Silicon Photonics for Artificial Intelligence and Neuromorphic Computing

Bhavin J. Shastri^{1,2}, Thomas Ferreira de Lima², Chaoran Huang², Bicky A. Marquez¹, Sudip Shekhar³, Lukas Chrostowski³, and Paul R. Prucnal²

¹Department of Physics, Engineering Physics & Astronomy, Queen's University, Kingston, ON K7L 3N6, Canada ²Department of Electrical Engineering, Princeton University, Princeton, NJ 08544, USA ³Department of Electrical and Computer Engineering, University of British Columbia, Vancouver, BC V6T 1Z4, Canada shastri@ieee.org

Abstract: Artificial intelligence and neuromorphic computing driven by neural networks has enabled many applications. Software implementations of neural networks on electronic platforms are limited in speed and energy efficiency. Neuromorphic photonics aims to build processors in which optical hardware mimic neural networks in the brain. © 2021 The Author(s)

The field of neuromorphic computing aims to bridge the gap between the energy efficiency of von Neumann computers and the human brain. The rise of neuromorphic computing can be attributed the widening gap between current computing capabilities and current computing needs [1], [2]. Consequently, this has spawned research into novel brain-inspired algorithms and applications uniquely suited to neuromorphic processors. These algorithms attempt to solve artificial intelligence (AI) tasks in real-time while using less energy. We posit [3] that we can make use of the high parallelism and speed of photonics to bring the same neuromorphic algorithms to applications requiring multiple channels of multi-gigahertz analog signals, which digital processing struggles to process in real-time.

By combining the high bandwidth and parallelism of photonic devices with the adaptability and complexity attained by methods similar to those seen in the brain, photonic neural networks have the potential to be at least ten thousand times faster than state-of-the-art electronic processors while consuming less energy per computation [4]. An example is nonlinear feedback control; a very challenging task that involves computing the solution of a constrained quadratic optimization problem in real time. Neuromorphic photonics can enable new applications because there is no general-purpose hardware capable of dealing with microsecond environmental variations [5].

Neuromorphic photonics approaches: Neuromorphic photonic [6] approaches can be divided into two main categories: coherent (single wavelength) and incoherent (multiwavelength) approaches. Neuromorphic systems based on reservoir computing [7], [8] and Mach-Zehnder interferometers [9], [10] are example of coherent approaches. In reservoir computing the predefined random weights of their hidden layers cannot be modified. An alternative approach uses silicon photonics to design fully programmable neural networks [5], with a so-called broadcast-and-weight protocol [11-13]. In this architecture, photonic neurons output optical signals with unique wavelengths. These are multiplexed into a single waveguide and broadcast to all others, weighted, and photodetected. Each connection between a pair of neurons is configured independently by one microring resonator (MRR) weight, and the wavelength division multiplexed (WDM) carriers do not mutually interfere when detected by a single photodetector. Consequently, the physics governing the neural computation is fully analog and does not require any logic operation or sampling, which would involve serialization and sampling. Thus, they exhibit distinct, favorable trends in terms of energy dissipation, latency, crosstalk and bandwidth when compared to electronic neuromorphic circuits [4]. The advantage of this approach over the aforementioned approaches is that it has already demonstrated fan-in, inhibition, time-resolved processing, and autaptic cascadability [14].

However, the same physics also introduce new challenges, especially reconfigurability, integration, and scalability. Information carried by photons is harder to manipulate compared to electronic signals, especially nonlinear operations and memory storage. Photonic neurons described here solve that problem by using optoelectronic components (O/E/O), which can be mated with standard electronics providing reconfigurability. However, neuromorphic photonic circuits are challenging to scale up because they do not benefit from digital information, memory units and a serial processor, and therefore requires a physical unit for each element in a neural network, increasing size, area and power consumption. Here, integration costs must also be considered, since the advantages of using analog photonics (high parallelism and high bandwidth) must outweigh the costs of interfacing it with digital electronics (requiring both O/E and analog/digital conversion).

Advances in Science and Technology to Meet Challenges: In a recent review [3] and roadmap article [15], we outlined some scientific and technological advances necessary to meet the challenges to envision a neuromorphic

processor outlined in [16]. Photonic processors have light sources, passive and active devices. Currently, there is no single commercial fabrication platform that can simultaneously offer devices for light generation, wavelength multiplexing, photodetection, and transistors on a single die; state-of-the-art devices in each of these categories use different photonic materials (SiN, Ge, InP, GaAs, 2D materials, etc) with incongruous fabrication processes (silicon-on-insulator, CMOS, FinFETs). Silicon photonics is becoming an ideal platform for integrating these devices while offering a combination of foundry compatibility, device compactness, and cost that enables the creation of scalable photonic systems on chip.

Materials: Energy efficient and fast switching optical and electro-optical materials are needed for non-volatile photonic storage and weighting, as well as high-speed optical switching and routing, with low power consumption. Neural non-linearities are already possible on mainstream platforms using electrooptic transfer functions [14], but new materials promise significant performance opportunities. Phase change materials (PCMs), and graphene and ITO-based modulators can also be utilized for implementing non-linearities. Plasmonic PCMs can bridge the optical and electrical signals, through the dual operation modes [17]. A general material design method is in urgent need to develop appropriate photonic materials for different photonic components [18].

Lasers and amplifiers: On-chip optical gain and power will require co-integration with active InP lasers and semiconductor optical amplifiers. Current approaches involve either III–V to silicon wafer bonding (heterogeneous integration) or co-packaging with precise assembly (hybrid approach) [19]. Quantum dot lasers are another promising approach as they can be grown directly onto silicon, but fabrication reliability does not currently reach commercial standards [20].

Electrical control: Co-integrating CMOS controller chips with silicon photonics to provide electrical tuning control/stabilization will be critical. Candidates include wire-bonding, flip-chip bonding, 2.5D integration (interposers), 3D stacking (through-silicon-vias), and monolithic integration. Each has performance and design tradeoffs [21].

System packaging: A photonic processor must be interfaced with a computer. It would need to be selfcontained, robust to temperature fluctuations, and with electrical inputs/outputs [5]. Currently, manufacturers do not assemble electrical/thermal elements and chip-to-fiber interconnects.

Algorithms: Significant advances will be required to map abstract neural algorithms to photonic processor to usher these platforms into the commercial space. So far, only individual devices and small control circuits are described in the literature. The goal is to enable neural network programming tools (TensorFlow) to directly reconfigure a neuromorphic photonic processor [5].

Conclusion: The physical limits of Dennard scaling is galvanizing the community to put forward candidates for next generation computing, from bio to quantum computers. Photonics and in particular neuromorphic photonics are a formidable candidate for analog reconfigurable processing. We expect the development of this field to accelerate as neuroscience makes further leaps towards our understanding of the nature of cognition and artificial intelligence demands more computational resources for machine learning. As photonics technology matures and becomes more accessible to academic groups and small companies, we hope and expect this acceleration to continue.

References

- [1] P. A. Merolla et al. Science 345, 668 (2014).
- [2] M. Davies et al. IEEE Micro 38, 82 (2018).
- [3] B. J. Shastri et al. Nat. Photon. 15, 102 (2021).
- [4] M. A. Nahmias et al. IEEE J. Sel. Top. Quantum Electron. 26, 7701518 (2020).
- [5] T. Ferreira de Lima et al. J. Lightw. Technol. 37, 1515 (2019).
- [6] P. R. Prucnal and B. J. Shastri. Neuromorphic Photonics (CRC Press, 2017).
- [7] D. Brunner et al. Nat. Commun. 4, 1364 (2013).
- [8] K. Vandoorne et al. Nat. Commun. 5, 3541 (2014).
- [9] Y. Shen et al. Nat. Photon. 11, 441 (2017).
- [10] T. W. Hughes et al. Optica 5, 864 (2018).
- [11] A. N. Tait et al. J. Lightwave Technol. 32, 4029 (2014).
- [12] A. N. Tait et al. Sci. Rep. 7, 7430 (2017).
- [13] B. A. Marquez et al. J. Phys. Photonics, 3, 024006 (2021).
- [14] A. N. Tait et al. Phys. Rev. Appl. 11, 064043 (2019).
- [15] K. Berggren et al 2020 Nanotechnology 32, 012002 (2020).
- [16] T. Ferreira de Lima et al. Nanophotonics 9, 4055 (2020).
- [17] N. Farmakidis et al. Sci. Adv. 5, eaaw2687 (2019).
- [18] W. Zhang et al. Nat. Rev. Mater. 4, 150 (2019).
- [19] D Liang and J. E. Bowers, Nat. Photon. 4, 511 (2010).
- [20] S. Chen et al. Nat. Photon. 10, 307 (2016).
- [21] A.H. Atabaki et al. Nature 556, 349 (2018).