Single chip photonic deep neural network with ADDITIONICS accelerated training

Saumil Bandyopadhyay¹, Alexander Sludds¹, Stefan Krastanov¹,

Ryan Hamerly^{1,2}, Nicholas Harris¹, Darius Bunandar¹, Matthew Streshinsky³, Michael Hochberg⁴, Dirk Englund¹

¹Massachusetts Institute of Technology ²PHI Labs, NTT Research ³Nokia Corporation ⁴Luminous Computing



The "Cambrian explosion" of DNN hardware







Waterloo CMOS Design and Reliability











Preprint manuscript available at https://arxiv.org/abs/2208.01623

|4ii

DNNs are getting better because they are getting bigger



DNN scaling is *bottlenecked* by energy consumption, which is ~1 pJ/OP for CMOS.

Preprint manuscript available at https://arxiv.org/abs/2208.01623

New DNN hardware is needed for ultra-low latency processing



New applications require ultra-low latency (µs-ns) processing of (optical) data

Preprint manuscript available at https://arxiv.org/abs/2208.01623

New DNN hardware is needed for ultra-low latency processing



New applications require ultra-low latency (µs-ns) processing of (optical) data This requires time-of-flight (clockless) compute on a single, integrated chip

Preprint manuscript available at https://arxiv.org/abs/2208.01623

An end-to-end optical DNN processor on a single chip



Preprint manuscript available at https://arxiv.org/abs/2208.01623

Шіг

An end-to-end optical DNN processor on a single chip



Preprint manuscript available at https://arxiv.org/abs/2208.01623

An end-to-end optical DNN processor on a single chip



Preprint manuscript available at https://arxiv.org/abs/2208.01623

Application-specific photonic integrated circuit



132 model parameters on-chip, 90 layers of optical

components Preprint manuscript available at https://arxiv.org/abs/2208.01623

PliT

Application-specific photonic integrated circuit



132 model parameters on-chip, 90 layers of optical

components

Preprint manuscript available at https://arxiv.org/abs/2208.01623

Illiī

Coherent matrix multiplication unit





 $T_{ij}(\theta, \phi) = ie^{i\theta/2} \frac{e^{i\phi}\sin(\theta/2)}{\left[e^{i\phi}\cos(\theta/2) - \sin(\theta/2)\right]} \quad U = D \prod T_{ij}(\theta, \phi)$

Programmable meshes of Mach-Zehnder interferometers compute linear algebra through passive optical interference

Coherent matrix multiplication unit



Nonlinear optical function unit



Nonlinear optical function unit



Coherent optical nonlinear functions without digitization or amplification

Model training benefits from optical acceleration



Model training also requires repeated forward inference.

Model training benefits from optical acceleration



Model training also requires repeated forward inference.



Optically acceleration can enable low-latency model training.

Plii

In situ training on optical hardware



Forward inference is **optically accelerated** during *in situ* training, which achieves the **same accuracy** as a digitally trained system.



• Single-shot optical inference in a deep neural network on a photonic chip with time-of-flight (**500 ps**) limited latency



- Single-shot optical inference in a deep neural network on a photonic chip with time-of-flight (**500 ps**) limited latency
- All-optical processing of data without digitization between layers, eliminating the latency introduced by optical-to-electronic conversion
 - Coherent matrix multiplication in the optical domain
 - Coherent optical nonlinear functions without electronic amplification



- Single-shot optical inference in a deep neural network on a photonic chip with time-of-flight (**500 ps**) limited latency
- All-optical processing of data without digitization between layers, eliminating the latency introduced by optical-to-electronic conversion
 - Coherent matrix multiplication in the optical domain
 - Coherent optical nonlinear functions without electronic amplification
- Entirely fabricated in a commercial foundry



- Single-shot optical inference in a deep neural network on a photonic chip with time-of-flight (**500 ps**) limited latency
- All-optical processing of data without digitization between layers, eliminating the latency introduced by optical-to-electronic conversion
 - Coherent matrix multiplication in the optical domain
 - Coherent optical nonlinear functions without electronic amplification
- Entirely fabricated in a commercial foundry
- In situ training takes advantage of near-instantaneous inference, reducing latency and power consumption of model training.

Plif

Acknowledgments



Alexander Sludds



Prof. Stefan Krastanov



Dr. Ryan Hamerly



Dr. Nicholas Harris



Dr. Darius Bunandar



Dr. Matthew Streshinsky



Dr. Michael Hochberg



Prof. Dirk Englund







