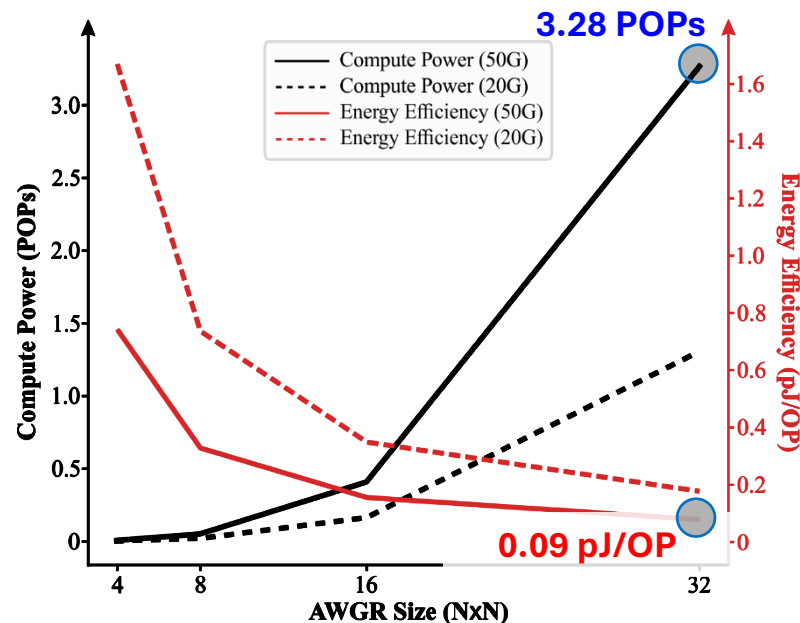
# 262 TOPs Photonic Tensor Core



# DDoS Attack Detection



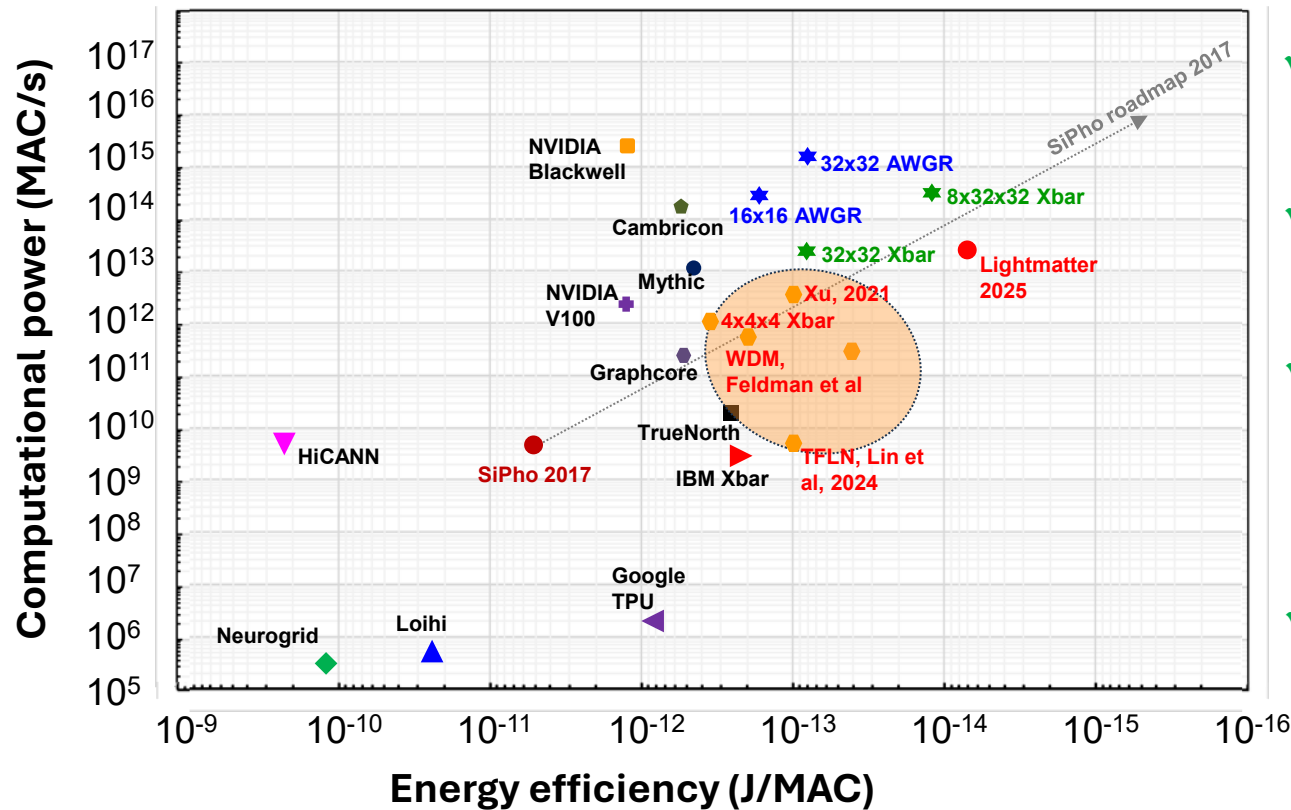
- ✓ **Compute** scales with  $O(N^3)$
- ✓ **hardware** scales with  $O(N^2)$



C. Pappas et al, OFC PDP 2024

C. Pappas et al, JLT 2025

# Break the PetaMAC/sec barrier

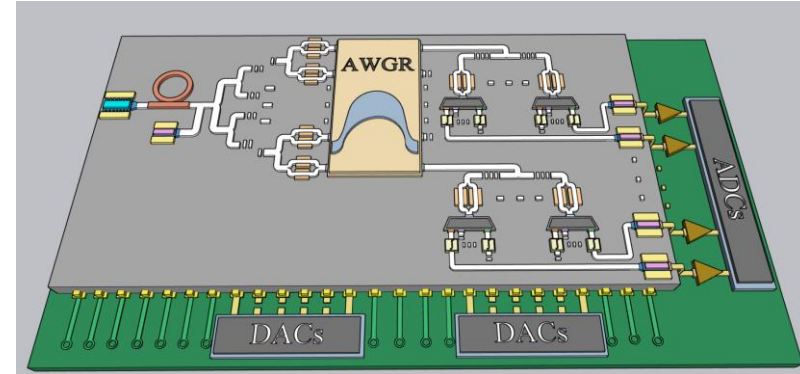
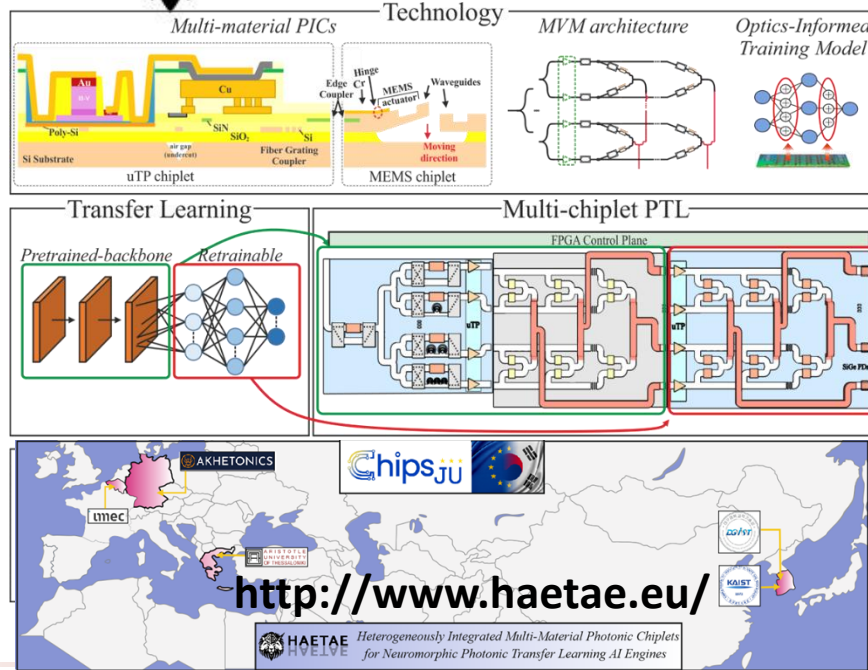


- ✓ On-chip reconfigurable PNNs up to **10 TOPS**
- ✓ Lightmatter Photonic Tensor Core: **65 TOPS**
- ✓ 4x4x4 Xbar **1 TOPS on-chip**: 32x32x8 Xbar projected to reach **400 TOPS**
- ✓ **262 TOPS demonstrated** – 32x32@50G AWGR projected to **>3 TOPS !**

# On-chip in 3 years from now



# HAETAELIVE



[C. Pappas et al, "A 262 TOPs Hyperdimensional Photonic AI Accelerator", submitted at APL Photonics, arXiv pre-print]

# Acknowledgements



<https://sipho-g.highlite-h2020.eu/>



**GATEPOST  
PROJECT**



**HAETAELVE**



*Thank you!*



Visit us at <http://winphos.web.auth.gr/>

