# Silicon Photonics Co-Design for Deep Learning

Qixiang Cheng, Member, IEEE, Jihye Kwon, Student Member, IEEE, Madeleine Glick, Senior Member, IEEE, Meisam Bahadori, Luca P. Carloni, Fellow, IEEE, and Keren Bergman, Fellow, IEEE

Abstract— Deep learning is revolutionizing many aspects of our society, addressing a wide variety of decision-making tasks from image classification to autonomous vehicle control. Matrix multiplication is an essential and computationally intensive step of deep learning calculations. The computational complexity of deep neural networks requires dedicated hardware accelerators for additional processing throughput and improved energy efficiency in order to enable scaling to the larger networks in upcoming applications. Silicon Photonics is a promising platform for hardware acceleration due to recent advances in CMOS compatible manufacturing capabilities, which enable efficient exploitation of the inherent parallelism of optics. This article provides a detailed description of recent implementations in the relatively new and promising platform of silicon photonics for deep learning. Opportunities for multiwavelength microring silicon photonic architectures co-designed with FPGA for pre- and post- processing are presented. The detailed analysis of a silicon photonic integrated circuit shows that a co-designed implementation based on the decomposition of large matrix vector multiplication into smaller instances and the use of nonnegative weights could significantly simplify the photonic implementation of the matrix multiplier and allow increased scalability. We conclude the paper by presenting an overview and a detailed analysis of design parameters. Insights for ways forward are explored.

*Index Terms*—silicon photonics; deep learning; neural network; photonic integrated circuit; microring resonator.

#### I. INTRODUCTION

DEEP learning is an extraordinarily popular machine learning technique that is revolutionizing many aspects of our society. Machine learning addresses a wide variety of decision-making tasks such as image classification [1], audio recognition [2], autonomous vehicle control [3], and cancer detection [4]. Matrix multiplication is an essential but time consuming operation in deep learning computations. It is the most time intensive step in both feedforward and backpropagation stages of deep neural networks (DNNs) during the training and inference, and dominates the computation time and energy for many workloads [1-3, 5, 6]. Deep learning uses models that are trained using large sets of data and neural networks with many layers. Since DNNs have high computational complexity, recent years have seen many efforts to go beyond general-purpose processors and towards dedicated

This work is supported by National Science Foundation (NSF) (grants CCF-1640108), and Semiconductor Research Corporation (SRC) (grant SRS 2016-EP-2693-A).

accelerators that provide superior processing throughput and improved energy efficiency.

It has been known for quite a while that matrix-vector multiplication can be performed by optical components taking advantage of the natural parallelism of optics to reduce computation time from  $O(N^2)$  to O(1) [7, 8]. Implementing these optical matrix-vector multipliers, however, has required the use of bulky inefficient optical devices. In the last several years the field of silicon photonics has made major progress to meet the massive needs of data centers and cloud computing. With silicon photonics, optical components and photonic integrated circuits are fabricated leveraging CMOS-compatible silicon manufacturing techniques to enable small-footprint, low-cost, power-efficient data transfers.

Optical matrix-vector multipliers (OMMs) based on silicon photonics represent a promising approach to address the challenge of compute-intensive multiplication in DNNs. An optimal solution must take into account the advantages and drawbacks of the silicon photonic technology along with the requirements of the application. Silicon photonics offers excellent co-design capabilities with off-chip control implemented by FPGAs to achieve accelerated computational gains. To analyze these capabilities in detail, we present the codesign of a DNN in conjunction with the OMM, developing an optical-electrical co-design infrastructure using FPGA control. The FPGA is used for (1) pre/post processing and (2) photonic device control. We identify opportunities for OMM architectures based on multiwavelength silicon microring resonators. We analyze and generalize the metrics of the microrings for linearity and reduced sensitivity to perturbations. The OMM can be used for time consuming, computationally expensive matrix multiplication. In the case of DNNs that are too large to be processed on a single optical chip, we explore methods to divide the computation, by using the parallelism at the system level to enable scaling to very large neural networks. In addition, we show how DNNs based on nonnegative weights significantly simplify the photonic implementation of the matrix multiplier and allow increased scalability.

The remainder of this paper is organized as follows. In Section II, we present a brief background of advances in deep learning and in silicon photonics. In Section III, we give an overview of and discuss trade-offs in the state-of-the-art

Qixiang Cheng, Madeleine Glick, Meisam Bahadori, and Keren Bergman are with Department of Electrical Engineering, Columbia University, New York, NY 10027, USA. (e-mail: qc2228@columbia.edu).

Jihye Kwon, and Luca P. Carloni are with the Department of Computer Science, Columbia University, New York, NY 10027, USA.

research in the implementation of silicon photonics for deep learning. Based on the analysis above, in Section IV, we propose a co-designed system for deep learning. We first present a detailed analysis of the design parameters and metrics for a silicon photonic integrated circuit (PIC) that implements an optical matrix multiplier. We generalize the role and characteristics of the silicon microrings, analyzing their limitations (including thermal sensitivities) in order to explore opportunities for optimized OMM structures. We then discuss system-level approaches towards electronic/photonic co-design for improved performance. At the end of this section, we provide insights on future directions and opportunities based on our analysis and the current state-of-the-art and application requirements. Section V concludes the paper.

## II. BACKGROUND

#### A. Deep learning

The fundamental concept of machine learning is that the core computation algorithm is not fully provided by a programmer, but automatically generated or improved by a computer system through experience [9]. The learning system explores a given class of computation models to determine the most suitable model among them based on the training data. One of the model classes that has gained widespread popularity is the DNN, which is the artificial neural network (ANN) with many layers in the network [10]. Inspired by the human brain, the concept of the ANN was first proposed in the 1940s [11]. More recently, with the increased volume of data, computing capability, and research interest, numerous ANNs have shown outstanding performance in machine learning tasks across various application domains [1-4]. Deep learning refers to machine learning using deep ANNs, also called DNNs. Two fundamental classes of ANNs are multilayer perceptrons (MLPs) and convolutional neural networks (CNNs).

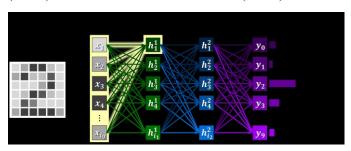


Fig. 1 A multilayer perceptron (MLP) for handwritten digit classification. The network consists of 4 layers (the input layers, two hidden layers, and the output layer) where each layer contains a number of nodes (also called neurons).

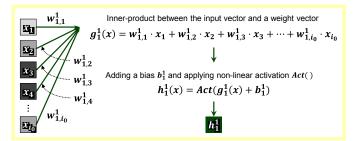


Fig. 2 The computation for a single node (the first node in the layer h1) in MLP.

MLPs, also known as fully-connected networks (FCNs), are the quintessential DNNs [10]. An MLP represents a function defined by a network consisting of multiple layers of nodes, which are also called neurons or perceptrons. For example, Fig. 1 shows an MLP for the task of recognizing a handwritten digit. The input image is represented as an array of pixel intensity values which are often normalized. The neural network behaves as a function that maps the input image to the probability score for each of the ten digits  $(0,1,2,\cdots,9)$ . Let  $i_0$  denote the number of pixels in the input image. Then, for an input array  $x \in \mathbb{R}^{i_0}$ , the neural network (shown in the box in Fig. 1) outputs  $y(x) \in \mathbb{R}^{10}$ , as follows. The input layer x contains  $i_0$  nodes  $x_1, x_2, \cdots, x_{i_0}$ . This layer is fully-connected to the first hidden layer  $h^1$ , which contains  $i_1$  nodes; each node  $h^1_k$   $(1 \le k \le i_1)$  is computed as

$$\begin{aligned} h_k^1\big(x_1,x_2\cdots,x_{i_0}\big) &= Act\big(g_k^1\big(x_1,x_2,\cdots,x_{i_0}\big) + b_k^1\big) & (1) \\ g_k^1\big(x_1,x_2,\cdots,x_{i_0}\big) &= \mathbf{w}_{k,1}^1\cdot x_1 + \mathbf{w}_{k,2}^1\cdot x_2 + \cdots + \mathbf{w}_{k,i_0}^1\cdot x_{i_0} & (2) \end{aligned}$$

where Act() denotes an element-wise nonlinear activation function (e.g., ReLU, sigmoid, softmax, tanh),  $b_k^1 \in \mathbb{R}$  is a bias, and  $w_{k,j}^1$  ( $1 \le j \le i_0$ ) represents the weight of the connection between node  $x_j$  and  $h_k^1$  (Fig. 2). Each hidden layer is fully-connected to the next layer, and the last layer in the network is the output layer containing ten nodes. The softmax function is often used for nonlinear activation of the output layer since it can be interpreted as a probability distribution.

The process of computing the output of a neural network as described above is called feedforward propagation. The information stored in the input layer propagates toward the output layer. How it propagates depends on the neural networks structure, weights, biases, and activation functions. During the training phase of supervised machine learning, given a large number of (input, output) instances, the values of weights and biases are updated through the gradient descent method, also called back-propagation [12]. Then, in the inference phase, a trained network is used to predict the output for a new input instance. With this approach, MLPs were among the first and most successful nonlinear learning algorithms [10]. The nonlinear activation plays a key role in ANNs. Without the ANN, the function expressed by an MLP is a composition of linear functions (which is linear). By inserting the nonlinear activation, such as ReLU or tanh, the resulting function becomes a composition of nonlinear functions, which can express much more complicated concepts.

In addition, the universal approximation theorem states that any continuous function defined on a compact set can be approximated by an MLP with a single hidden layer [13, 14]. Nevertheless, it does not address how many nodes are required in the hidden layer, or how to learn the weights and biases of such an MLP. Empirically, the accuracy of the trained networks improves as the number of nodes per layer increases, and as the number of layers increases. This motivated the advancement of DNNs.

CNNs were first proposed by LeCun et al. in 1989 for handwritten digit recognition [15], and they have outperformed many proposed MLPs, especially for more complex tasks such as colored image classification. Fig. 3 illustrates the overview of a CNN for image classification. The input image is stored across three channels, each representing the Red, Green, or Blue intensities. As shown in Fig. 4a, a convolutional layer (Layer L+1) usually contains multiple channels, and the values of nodes in each channel are computed using the information from all channels in the previous layer (Layer L). Fig. 4b gives a closer look at the connection between an input channel from the previous layer (Channel A in Layer L) and an output channel in a convolutional layer (Channel B in Layer L+1). A convolution kernel (of size  $3 \times 3$  in the example) dedicated to this connection defines how to obtain a value for each node in the output channel from a small neighbor (of size  $3 \times 3$ ) in the input channel (Fig. 4c). The kernel slides both vertically and horizontally on the input channel to cover all nodes in the channel, and the convolution result is propagated to the node in the associated position in the output channel. The amount by which the kernel slides is called the stride and this is often set to 1. Around the boundary of the input channel, additional nodes of the value zero can be padded before the convolution. When the kernel is of size  $R \times R$ , a padding of size  $\left|\frac{R}{2}\right|$  is commonly applied. At each node in the output channel, a bias and activation function are applied to the summation of the corresponding convolution results. To summarize, the value of a node  $Z_{c,d}^{L+1,B}$  at (c,d)-coordinate on channel B in layer L+1is computed as

$$\begin{split} z_{c,d}^{L+1,B}(z^{L,1},z^{L,2},\cdots z^{L,N_L};k^{L,1;L+1,B},k^{L,2;L+1,B},\cdots,k^{L,N_L;L+1,B}) \\ &= Act \left( \Sigma_{A=1}^{N_L} v_{c,d}^{L+1,B}(z^{L,A};k^{L,A;L+1,B}) + b_{c,d}^{L+1,B} \right) \\ & v_{c,d}^{L+1,B}(z^{L,A};k^{L,A;L+1,B}) \\ &= \Sigma_{\alpha=1}^{M} \Sigma_{\beta=1}^{M}(k_{\alpha,\beta}^{L,A;L+1,B} \cdot z_{c-\left\lfloor \frac{M}{2} \right\rfloor + \alpha,d-\left\lfloor \frac{M}{2} \right\rfloor + \beta}^{M}) \end{split} \tag{4}$$

where Act() is an activation function,  $b_{c,d}^{L+1,B} \in \mathbb{R}$  is a bias associated with this output node,  $z^{L,A}$  denotes channel A in the previous layer L,  $N_L$  represents the number of channels in layer L, and  $k^{L,A;L+1,B}$  refers to the convolution kernel of size  $M \times M$  defined for the connection between the channel A in layer L and the channel B in layer L+1.

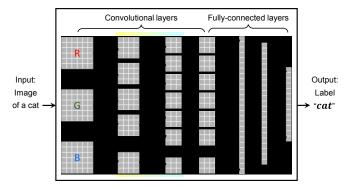


Fig. 3 A convolutional neural network (CNN) for image classification, consisting of convolutional layers followed by fully-connected layers.

Convolutional layers are elaborated in Fig. 4, and the computation for fully-connected layers is depicted in Fig. 2.

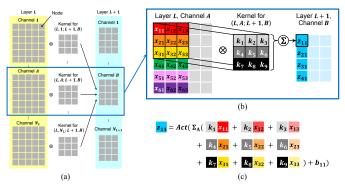


Fig. 4 Overview of convolutional layers. A convolutional layer consists of one or more channels where each channel contains a number of nodes. (a) Every channel in the previous layer is connected to each channel in the next layer. (b) The connection between one input channel (in Layer L) and one output channel (in Layer L+1). The convolution between a set (in the black square) of nodes in the input channel and the convolution kernel  $(k_1, \dots k_g)$  contributes to one node in the output channel. (c) The computation for a single node in the output channel.  $\Sigma_A$  denotes the summation over all input channels A.

Optionally, a convolutional layer may be followed by a pooling layer that reduces the size of the representation by pooling neighbors of  $R \times R$  nodes, where R often takes a small value such as 2,3,4 or 5. Most commonly used pooling functions are *maximum* (i.e., taking the maximum value from the  $R \times R$  neighborhood), average, median, and stochastic.

The network size and computational complexity of state-ofthe-art DNNs have generally increased over the decades. Meanwhile, much research has been also conducted on the accelerated and efficient computation of DNNs [16]. For both MLPs and CNNs, the core computation requirements during feedforward propagation are inner products of two vectors, or matrix-vector multiplications [5, 6, 17]. Both the weight product function g() for fully-connected layers (in Equation 2) and the convolution function v() for convolutional layers (in Equation 4) can be naturally translated into vector-vector or matrix-vector multiplications. GPUs have been extensively exploited to accelerate this type of computation, mainly leveraging their inherent feature of single-instruction multipledata parallelism [18, 19]. In addition, there has been growing interest in designing custom hardware accelerators and reconfiguring the DNNs for higher efficiency [20]. Haensch et al. have proposed in-memory analog computation for DNNs and have analyzed nonvolatile memory material candidates [5]. Amiri et al. have proposed a multi-precision CNN framework on an FPGA-CPU heterogeneous device [21].

#### B. Silicon photonics

While graphical processing units (GPUs), field-programmable gate arrays (FPGAs), and application-specific integrated circuits (ASICs) have received extensive interest for developing dedicated hardware accelerators in deep learning calculations [22-24], photonics has been long recognized as a promising alternative to address the fan-in and fan-out problems for linear algebra processors [25, 26]. A few unparalleled features motivate the exploration of a photonic

implementation. (1) The power consumption for data transfer that accounts for a large portion in electronic applicationspecific integrated circuits (ASICs) [16] can be greatly reduced by leveraging state-of-the-art optical transceivers. In addition, once a neural network is trained, the matrix configuration can be passive and optical signals can be processed with no additional power consumption [27]. (2) The operation bandwidth of such an OMM could potentially match that of the photodetection rate (typically in 100 GHz), which can be at least over an order of magnitude faster than the electronic system (typically restricted to the clock rate of a few GHz). (3) The OMM could have significantly lower latency, since the electronic hardware accelerators still rely on electronic transport that is bounded by the speed and power limits due to RC parasitic effects. Early demonstrations of photonic solutions were implemented with bulky free-space optics [25, 26], which required rigorous calibration for phase matching and have extreme scaling difficulties. Current photonic integration platforms provide opportunities for highly scalable solutions that improve energy efficiency and significantly reduce overhead of assembly, calibration, synchronization, and management [28].

Over the last two decades, silicon has been shown to be an excellent material platform for fabricating photonic devices, and processes have been developed to permit the reuse of CMOS manufacturing infrastructure to build complex PICs. It is, therefore, not surprising that silicon photonics is now widely as a key technology in next-generation communications systems and data interconnects [29]. On the one hand, following the example of the electronic fabless semiconductor industry, process design kit (PDK) libraries are being developed and standardization is being encouraged by the silicon photonics industry and users for broader accessibility [30, 31]. On the other hand, component customization is driven by a number of research groups and companies that design a large variety of specialized photonic components [32-34]. The ability to include increasing numbers of a wide range of optical components at the wafer scale has led to a powerful class of silicon-based PICs [35]. Such integration technology fundamentally improves circuit-level performance by reducing the complexity in assembly, calibration, and synchronization. As it matures, sustained increases in the functionality, performance, and reliability of circuits are enabled. This, in turn, stimulates new research directions leveraging the largescale photonic integration capabilities [27, 36-40]. Lightwave signals have been manipulated in their intensity and phase at the space, wavelength, polarization, and mode dimensions, for data transmission [33, 41], switching [42-44], and processing [27, 37, 40], in both digital [33, 41] and analog formats [27, 39, 40].

In addition, in recent years the ecosystem of silicon photonics has been extended to enable further functionality. The ability to add CMOS-compatible materials, such as Germanium (Ge), Ge-rich GeSi, and Silicon Nitride (SiN), to the Silicon-on-Insulator (SOI) platform has significantly enriched the component library and enhanced circuit-level performance. Notable examples include the Ge-on-Si photodiodes [45], high-

speed GeSi modulators [46], and the ultralow loss Si/SiN multilayer structure [47]. The development of heterogeneous integration [48, 49] as well as breakthroughs on direct growth of III-V quantum dot materials on silicon substrates [50] further complete the ecosystem, enabling a *System-on-Chip*.

The Mach-Zehnder interferometer (MZI) and the micro-ring resonator (MRR) are two of the most common functional building blocks in many photonic systems, such as modulators [32, 51, 52], filters [34, 53], multiplexers [54, 55], switches [56-58], and computing systems [27, 59, 60]. The MZI was first proposed over a century ago to determine the relative phase shift variations between two collimated beams derived by splitting light from a single source. Later work extended this concept to manipulate the probability of light arriving at either port, by precisely controlling the phase difference between the two arms [61]. Integrated MZIs generally consist of two 3 dB couplers with phase shifters embedded in each of the two arms. Detailed design considerations can be found in [52, 53, 62]. An MRR consists of an optical waveguide which is looped back on itself and coupled waveguides. Resonance occurs when the optical path length of the resonator is exactly a whole number of wavelengths and thus multiple resonances are supported. The spacing between these resonances is called the free spectral range (FSR). Similarly, a phase shifter can be embedded in the resonator to tune the optical path length in order to shift the resonance spectrum. The properties of MRRs are extensively described in the literature [63, 64], as well as their design considerations, performance metrics, and potential challenges [29, 32, 34, 54, 63, 64]. We discuss applications of the MRR in more detail below.

## III. SILICON PHOTONICS FOR DEEP LEARNING

This emerging area of research has been stimulated by recent results in which silicon photonics has been utilized to implement optical neural networks based on a spatial multiplexing technique with coherent interference [27], and a spectral multiplexing technique with wavelength filters [60]. In this section, we give a detailed overview of this recent progress in programmable silicon photonics for deep learning hardware accelerators.

# A. Linear MZI-based meshing optics with orthogonal spatial modes

Pioneered by the work of *Reck et al.* [65] showing that a mesh of 2×2 beam splitters and phase shifters in the form of a Mach-Zehnder interferometer (MZI) can be programmed to enable independent control of amplitude and phase of light for a set of optical channels, various novel architectures and design principles based on a cascade of MZIs have been proposed and demonstrated for both classical and quantum applications [27, 37, 39, 66-68]. These works are also referred to as "programmable linear optic processors" [69]. Phase shifters that are embedded in the arms of MZI units are used to control the interference of beams at the combining stage, while a pair of external phase shifters is employed in order to set a differential output phase. This allows the control of relative

amplitude and phase of the beams at each stage and thus the programing of the mesh. With specific interconnection patterns, universal linear optical components can be obtained [66-68, 70, 71].

Whereas most of the mesh networks are explored as universal linear optics for unitary operations [37, 65-67, 70, 71], Miller proposed a design method that implements arbitrary, nonunitary matrices, as shown in Fig. 5a [68]. This approach describes a self-configuring universal linear mesh that employs a set of orthogonal beams. The mathematics behind this design demonstrates that any linear optical device can be factorized using the singular value decomposition (SVD), as  $D = V\Sigma U^{\dagger}$ , where V and  $U^{\dagger}$  are unitary matrices and  $\Sigma$  is the diagonal matrix [68]. Theoretically, the universal unitary matrices of V and  $U^{\dagger}$  can be implemented following the designs proposed by Reck et al. [65] (Fig. 5b) and Clements et al. [66] (Fig. 5c), and the diagonal matrix  $\Sigma$  can be represented by an array of modulators that can set amplitude and phase [68], as illustrated by Fig. 5a. The unitary matrices of V and  $U^{\dagger}$  can be further decomposed to analytically define the values of beam splitters, i.e. phase settings of MZIs [65, 66].

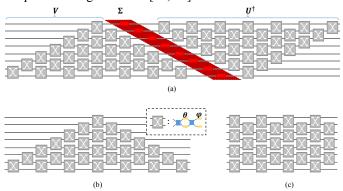


Fig. 5 (a) Universal linear mesh network comprising two unitary matrices and a diagonal matrix to set amplitude and phase, as proposed in [68]. Universal unitary matrix proposed by [65] in (b) and by [66] in (c).

The recent work by *Shen* and *Harris et al.* proposed a novel architecture (Fig. 6) for an optical neural network that offers hardware acceleration for deep learning applications [27]. Vectors were encoded in the intensity and phase of light and then fed into each layer of the network, which was comprised of an optical interference unit (OIU) and an optical nonlinearity unit (ONU). While the ONU function was emulated on a computer to act as a saturable absorber, the OIU was implemented using a silicon PIC to perform the optical matrix multiplications following Miller's design, which leverages the SVD decomposition [68]. This optical device consists of 56 programmable MZI units, each of which has two 50:50 power splitters and two pairs of phase shifters parameterized by  $(\theta, \phi)$ . The power splitters/combiners are realized by directional couplers and the  $\pi/2$  phase difference between the two outputs ensures the unitary property of its transformation. As a nonapplication-specific PIC, one matrix transformation requires two passes through the chip for: (1)  $V\Sigma$  and (2)  $U^{\dagger}$ . The required orthogonal beams are implemented by a set of coherent spatial modes. This device does not use on-chip detectors for

self-alignment. However, other generic approaches for setting up meshes can be leveraged to enable the calibration of phase disorders due to fabrication variations, such as the one described in [72]. In addition, the broadband nature of MZIs does not have a strong requirement for local phase stabilization, although on-chip thermal crosstalk could be a significant cause of phase errors.

Neural network training algorithms [73] can be leveraged to train the matrix parameters for different layers. Each layer contains a set of weights, which can be decomposed into phase settings and then programmed into the OIU. By implementing a two-layer optical neural network with 4 neurons per layer, a primitive task for vowel recognition was executed and achieved an accuracy of 76.7% [27]. Compared to the accuracy of 91.7% by execution with a conventional 64-bit digital computer, the key limiting factor for the accuracy of the optical neural network can be attributed to the computational resolution. The phase-encoding noise and the photo-detection noise are believed to be the primary factors causing reduced resolution [27]. This is also reflected in the fidelity analysis showing that the percentage error for each output of the SU(4) unitary matrix core is approximately 2.24% [27], which bounds the system's effective resolution. Suppressing on-chip thermal crosstalk, and lowering photo-detection noise would thus lead to a superior computational resolution of the network.

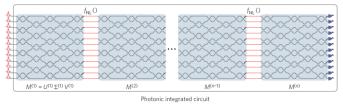


Fig. 6 All-optical architecture for integrated neural network [27].

The work described above shows an impressive example of applying silicon photonics to deep learning applications; yet, three factors in particular might bound the practicability of this approach. (1) Limited scalability of neurons. Let N denote the number of neurons. The optical depth (the number of MZI units traversed through the longest path) for the unitary matrix is given as 2N - 3 and as N in the scheme by Reck et al. [65] and by Clement et al. [66], respectively. This, therefore, leads to a total optical depth of 2N - 1 (with output reflected for a more compact layout [68]) and of 2N + 1, respectively, for the optical device that implements the arbitrary linear transformation using SVD encoding where the diagonal matrix  $\Sigma$  is implemented by an array of MZIs. Note that although the device using Reck et al. design has a slightly smaller optical depth, the Clement et al. layout is shown to be more tolerant to component loss in realistic interferometers, maintaining high fidelity [66]. The optical depth increases linearly with the number of neurons (N) by a factor of 2 which directly translates into additional loss in silicon photonics integrations. This additional loss could quickly outpace the optical power link budget and significantly deteriorate the system signal-to-noise ratio, thus limiting the computational resolution. (2) Error accumulation. Whereas the on-chip thermal crosstalk can be

suppressed, the finite encoding precision on phase settings will remain as the fundamental limitation for the optical neural networks with high computational complexity. The phase errors, in particular, accumulate when the lightwave signal traverses the MZI mesh with an optical depth of 2N + 1. In addition, such errors propagate through each layer of the network, which ultimately restricts the depth of the neural network. (3) Complex encoding scheme of matrix. The SVD method provides a perfect solution to decompose an arbitrary linear transformation. However, mapping the trained matrix parameters to the phase settings of the MZI mesh consumes additional computational power.

# B. Microring weight banks for spiking networks

Inspired by neuroscience in which biological neurons communicate by short pulses, spike processing, with this integrate-and-fire neuron model, has been proposed to exploit its massive parallelism potential in computation [74]. The cornerstone of the communication protocol is the spike coding scheme, which is digital in amplitude, and analog in pulse timing [75]. Input spikes from multiple sources are multiplied by a set of weight factors and temporally integrated to trigger a neuron firing a single output spike if the threshold is satisfied [76]. It has been recently recognized that photonics can be a powerful alternative to the microelectronic platform to implement such a spike processing system, given the significant advancement in both excitable lasers for the nonlinear processing, and analog PICs for the linear processing [77].

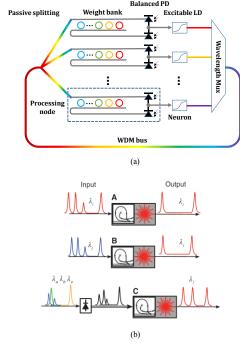


Fig.7 (a) Broadcast-and-weight spiking network proposed by [60]. (b) Classification of semiconductor excitable lasers [77].

An on-chip optical architecture, named broadcast-and-weight, was proposed by *Tait et al.* to implement scalable photonic spike processing networks to connect parallel neurons [60]. As illustrated in Fig. 7a, each spiking laser represents a neuron, and the optical neural network connects the output of

each neuron to multiple other neurons making use of wavelength division multiplexing (WDM). In contrast to the spatial multiplexing approach, channelization of the spectrum can somewhat simplify the interconnect network of neurons, as WDM channels can coexist in a single bus waveguide channel without interfering. The group of neurons that each utilizes a distinct wavelength share a common bus waveguide, as shown in Fig. 7a. The broadcast can be simply realized by passively splitting the bus waveguide to connect each of the neurons, enabling the all-to-all connection [60]. Each neuron is attached to a weight processing unit which is used to execute the linear transformation function for the N incident WDM signals that represent N neural nodes including itself. In this case, being capable of independently manipulating each weight is critical for creating differentiation among WDM channels. The silicon add-drop MRR is a natural choice due to its wavelength selective nature, as well as its cascadability, and continuous power-ratio-tunable feature [54]. The bank of cascaded MRRs, as an array of reconfigurable add-drop filters, imprint the weight coefficient to each corresponding channel. In a network of N neurons with N wavelength channels, each neuron incorporates a bank of N MRR filters, leading to a total number of  $N^2$  MRRs. The through port and drop port of the cascaded MRRs are respectively connected to create two subsets of weighted power, each connected to one of the balanced photodiode pair that performs the summation by incoherently aggregating the total incident optical power. The layout of the balanced photodiode subsequently enables subtraction between the two subsets of weighted powers for inhibitory weighting. The weighted sum is then used to excite a spiking laser neuron and three classifications of semiconductor excitable lasers are shown in Fig. 7b [77]. When the temporal integration of weighted pulses can push the gain above the lasing threshold, the neuron releases a spike. Otherwise, the system stays at rest.

As a key constituent element, the MRR weight bank has been carefully studied [78-81], since its scalability and tunability are closely tied to the performance limits of the optical neural network. Quantitative analysis was provided to measure the scaling of channel count, N, for an MRR filtering bank, illustrating the limiting factors of inter-channel crosstalk, insertion loss and more importantly, the bus length that causes coherent interactions between adjacent MRRs [78]. Similar to the MRR devices in data communication links, the interchannel crosstalk and cascading loss are the two fundamental constraints for system scale-up [54]. However, in contrast to the (de-)multiplexing-oriented designs that have only one common bus, the bus length becomes a key factor in the weight bank design that brings about multi-MRR coherent interactions due to the two bus configuration. This inevitably introduces another dimension of design complexity. Such inter-channel interference also deteriorates the independent control of the WDM channels, as the weights cannot be linearly separated. A more rigorous calibration process can be undertaken to improve these impairments in the WDM channels. Any power leakage or loss can be counter-balanced by adjusting the corresponded MRR coupling ratio. However, the degradation of the MRR

weight tuning range eventually becomes irreparable [78]. For a given system error  $\sigma$ , the tuning range is a critical factor that determines the network's computational resolution, as shown below

A few efforts have been made to optimize the device design and control plane for microring weight banks in silicon photonic integration platforms [79-81]. A continuous range of complementary (+/-) weighting has been demonstrated and recent work shows an effective weight setting accuracy of 5.1 bits [81], which is defined by  $log_2[(\mu_{max}-\mu_{min})/\sigma]$  where  $(\mu_{max}-\mu_{min})/\sigma$  $\mu_{\min}$ ) is the tuning range, and  $\sigma$  is the measured system error. The chip performance in this experiment is facilitated by photoconductive heaters which provide online feedback of photo-induced resistance to estimate the filter transmission. Considering that MRRs are particularly sensitive to thermal drift [82], the real-time feedback control loop, which tracks thermal fluctuations, including ambient temperature change, self-heating effects, and thermal crosstalk, plays a major role in such a multi-resonator system. It, therefore, provides superior performance compared to the feedforward control scheme [79, 80], which relies on fixed pre-built references.

Whereas a set of MRRs sandwiched by two buses that drop power into a balanced photodetector offer complementary (+/-) weight factors, the closed WDM link makes it difficult to monitor the isolated transmission state for each wavelength channel. Altering the weight factor via shifting the resonance spectrum of individual MRR unit arranged in a cascading scheme would significantly constrain its tuning range, given that all channelized MRR filters coexisting on the same bus have to tightly fit within one FSR. The embedded photoconductive heaters within MRRs, provide a limited but adequate solution for neuromorphic applications, [77, 81]. However, the adoption of photoconductive effects in the analog computing system may not sufficiently deliver the requirements for optical matrix multiplication with higher resolutions.

### C. Discussion

Both of the aforementioned approaches aim at processing an entire ANN application or an entire matrix-vector multiplication on a single optical device. Whereas those approaches may have advantages in the processing speed, the capability of the optical device strictly limits the size of the ANN to be processed. For instance, the optical neural network architecture proposed by Shen et al. consists of two layers, each with four neurons, for a primitive machine learning task of classifying four vowels in speech [27]. However, many machine learning tasks in practice involve learning more complex functions that take in a large number of inputs. For a handwritten digit recognition task, the number of input neurons are  $28 \times 28 = 784$ , one for each pixel of the input image, and the number of output neurons are 10, which equals the number of candidate digits [15]. For breast cancer detection, an MLP with 30 input neurons, 500 neurons in each of the three hidden layers, and 2 output neurons was used to achieve the detection accuracy of 99% [83]. The computation for this MLP includes the multiplication between a matrix of size  $500 \times 500$  and a vector of dimension 500. It is not feasible or practical to fully

optically implement such large neural networks or matrixvector multiplications using the above approaches due to their limited scalability.

## IV. SILICON PHOTONICS CO-DESIGN FOR DEEP LEARNING

Co-design of silicon photonic and electronic circuits provides new opportunities for efficient computation of deep learning. Silicon photonics has the potential for high-speed analog matrix multiplication. However, the computational requirement for ultra-large DNNs with high accuracy demand may exceed the capability of a single PIC, for high-complexity computing tasks. Our co-design approach, described in this section, explores practical and scalable solutions to process such large neural networks while employing feasible optical devices.

Figure 8 illustrates an overview of the proposed co-designed system with the electronic circuitry that processes DNNs at the system level, and a PIC that performs optical matrix-vector multiplication of fixed-size inputs. The PIC takes in a matrix K of size  $M \times N$  and a vector x of size N, and outputs a vector y of size M such that  $y = K \cdot x$ .

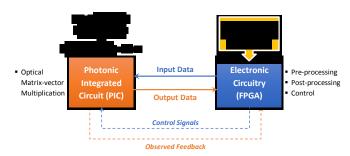


Fig. 8 Overview of the proposed co-designed system for deep learning.

## A. Silicon photonic OMM

Channelization in the wavelength domain avoids the phase-sensitive designs that require the control of relative phases from different nodes for coherent interference effects. Therefore, the wavelength multiplexing technique provides an elegant solution to address the many-to-one coupling (fan-in), which is a typical problem in neuron networks. Combined with tunable add-drop MRR technology, the direct mapping from the weight matrix to the power coupling ratio of wavelength filters can also eliminate the complex encoding phase. Such simplicity would further boost the validity of optical neural networks as hardware accelerators for deep learning applications.

The high thermo-optic coefficient of silicon (1.8×10<sup>-4</sup> K<sup>-1</sup>) creates a double-edged sword for silicon MRR elements. While it allows effective manipulation of light by the thermo-optic effect, due to its narrow-band nature, this thermal susceptibility can be detrimental to device performance. Therefore, accurate monitoring and control mechanisms are normally required. A number of energy-efficient yet precise locking schemes have been demonstrated [84-88], for data communications. The work by *Tait et al.* sheds some light on the feedback control of an analog MRR system, which relies on an estimate of filter transmission [81]. The plasma dispersion effect via either

carrier depletion or injection can be leveraged to provide nanosecond-scale tuning mechanism [29]. However, fabrication variations [89], in addition to the self-heating effect and ambient temperature change [82], most often require an additional thermo-optic phase shifter. The electro-optic tuning mechanism also requires attention to the induced electro-absorption loss that compromises the extinction ratio of resonance, thus the resolution of computation. This additional loss disturbs the balance between the coupling power and the round trip loss in the ring cavity from the critical coupling point. The operation condition for critical coupling is discussed in the following subsection.

In general, the on-chip thermal crosstalk is a primary culprit of the system instability for MRR-based silicon photonic circuits. Hereby, we start with an analytical model of add-drop MRRs to provide an insight into constraints on weight resolution due to thermal crosstalk. We focus on the thermoptic phase shifting effect since it is a lossless tuning mechanism but imposes the most thermal impairments. We define the MRR weight sensitivity and discuss an approach to increase system stability. The ability to utilize only nonnegative values for training weight factors opens new opportunities to refine the ring locking scheme in the analog domain. A new class of highly-accurate, yet scalable OMMs that are based on add-drop MRRs for deep learning can thus be obtained.

## 1) Thermal-crosstalk restricted weight resolution

# 1.1) Weight definition and sensitivity of add-drop MRR

An add-drop ring resonator refers to a circular ring structure that couples to two straight waveguides, as schematically shown by the inset in Fig. 9a. The optical transfer function of the drop and through port can be expressed as [64]:

$$D(\phi) = \left| \frac{-\kappa_1 \kappa_2 L^{0.25} \exp(-j\phi/2)}{1 - t_1 t_2 \sqrt{L} \exp(-j\phi)} \right|^2 \tag{5}$$

and

$$T(\phi) = \left| \frac{t_1 - t_2 \sqrt{L} \exp(-j\phi)}{1 - t_1 t_2 \sqrt{L} \exp(-j\phi)} \right|^2$$
 (6)

where  $t_1$  and  $\kappa_1$ ,  $t_2$  and  $\kappa_2$  are the self-coupling and cross-coupling coefficient for the input and drop coupling region, respectively. L is the round-trip optical power attenuation of the ring. We assume  $t^2 + \kappa^2 = 1$  which allows the loss introduced by the couplers to be included in L.  $\phi$  is the relative optical phase shift inside the ring:

$$\phi = \frac{(\lambda - \lambda_{res})}{FSR} \times 2\pi \tag{7}$$

where  $\lambda_{res}$  is the ring resonance wavelength, and *FSR* is the free spectral range of the resonance spectrum. We define the weighting function,  $\mu$ , using the through port of the add-drop MRR, considering the negligible through loss and its flexibility in cascading. Equation 6 can be rewritten as:

$$T(\phi) = \frac{T_0 + \left(\frac{2F}{\pi}\sin\left(\frac{\phi}{2}\right)\right)^2}{1 + \left(\frac{2F}{\pi}\sin\left(\frac{\phi}{2}\right)\right)^2} \tag{8}$$

where

$$T_0 = \frac{(t_1 - t_2 \sqrt{L})^2}{(1 - t_1 t_2 \sqrt{L})^2} \tag{9}$$

and F is the finesse of the ring, given by

$$F = \pi \frac{\sqrt{t_1 t_2 \sqrt{L}}}{1 - t_1 t_2 \sqrt{L}} \approx \frac{FSR}{\Delta \lambda_{3dB}}$$
 (10)

where  $\Delta \lambda_{3dB}$  is the optical bandwidth of microring. Note that the approximation in Equation 10 holds only when F >> 1. Under critical coupling, the coupled power is equal to the power loss in the ring cavity, i.e. satisfying the relation  $t_1 = t_2 \sqrt{L}$ , hence the transmission drops to zero,  $T_0 = 0$ . Therefore, the design to operate at critical coupling mode enables the maximum extinction ratio (i.e. dynamic range) for the power transfer at the through port. The weighting function,  $\mu$ , can thus be given as:

$$\mu(\phi, F) = \frac{T(\phi)}{T(\pi)} = \frac{T_0 + \left(\frac{2F}{\pi}\sin\left(\frac{\phi}{2}\right)\right)^2}{T_0 + \left(\frac{2F}{\pi}\right)^2} \times \frac{1 + \left(\frac{2F}{\pi}\right)^2}{1 + \left(\frac{2F}{\pi}\sin\left(\frac{\phi}{2}\right)\right)^2}. \tag{11}$$

$$\frac{1}{0.8}$$

$$\frac{1}{0.8}$$

$$\frac{1}{0.9}$$

Fig. 9 (a) Weight definition of the through port of an add-drop MRR for critical coupling condition with various finesse. (b) Sensitivity of the weights for different finesse. (c) Optical phase shift at maximum sensitivity point. (d) Maximum sensitivity of weights as function of finesse. A linear trend is observed.

We define F as a variable of  $\mu$  since finesse stands out as a key parameter of both the ring sensitivity and the scalability of number of neurons [77]. The finesse is a measure of the sharpness of resonances relative to their spacing (FSR) and represents, within a factor of  $2\pi$ , the number of round-trips made by light in the ring before its energy is reduced to 1/e of its initial value [63]. Therefore, from this point of view, the round trip loss, L, as well as the coupling coefficients in the coupling regions of the ring,  $t_1$  and  $t_2$ , are loss factors that can be manipulated to alter F, which is also reflected in Equation 10. For datacom applications, MRRs are generally designed with radii in the region of 5-10 um to avoid undesired high bending losses (i.e., radiation and scattering) while maintaining reasonably large FSR, therefore limiting the finesse to the order of tens [64]. Special designs, however, can lead to finesse with values of a few hundreds [90, 91]. Further details are discussed in Section IV.C.1. In Fig. 9a we plot the weight factor as a

function of  $\phi$  for F values in the range of 10 to 100 (10, 20, 50, and 100), assuming the critical coupling operation ( $T_0 = 0$ ).

It is not a surprise to see that a sharper resonance, i.e. larger F, gives rise to a more abrupt change in the weight as a function of  $\phi$ . It has been shown that the thermo-optic response (i.e. optical phase shift,  $\Delta \phi$ ) of the microring is a linear function of heating power ( $\Delta P$ ) [82]. However, for a thermal perturbation ( $\Delta P$ ),  $\Delta \mu$  varies depending on the weight  $\mu$ , due to the nonlinear behavior of the optical transfer function. We thus define the sensitivity of the weights as the slope of the weight:

$$\frac{\Delta\mu}{\Delta\phi} \approx \frac{\partial\mu}{\partial\phi} = \frac{(1+a^2)(1-T_0)}{T_0+a^2} \frac{0.5a^2\sin\phi}{\left(1+\left(a\sin\frac{\phi}{2}\right)^2\right)^2}$$
(12)

where  $a=2F/\pi$ . Combining with Equation 11, we can plot the sensitivity as a function of weight for various values of F as in Fig. 9b for critical coupling ( $T_0=0$ ). One can see that a lower sensitivity exists in the weighting function with a smaller F, where the change in weight is milder over  $\Delta \phi$ . This can be understood by realizing that a lower F results in a wider resonance linewidth, hence the weight has a smaller gradient as seen in Fig. 9a. The optical phase settings at the maximum sensitivity (i.e.  $\partial^2 \mu / \partial \phi^2 = 0$ ) as a function of F is further given by:

$$\phi_{\text{max}} = 2 \tan^{-1} \sqrt{\frac{\frac{3\alpha^2 + 2 - \sqrt{9\alpha^4 + 4\alpha^2 + 4}}{\alpha^2 - 2 + \sqrt{9\alpha^4 + 4\alpha^2 + 4}}}$$
(13)

and illustrated in Fig. 9c, indicating the weight variations are most sensitive close to the resonance, which agrees with the trend illustrated by Fig. 9a due to the nonlinear power transfer in MRRs. Figure 9d indicates that the maximum sensitivity of the weight has a linear dependence on the finesse of the MRR, again showing that larger finesse leads to worse sensitivity. To facilitate the quantitative analysis on the bounded effective resolution, we use a first order Taylor expansion of  $\partial \mu/\partial \phi$  assuming that  $a^2 \gg 1$  (see Appendix II) to show this. The result is:

$$\left|\frac{\partial \mu}{\partial \phi}\right|_{\text{max}} \approx \frac{9}{16\sqrt{3}} \ \alpha = \frac{3\sqrt{3}}{8\pi} \ F = 0.2067 \ F.$$
 (14)

#### 1.2) Thermal crosstalk induced weight error

Thermal crosstalk occurs due to the proximity of rings to each other. The linear dependence of the temperature changes on the heater power results in a linear perturbation relation of the ring's temperature:

$$\Delta T_i^{xtalk} = \sum_{j=1, j \neq i}^{N} \chi_j P_{H,j} (\mu_j)$$
 (15)

where  $P_H$  is the heating power of other rings for setting their corresponding weights. This change of temperature translates into a change in the optical phase inside the ring:

$$|\Delta\phi| = |\Delta\lambda_{\rm res}| \frac{2\pi}{FSR} \approx 0.07 \times |\Delta T^{xtalk}| \times \frac{2\pi}{FSR}.$$
 (16)

We use 0.07 nm/K as the typical resonance thermal sensitivity of silicon microrings [82]. Thermal crosstalk can be considered as a biased (deterministic) perturbation; hence, it affects the average value of the error,  $|\overline{\Delta\mu}|$ . Since the optical phase shift due to thermal crosstalk is a direct consequence of the weight

of other rings whereas the weight sensitivity is dependent on the weight of interest, these two factors are uncorrelated and both can simultaneously occur at their worst cases. Therefore, the maximum weight error due to thermal crosstalk can be written as:

$$\max |\overline{\Delta\mu}| = \left|\frac{\partial\mu}{\partial\phi}\right|_{\max} \times |\Delta\phi|_{\max} = 0.091 \times \frac{|\Delta T^{xtalk}|_{\max}}{\Delta\lambda_{3dB}}.$$
 (17)

Considering adjacent MRR elements as thermal crosstalk sources and that the maximum phase shift inside each adjacent ring is  $\pi$ , the maximum temperature change due to thermal crosstalk from an adjacent ring can be given as:

$$\Delta T_{max}^{xtalk} \approx 7.143 \, FSR \times \alpha_T$$
 (18)

where  $\alpha_T$  is the fraction of the thermal energy from adjacent rings. The weight error then aggregates as:

$$\max|\overline{\Delta\mu}| = 0.65 F \times \sum_{i} \alpha_{T,i}. \tag{19}$$

The solution of heat diffusion in 2D space of the chip has a form of [82]:

$$q(r) = q_{ring} \times \frac{\ln(\frac{R_{\infty}}{r})}{\ln(\frac{R_{\infty}}{R})}$$
 (20)

where q is thermal energy density (proportional to the change in temperature at each location), r is the distance to the crosstalk source, R is the radius of the ring, and  $R_{\infty}$  can be viewed as the boundary of the chip, as shown in Fig. 10a. Figure 10b shows the validation of this analytic equation with COMSOL simulation [82, 92] for  $R = 10 \,\mu m$  and  $R_{\infty} = 1 \,mm$ . As expected, the heat density decreases at farther distances from the MRR's heater, but the 2D heat diffusion shows a rather strong thermal crosstalk impact (e.g. 50% at 100 µm proximity). Note that in an actual photonic chip the heaters are most commonly located on top of the MRR so that the heat can also diffuse vertically. Since the thickness of the heater,  $t_H$ , is typically much smaller than the footprint of the heater ( $\approx 100$ nm [93]), most of the heat generated by the heater diffuses vertically (out of plane) instead of horizontally (in-plane). Therefore, the fractional in-plane heat crosstalk from one ring to another can be estimated by:

$$\alpha_T \approx \frac{t_H}{2R} \times \frac{\ln(\frac{R_{\infty}}{r})}{\ln(\frac{R_{\infty}}{r})}$$
 (21)

and thus:

$$\max|\overline{\Delta\mu}| = 0.65 F \times \frac{t_H}{2R} \times \sum_{i} \frac{\ln\left(\frac{R_{\infty}}{r_i}\right)}{\ln\left(\frac{R_{\infty}}{r_i}\right)}.$$
 (22)

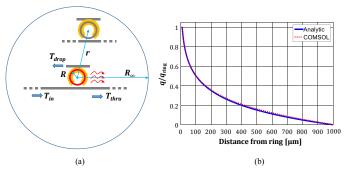


Fig. 10 (a) Schematic of thermal crosstalk between adjacent MRRs.  $R_{\infty}$  denotes the boundary of the chip. Thermal crosstalk arises from in-plane diffusion of heat and gets worse at closer proximity. (b) Comparison of analytic 2D equation

for heat diffusion with finite element results in COMSOL. The logarithmic behavior for the heat diffusion is confirmed. Note that the heat density, q, is proportional to the temperature change,  $\Delta T$ , and can be considered a measure of thermal crosstalk.

## 1.3) Weight resolution

The resolution determines the minimum possible steps for setting weights with the highest certainty. If  $\mu$  is the calibrated weight in the ideal case and  $\hat{\mu}$  is the weight in the presence of perturbations, we can write:

$$\hat{\mu} = \mu + \Delta \mu(t) = \mu + \overline{\Delta \mu} + \delta \mu(t) \tag{23}$$

where  $\Delta\mu(t)$  is the error of the weight. This error can be decomposed into a stationary (deterministic) average denoted by  $\overline{\Delta\mu}$  and a random noise like term denoted by  $\delta\mu(t)$ . We consider the resolution is set by the maximum root mean square error given by  $\max|\Delta\mu(t)|=\max|\overline{\Delta\mu}|+\sigma_{\mu}/2$  where  $\sigma_{\mu}^2=\overline{\delta\mu^2(t)}$  is the standard deviation of the noise-like error. The resolution is then written as:

Resolution = 
$$\frac{1}{\max|\Delta\mu(t)|} = \frac{1}{\max|\overline{\Delta\mu}| + \frac{\sigma\mu}{2}}$$
. (24)

In such a system, it is reasonable to assume the thermal crosstalk induced error (i.e.  $\overline{\Delta\mu}$ ) is dominant over the photodiode noise,  $\sigma_{\mu}$ . For an MRR element in an array, its two adjacent rings are considered as the dominant sources of thermal crosstalk. Therefore, referring to Equation 22, we can plot the contours of effective bit resolution for an MRR unit as a function of both the unit pitch and its finesse for R=10~um, and  $R_{\infty}=1~mm$ , as shown in Fig. 11.

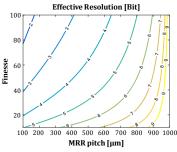


Fig. 11 Contours of effective bit resolution for the weight of MRR due to thermal crosstalk as function of finesse and proximity.

Note that this model is more accurate for small thermal perturbations; however, the combination of Equations 22 and 24 still serves as a qualitative analysis on how the pitch size of MRR weighting elements and their finesse bound the effective resolution, even when the thermal crosstalk is strong. Feedforward calibration can somewhat alleviate the thermal crosstalk restrictions, yet the calibrated system accuracy heavily depends on the weight settings of adjacent MRR units. A scalable OMM with the capability of high resolution thus calls for a new design approach and the capabilities of computing using only nonnegative weight factors open up a new design philosophy, as discussed in the following subsection.

#### 2) Hitless weight-and-aggregation architecture

We propose a co-designed architecture for optical matrix multipliers which are specially customized for highly-accurate, scalable and nonnegative weight matrices. The hitless weightand-aggregation design essentially describes an interconnect architecture that allows computational nodes (neurons) to carry arbitrary input vectors and to be independently weighted and summed. Such a many-to-one network is formed on the basis of channelization of the spectrum, creating physical and logical connections between input and output vectors. We put forward a hitless weighting structure by employing the colored channels in parallel rather than cascading them. This design isolates each weight on each connection and makes the tuning of MRR filters truly independent, i.e. not interfering with other channels. Such a hitless design also decouples the weighting and summation functions by allocating dedicated functional blocks, both of which employ MRR units, thus allowing independent optimization to decouple the constraint between the scalability of neurons and the weight sensitivity. The nonnegative weights are defined using the optical transfer function of the MRR through port, while the drop port is used as a monitoring outlet to provide real-time feedback for the weight control loop.

An M×N OMM consisting of M N×1 vector multipliers is illustrated in Fig. 12. Distinct continuous-wave (CW) wavelengths (representing N neurons) can be implemented by either M sets of N wavelength-multiplexed laser arrays [94] or optical frequency comb lasers [33, 95], or one set of lasers passively split into M copies. The nonnegative weight factors obtained from the trained matrix parameters are mapped to the coupling ratios and imprinted to the CW signals using multiring weighting blocks. The colored signals that carry the same set of weights are routed to all outputs. The N input vectors are then formed by a set of intensity modulators to the fanned-in WDM signals, before combining to form the M output vectors. The aggregation will be performed by another dedicated set of N high-finesse MRRs that are critically coupled to the WDM bus waveguide. The wavelength-multiplexed data streamed into the bus are optically summed by a photodiode, in which the photocurrent represents the total optical power. The M output vectors are then sent for nonlinear processing. Design considerations for each functional block are detailed in the following subsections.

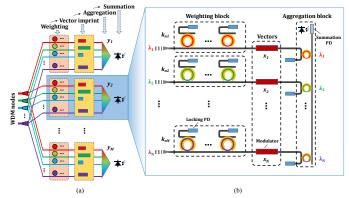


Fig. 12 (a) Hitless weight-and-aggregation architecture for  $M \times N$  vector-matrix multiplier. (b) One unit out of M for  $N \times 1$  vector to be multiplies by  $1 \times N$  matrix.

#### 2.1) Hitless architecture for nonnegative weight factors

The design philosophy for the MRR-based weighting block

is different from the conventional approach, in which tuning a filter in a link where WDM signals coexist controls the power coupling of the desired wavelength. The drop spectrum of such a MRR filter also sees other channels on the bus and thus the inevitably interferes adjacent channels. interference not only limits the weight tuning range but also acts as an unbiased perturbation to the weight that bounds the resolution for the nonnegative OMM system. Thus, a large channel spacing is required which trades off the system scalability.

Instead of utilizing the cascading layout of MRRs, the hitless design exploits a parallel arrangement of the weighting filters, shown in Fig. 12. This strategy stabilizes the weighting block within each wavelength branch before multiplexing onto the WDM bus, ensuring full tuning independence. Therefore, the design considerations for the MRR weighting filters can be narrowed down to a sole factor, i.e. sensitivity. As defined by Equation 8 and Fig. 9b, a small finesse is favored. Note that a trade-off exists since higher optical phase change is required to set the MRR to a specific weight for smaller finesse, which translates into higher heating power and, in turn, makes the thermal crosstalk worse. However, Equation 9 still provides the worst possible scenario for the thermal crosstalk effects.

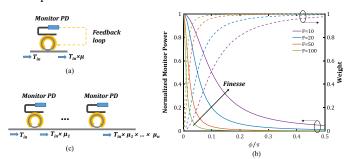


Fig. 13 (a) Single MRR weighting element with monitor PD. (b) The normalized monitor power as well as corresponded weight factor as a function of  $\phi$ . (c) Multi-MRR weighting element.

While the filter-through port is used to define the weighting function, the drop port connects to a monitor photodiode (PD), shown in Fig. 13a, providing a highly accurate feedback control loop for precise ring power locking. Figure 13b plots the normalized monitor power for this structure as a function of  $\phi$ , together with the corresponded weight factors. The locking accuracy could be compromised at power levels approaching zero (weight factors approaching one), given the existence of photodiode shot noise. To obtain a more linear transmission response, the ring spectrum tail can be omitted at the sacrifice of a slightly reduced weighting range.

The precise locking scheme would require a calibrated process, which sets up a look-up table (LUT) that maps the weight factor to the monitored optical power for each filter. By periodically polling the power monitor and comparing to the LUT, the locking scheme can effectively offset thermal perturbations, including on-chip thermal crosstalk, and ambient temperature fluctuations. The locking accuracy, which could translate into weight resolution, can be limited by the PD shot noise, the finite precision that offers by the DAC/ADC, as well as the polling and locking rate.

# 2.2) Multi-ring weighting block for reduced sensitivity

By utilizing multiple MRR filters as illustrated in Fig. 13c, the weight sensitivity can be further relaxed. The overall weighting function,  $\mu_o$ , for *n* cascaded ring filters can be given

$$\mu_0 = \mu_1 \cdot \mu_2 \cdot \dots \cdot \mu_n. \tag{25}$$

 $\mu_o = \mu_1 \cdot \mu_2 \cdot \dots \cdot \mu_n. \tag{25}$  For simplicity, we assume  $\mu_1 = \mu_2 = \dots = \mu_n = \mu$ , in which case  $\mu$  is given by Equation 11 and the phase settings are the same for all MRR filters. Figure 14a plots  $\mu_0$  as a function of  $\phi$ , for n=1, 2, 3, 4, 5, with F=10. It can be seen that the weighting function gets increasingly linear as *n* increases.

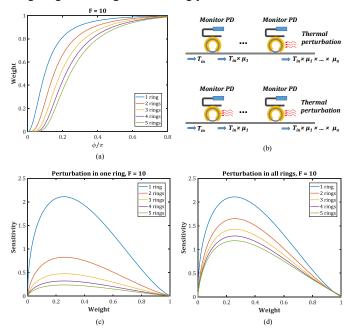


Fig. 14 (a) Weight as a function of  $\phi$ , for n=1, 2, 3, 4, 5, with F=10. (b) Illustration for one perturbed ring and all perturbed rings in a multi-ring weighting block. Weight sensitivity as a function of  $\phi$  with thermal perturbation in (c) one ring and (d) all ring simultaneously.

We can analyze two cases for the weight sensitivity of the multi-ring system: 1) One ring is perturbed thermally, and 2) all rings are perturbed thermally at the same time, shown by Fig. 14b. When the OMM setting leads to one or multiple heat sources on a chip, the dominant thermal effect is considered to be from adjacent rings. It is thus reasonable to take the one perturbed ring as the lower boundary for weight sensitivity. We have:

$$\frac{\partial \mu_0}{\partial \phi} = \frac{\partial \mu_1}{\partial \phi} \cdot \mu_2 \cdot \dots \cdot \mu_n. \tag{26}$$

This can be readily solved by referring to Equation 12. The weight sensitivity with thermal perturbation in one ring can thus be plotted and is shown in Fig. 14c, in which the single ring case is included for direct comparison. It can be seen that the two-ring system suppresses the weight sensitivity significantly, but the trend continues with a decreasing decrement when the number of rings increase. For the case that thermal perturbation occurs in all rings, we can have:

$$\frac{\partial \mu_o}{\partial \phi} = \sum_{i=1}^n \frac{\mu_o}{\mu_i} \frac{\partial \mu_i}{\partial \phi} = \frac{\partial (\mu^n)}{\partial \phi}.$$
 (27)

Figure 14d plots this case representing the upper boundary for

weight sensitivity. It can be seen that the system still gains tolerance to thermal perturbation compared to the single ring case. Considering the additional cost, footprint, and complexity introduced by the multi-ring system, the lower number two is preferred. Therefore, for the implementation of an M×N vectormatrix multiplier, the total number of MRRs is 3M·N including both weighting MRRs and aggregation MRRs. The total number of PDs is M.

Although the multi-ring system exhibits lower weight sensitivity, overcoming the limitation of the finite precision for the DAC with which an optical phase can be set is still a challenge. In an n-ring weighting block, the minimum step in the weight,  $\delta\mu$ , bounded by the DAC resolution for setting the optical phase of each ring yields a weight  $\hat{\mu} = \mu \pm \delta\mu$ ; hence the overall weight is  $\hat{\mu}_o = (\mu \pm \delta\mu)^n \approx \mu^n (1 \pm n \, \delta\mu/\mu)$ . Therefore, the error given by  $n \, \mu^{n-1} \delta\mu$  can be at its worst (i.e.  $n \, \delta\mu$ ) when  $\mu$  is close to 1. A smaller error than  $\delta\mu$  is achieved only for weights for which  $n \, \mu^{n-1} < 1$ . For a two-ring weight block, the worst error is  $2\delta\mu$  which can occur for any weight.

## 2.3) Aggregation and summation

In contrast to the MZI-based OIU for matrix multiplication where the input vectors are imprinted before feeding into the OIU [27], we process the vector imprint after the weighting stage. This is because the weight factor, i.e. coupling ratio, is locked by the dropped power as illustrated in Fig. 13, and the streamed input vectors with power fluctuations would deteriorate the locking accuracy. Therefore, the proposed processing flow as shown in Fig. 12 resolves this issue. The input vectors are imprinted via high-speed intensity modulators [96]. A linear intensity modulator, such as the Mach-Zehnder modulator, is favored [52]. As we analyzed in the following subsection, high computation accuracy can be obtained when the input vectors have the same resolution as the weights.

The weighted input vectors can subsequently be aggregated into the WDM bus through dedicated ring filters. As shown by Fig. 15a, the locking scheme for the aggregation MRRs operates differently, where the through power is always locked at the minimum state for a total power drop. This non-tunable feature ensures the maximal spectral efficiency regarding the number of wavelengths that can reside in the WDM bus.

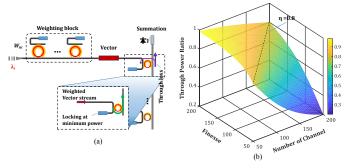


Fig. 15 (a) Operating principle of the aggregation MRRs. (b) Through power ratio as a function of both finesse and number of channels.

Since the aggregation ring filters act only as wavelength multiplexers, a large finesse is favored in order to achieve high scalability in the number of wavelength channels, i.e. number of neurons. For a given finesse, the number of channels that can be carried within one FSR is determined by the channel spacing. A trade-off exists for the channel spacing as it also determines the inter-channel crosstalk when the dropped signals pass through neighboring rings towards the summation PD on the bus. This leads to a through loss as illustrated by Fig. 15a. We can rewrite Equation 7 as:

$$\phi = \frac{(\lambda - \lambda_{res})}{FSR} \times 2\pi = \frac{2\pi}{N_{\lambda}}$$
 (28)

where  $(\lambda - \lambda_{res})$  and  $N_{\lambda}$  are the channel spacing and number of channels, respectively. A large portion of the power loss gets dropped to the locking PD. This however does not compromise the weighting resolution. If we limit the through power ratio to  $\eta$ , we have:

$$N_{\lambda} = 2\pi/\cos^{-1}\left(1 - \frac{2}{\frac{4F^2}{\pi^2}} \frac{T_0\left(1 + \frac{4F^2}{\pi^2}\right) - \left(T_0 + \frac{4F^2}{\pi^2}\right)\eta}{\left(T_0 + \frac{4F^2}{\pi^2}\right)\eta - \left(1 + \frac{4F^2}{\pi^2}\right)}\right). \tag{29}$$

We can then plot a 2D contour for  $\eta$  as a function of both finesse and number of channels, as shown in Fig. 15b. Here, we assume the induced loss is dominated by the adjacent channel. It can be seen that for  $\eta$ =0.8, which translates into ~1 dB through loss,  $N_\lambda \approx F$ . It should be noted that the insertion loss for all wavelength channels should be equalized by adjusting the individual input power, in order to allow each neuron to have the same maximum weight at summation. In addition, due to the multi-ring weighting block, the system can achieve higher order crosstalk suppression for the "0" weight.

## B. System-level co-design

In order to take full advantage of both the optical speed-up and electronic manipulation of the parallelism and memory, interactions between the two technologies require careful attention, especially when one processes digital signals and the other analog signals. We identify the system-level challenges for the co-design as following: (1) Computation breakdown to match the interface. Processing a DNN may require matrixvector multiplications for ultra-large matrices and vectors. The electronic circuitry should preprocess the DNN, breakdown the computation to smaller matrix-vector multiplication instances, send the request to a silicon photonic circuit, and post-process the results. (2) Minimization of the number of updates for the input matrix to the OMM. For each instance of matrix-vector multiplication requests, changing the values represented by the OMM microrings introduces a nonnegligible delay. Thus, to make the most of the high capacity of optical interconnects, it is desirable to have the elements of the input matrix to the OMM constant over a sequence of matrix-vector multiplication requests sent from the electronic device. (3) Analyzing the computation precision and nonnegative networks. As discussed in the previous subsections, photonics is most suitable with nonnegative weights which can be directly mapped to the power ratios. The capability of defining weights using only nonnegative values would significantly simplify the design, fabrication, and control for the optical programmable processers. However, conventional training algorithms are developed using complementary (+/-) weight factors. Thus, it is important to investigate how the resolution level, nonnegative mode, and network size affect the accuracy of a neural network

for a target task. (4) System-level scheduling and orchestration. To maximally utilize both types of devices, the latency of each device should be taken into account during the system-level scheduling and orchestration.

# 1) Fully-connected layers: computation breakdown

For a fully-connected layer  $h^{j+1}$  of size  $i_{j+1}$  from the previous layer  $h^j$  of size  $i_j$ , let  $W^{j+1} \in \mathbb{R}^{i_{j+1} \times i_j}$  denote the weight matrix. Note that  $i_j$  and  $i_{j+1}$  can be much larger than N. Given an activation function Act() and bias  $b^{j+1}$ , the layer  $h^{j+1}$  can be computed as follows:

$$h^{j+1} = Act(g^{j+1} + b^{j+1})$$

$$g^{j+1} = W^{j+1} \cdot h^{j}$$
(30)
(31)

To compute  $g^{j+1}$  using the aforementioned PIC, we can partition the input into matrices of size  $N \times N$  and vectors of size N as follows for  $0 \le k \le \frac{i_{j+1}}{N}$ :

$$g_{k+1}^{j+1} = g_{k+2}^{j+1} = \sum_{\ell=0}^{i_{j}} \begin{pmatrix} W^{j+1} & W^{j+1} & \cdots & W^{j+1} \\ k+1,\ell+1 & k+1,\ell+2 & \cdots & k+1,\ell+N \\ W^{j+1} & W^{j+1} & \cdots & W^{j+1} \\ k+2,\ell+1 & k+2,\ell+2 & \cdots & k+1,\ell+N \\ \vdots & \cdots & \ddots & \vdots \\ W^{j+1} & W^{j+1} & \cdots & W^{j+1} \\ k+N,\ell+1 & k+N,\ell+2 & \cdots & W^{j+1} \\ k+N,\ell+1 & \cdots & k+N,\ell+2 \end{pmatrix} \cdot \begin{pmatrix} h^{j} \\ \ell+1 \\ h^{j} \\ \ell+2 \\ \vdots \\ h^{j} \\ \ell+N \end{pmatrix}$$
(32)

The overview of this approach is also depicted in Fig. 16.

The total number of multiplications required to compute layer  $h^{j+1}$  from  $h^j$  is  $i_{j+1} \cdot i_j$ . With the above approach using OMMs of width N, the total number of OMM requests is  $\left\lceil \frac{i_{j+1}}{N} \right\rceil$ . This reduction by the factor of  $\frac{1}{N^2}$  is achievable because there is no waste of operations associated with the partitioning.

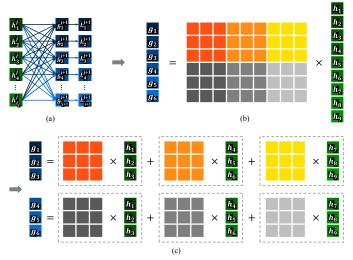


Fig. 16 Computation for fully-connected layers using optical matrix-vector multipliers (OMMs). (a) Fully-connected layers  $h^j$  (green) and  $h^{j+1}$  (blue).  $g^{j+1}$  is obtained as a result of the inner products between  $h^j$  and the weight vectors.  $h^{j+1}$  is obtained by applying the bias and activation to  $g^{j+1}$ . (b) Matrix-vector multiplication between the weight matrix (orange and gray) and  $h^j$  to obtain  $g^{j+1}$ . The superscripts are omitted for simplicity. (c) Computation equivalent to that of (b) but using OMMs with the input matrix size of  $3 \times 3$ .

# 2) Convolutional layers: minimization of the reconfiguration of OMM input matrices

Fig. 17 shows the convolution part of convolutional layers

computed using OMMs. The total number of multiplications required in computing one output channel is:

$$i_{in\ ch}\cdot (W-2)\cdot (H-2)\cdot N^2\tag{33}$$

where  $i_{in\_ch}$  denotes the number of input channels, W and H denote the width and height of an input channel, and  $N^2$  represents the size of the convolution kernel. With the above approach, the total number of OMM requests for computing one output channel is  $i_{in\ ch} \cdot (W-2) \cdot (H-2)$ .

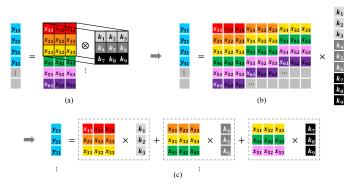


Fig. 17 Computation for convolutional layers using OMMs. The first column of the output channel (nodes  $y_{11}, \cdots$ ) and the first three columns of the input channel (nodes  $x_{11}, \cdots$ ) are shown in the above illustration. (a) Convolutions on a single channel of the input layer. The convolution results over all channels in the input layer will be summed up and mapped to the output channel after the bias and activation are applied. (b) Conversion of the convolutions to a matrix-vector multiplication. The computation for one column in the output channel can be performed by a single matrix-vector multiplication for each input channel. (c) Computation equivalent to that of (b) but using OMMs with the input matrix size of  $3 \times 3$ , which equals the size of the convolution kernel.

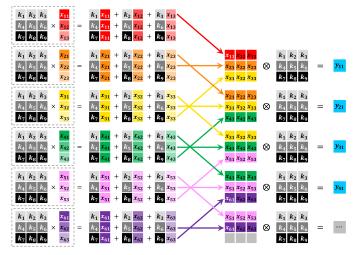


Fig. 18 Proposed computation of convolutional layers using an OMM without updating its input matrix values.

The above approach updates the matrix elements for each OMM request. On the other hand, we propose another approach illustrated in Fig. 18, which minimizes the number of updates of the input matrix for the OMM. This approach follows a similar direction to the weight stationary optimization technique of ANN accelerators, where the weight values stay in the local register file of processing elements of the hardware accelerators [97]. The fundamental goal of this optimization is to minimize the time for processing elements to be reading the weights. In our co-designed system, the weights must be

converted to analog signals and passed to the OMM to be set up for the computation. Thus, we aim at reducing the latency of the overall process by minimizing the number of the OMM's input matrix updates. This can be achieved by mapping the convolution kernel itself to the OMM's input matrix, when the size of the OMM's input matrix is larger than or equal to that of the convolution kernel, which often ranges between  $2 \times 2$  and  $5 \times 5$ . The convolution kernel weights form the input matrix and the network nodes form the input vectors. Then, the results of the first N = 3 matrix-vector multiplication instances in Fig. 10 contain the convolution result for  $y_{11}$ . The second, third, and forth matrix-vector multiplication results contain the convolution result for  $y_{21}$ . Consecutive N=3 results contain the convolution result for the corresponding output element. While processing the entire input channel, the input matrix for the OMM does not change. With this approach, the total number of OMM requests for one output channel is  $i_{in ch}$ .  $(W-2)\cdot H$ .

#### 3) Analysis on the nonnegative property and resolutions

Most neural networks used in practice have both positive and negative input values, weights, and node values. Thus, feedforward propagation of these networks, either during the training or inference, requires matrix-vector multiplications with both positive and negative values. Then, it is of interest to consider a mapping between the values in the range of [-1, 1] and the range of [0, 1] such that matrix-vector multiplication is preserved by this mapping. However, the theorem in Appendix I verifies that such mapping does not exist. There have been approaches to use only nonnegative input and weights to obtain a more understandable network with slight decrease in the accuracy [98]. Another approach performs nonnegative matrix factorization of the weights in order to reduce the input complexity, but the input values in this case can be both positive and negative [99].

To avoid matrix-vector multiplication with negative values, we train the neural networks using nonnegative input, weights, and nodes. In our experiment, we restrict not only the sign of the input and weights to be nonnegative but also the resolution used during inference. Figure 19 shows the estimated inference accuracy of 2-layer MLPs over a range of the resolution levels (the number of bits used to represent the input values and weights in a fixed-point format), and the network sizes (the number of nodes in the hidden layer of the MLP) trained in two different modes for the task of handwritten digit recognition: (a) conventional mode that supports negative input, weights, and nodes, and (b) nonnegative mode that normalizes the input to [0, 1], and constrains the weights and nodes to be nonnegative. One network for each mode and each level of the network size was trained using the MNIST train dataset [100], with 32-bit floating point representation [101]. The input image contains  $28 \times 28$  pixel values in the range of [0, 255], which were normalized to [-1, 1] or [0, 1] depending on the training mode. For activation functions, tanh was used in the hidden layer, and softmax was used in the output layer. After activation in the nonnegative mode, all negative values were rounded up to 0. All weights and biases were randomly initialized, and the weights for the nonnegative mode were initialized to [0, 1].

These weights and biases were updated using ADAM, which is a state-of-the-art stochastic back-propagation method [68].

Each of the trained networks was tested on the MNIST test dataset, with both the input values and weights converted to the fixed-point representation for each resolution level. We note that one instance of a trained network with a given network structure does not represent the most optimized network of that structure. Nevertheless, all networks in this test case were trained using the same approach with similar optimization efforts, aside from the training time which increases for larger networks. Thus, we refer to these networks in order to practically and roughly estimate the performance trend over various network sizes, resolution levels, and the training mode. As shown in Fig. 19, the test accuracy has generally improved as the network width increased and as the resolution level was enhanced. It turns out that the accuracy of networks trained in the conventional mode were more affected by the restricted resolution, whereas the accuracy of those trained in the nonnegative mode were more affected by the network width.

The test accuracy achieved by nonnegative networks are lower than that by the conventionally trained counterpart, but a larger nonnegative network can sometimes outperform a smaller conventional network. During the training in the nonnegative mode, the biases and activation functions were allowed to take negative values because in this co-design approach only the matrix-vector multiplications will be offloaded from the electronic device to the optical device. This seems to have enabled the network to cut out less relevant, or negatively related connections and to focus on positively related ones, resulting in comparable accuracy for large nonnegative networks.

The issue of positive and negative inputs is an interesting example of the approach to optimization required for co-design. As mentioned in section IV.A, photonics is implemented more readily with nonnegative values. This initial investigation indicates that, although in current practice both positive and negative values are used, using only nonnegative values for the matrix-vector multiplications can actually be advantageous in some circumstances.

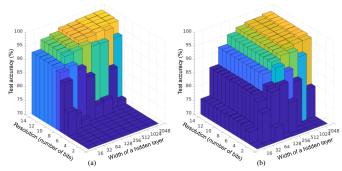


Fig. 19 Test accuracy of MLPs for handwritten digit recognition with varying resolutions and network sizes. (a) Networks trained in the conventional mode using negative values, 0, and positive values in the computation. (b) Networks trained in the nonnegative mode, where only 0 and nonnegative values are used during matrix-vector multiplications.

## 4) System-level scheduling to maximize the throughput

To accelerate the inference process of a trained neural network with OMMs, an FPGA-based co-designed system

down the computation, sends matrix-vector breaks multiplication requests to OMMs, and performs the remaining part of the computation including the nonlinear activation (which could also be done optically or via well-designed analog electronics as discussed in section IV.C.4). Figure 20 illustrates the overview of the proposed co-designed system that contains three specialized processors: the ANN processor, the input processor, and the output processor. For each OMM request, the ANN processor sends the input  $M \times N$  matrix K to the MRRs via DACs, and the input processor sends the input Ndimensional vector x to the modulators via DACs. The output processor collects the resulting M-dimensional vector y from the PDs via ADCs, and it also applies the bias and nonlinear activation function. The very recent demonstration on a 1-to-56 Gb/s ADC/DAC-based transceiver [102] paves the way for high-speed, low-energy ADC/DACs as the interface between the OMM and FPGA, without harming the throughput.

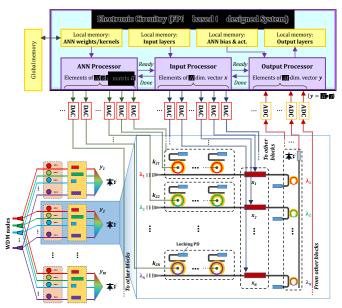


Fig. 20 System-level overview for the proposed co-design approach. The FPGA-based electronic system (on the top) invokes and controls the optical system (in the bottom). The MRRs that receive electrical signals from DACs act as electrical-to-optical converters, whereas the summation PDs perform the optical-to-electrical conversions. The summed signals are connected to the FPGA via ADCs. Details regarding the memory systems, which depend on the specific application, are abstracted in this figure.

Although the computation complexity of an OMM is in O(1), the DAC, MRR configuration, ADC, and the computation on the FPGA consume non-negligible latency. The goal of the system-level scheduling is to overlap these latencies to maximize the throughput. Figure 21 shows abstract timing diagrams with pipelined executions by the ANN, input, and output processors. Figure 21a illustrates the case of invoking a single OMM instance. As shown in Fig. 21b, the latency  $T_L$  of a period between consecutive OMM invocations can be expressed as:

$$T_L = T_M + T_{DA} + T_{AD} \tag{34}$$

 $T_L = T_M + T_{DA} + T_{AD}$  (34) where  $T_M$  denotes the latency of the DACs and MRR configuration,  $T_{DA}$  the latency of DACs, and  $T_{AD}$  the latency of ADCs. This holds as long as the ANN processor's latency  $T_A$ does not exceed  $T_{DA} + T_{AD}$ , and similarly, the input and output

processors' latencies  $T_I$  and  $T_O$  are less than or equal to  $T_M$  +  $T_{AD}$  and  $T_M + T_{DA}$ , respectively.

When consecutive OMM instances contain the same input matrix elements so that it is not needed to reconfigure the MRRs, the latency  $T_L$  of the period can be expressed as:

 $T_L = T_{DA} + T_{AD}$  (35) as shown in Fig. 21c. In both cases of Fig. 21b and 21c, the asymptotic throughput is proportional to  $\frac{1}{T_L}$  and the number of OMM devices that can be processed in parallel, and is inversely proportional to the total numbers of OMM invocations for fullyconnected or convolutional layers which have been discussed in the previous subsections.

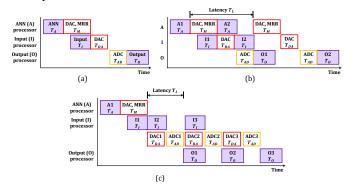


Fig. 21 (a) Timing diagram of invoking one OMM instance containing an input matrix and vector. (b) Timing diagram of invoking 2 OMM instances, each with an input matrix and vector. More invocations can be added on the right in a similar pattern. The latency of a period between consecutive invocations is denoted as  $T_{I}$ . (c) Timing diagram of invoking multiple OMM instances, where the first instance contains a new input matrix and vector and subsequent instances contain only new input vectors. The latency  $T_L$  has been reduced with respect to (b).

## C. Discussion

# 1) Silicon Ring Resonators: Finesse vs. Bandwidth

Silicon ring resonators with high finesse (up to a few hundreds) have been extensively demonstrated [90, 91]. However, these demonstrations aim for high quality factors and tend to have a relatively small 3 dB bandwidth. For the aggregation ring filters in this OMM system, a large 3 dB bandwidth is an equally important factor that allows high data rate vectors to be fanned in, for high computational speeds. It would be preferable for the operation bandwidth of such an OMM to match that of the photo-detection rate (typically at 100 GHz).

The recent demonstration of a submicron-scale MRR shows great potential for the aggregation ring filters with high finesse and large bandwidth [103]. It features a 3 dB bandwidth of 100 GHz and a finesse of 116, supporting up to 116 wavelength channels given a 1 dB through loss budget as discussed in section IV.A.2.3. This ultra-small ring resonator has the additional benefit of reducing the thermal tuning power, which is proportional to its size [103]. Another notable demonstration that combines an MRR based filter with grating-assisted contradirectional couplers frees the constraint of FSR [104]. The addition of grating-assisted couplers provides an extra degree of freedom for longitudinal mode selectivity. This design, therefore, paves the way for independent optimization of the

3 dB bandwidth, and potentially enables a full utilization of the transmission window in the silicon platform, yielding an extremely scalable OMM.

# 2) Optical Phase shifting technology

Phase shifter technology is key in the OMM. Thermo-optic phase shifting is preferred since it is the most commonly applied lossless mechanism in the silicon platform. The induced on-chip thermal crosstalk can be reduced by adding isolation trenches [105]. In addition, a selective silicon etch can be applied to the silicon substrate to undercut the waveguides. The selective etch localizes the heat and improves the heating efficiency [106]. The reduced heating power could in turn ameliorate on-chip thermal perturbations. The limited thermal frequency response (up to a few hundred KHz [82]) is, however, a limiting factor in latency, when dynamic reconfiguration for the OMM is required. For fast phase tuning, as aforementioned, electro-optic phase shifting leveraging the plasma dispersion effect is the most popular all-silicon technology [96]. It offers nanosecond-scale reconfiguration time, albeit with some performance penalty due to the electro-absorption loss. The E-O phase shifters would be straight-forwardly included in the weighting blocks with additional considerations for the excess electro-absorption loss.

With the advances in heterogeneous integration technology, other materials can be introduced on the silicon platform. Notable examples include III-V materials [49], graphene [107], and nonvolatile phase-change materials (PCMs) [108]. III-V materials exhibit high electro-optic phase modulation efficiency, which can be effectively combined with silicon waveguides using wafer-bonding techniques [49]. Thin layers of graphene can be deposited on top of the Si waveguide [107]. forming a capacitor that overlaps with the tail of the waveguide's optical mode. The application of voltage will then shift the Fermi level of graphene and enable inter-band transitions of charge carriers, and thus modulate the intensity of light travelling through the waveguide. The PCMs can introduce gigantic optical phase changes and most importantly, such phase changes are nonvolatile. This nonvolatility adds the capability of self-holding, maintaining optical states even in the absence of power input [109].

## 3) Power consumption and footprint of the OMM

The power consumption of the OMM is dominated by the tuning and locking of MRR elements. Current technology features a thermo-optic tuning efficiency of 1 nm/mW with doped-silicon micro-heaters [82], leading to a small power consumption of a few mW per weighting MRR. Femtojoule-level depletion-mode modulators in vertically doped micro-disk structures [32], featuring low operating voltage (0.5  $V_{PP}$ ), offer the possibility for ultralow power electro-optic OMMs. The power consumption would then derive from the undesired leakage current, approximately  $\sim\!\!6~\mu\mathrm{W}$  per element [32]. In future implementations, the phases could be set using the nonvolatile PCMs [109]. In that case, power would only be drawn during state transitions. A recent demonstration on a nonvolatile PCM-based photonic memory cell shows

programming energy and time of only 680 pJ and 250 ns, respectively [110].

A number of wavelength locking schemes have been proposed, including the use of the photoconductive effect [87], small dithering signals [84], radio frequency (RF) detection [86], additional partial drop rings [88], and monolithically integrated locking controllers in the 45 nm CMOS-SOI platform [85]. The locking power consumption has been demonstrated to be in the range of a few hundred µW [84, 85]. Furthermore, there has also been noteworthy research progress on athermal MRRs that could significantly overcome the temperature sensitivity [111-113]. Here, the key idea is to introduce an upper cladding that has a negative thermo-optic coefficient to counteract the T-O effect of silicon. Titanium dioxide (TiO<sub>2</sub>) holds the most promise as it exhibits a relatively strong negative thermo-optic coefficient and has been included in the CMOS-compatible fabrication process [112, 113]. This technique offers a path to extremely power-efficient OMM units.

Current implementation of MRR-based PICs for on-off switching (two-state) applications normally features a pitch size of 100  $\mu m$  [29]. Hundreds of MRR elements have been monolithically integrated on a single chip, within an area of a few tens of millimeter squares [29]. The temperature-insensitive MRRs could potentially reduce the footprint of the OMM significantly, even for high-resolution operations, as the pitch limitation due to thermal restrictions is offset. The size will then be merely limited by the pitch size of electrical bonding pads, which can be as small as 25  $\mu m$  to 40  $\mu m$  [114], thus enabling the footprint shrink of the OMM by over an order of magnitude.

#### 4) Nonlinear activation function

To implement a full neural network, as aforementioned, a nonlinear activation function is required in addition to the linear OMM units. For a nonlinear activation function implemented in optics, there are generally two types, implemented using: (1) electro-optic nonlinearity and (2) all-optical nonlinearity. The former type requires first converting an optically weighted signal into the electrical domain and then triggering the nonlinear activation function to have an optical outcome. Examples include semiconductor excitable lasers (type C in Fig. 6c) [77], and electro-absorption modulators [115]. This type of solution might impair the processing speed and cascadability of neural networks due to the movement of charge carriers and the optical-to-electrical conversion noise. The latter, all-optical solution holds greater promise. The most commonly used optical nonlinearities are saturable absorption, such as in the use of monolayer graphene absorbers [116] and two-photon absorption [117], and bistability in nonlinear photonic crystals [118] and optical superlattices [119]. The nonlinearity of ring resonators can also be exploited [120]. Currently, the optical nonlinear activation function is an important research topic which could be used in order to enhance the throughput of an optical neural network, thus lowering the system latency and power consumption. However, the monolithic integration of these nonlinear units with OMMs,

the efficiency of the nonlinear modulation, and the operational speed and accuracy are open challenges [121].

While the development of an all-optical on-chip neural network represents a longer-term goal, implementing the nonlinear activation function electrically is a promising alternative in the short term. The very recent work of building optical neural networks based on photoelectric multiplication also proposes to implement the nonlinear activation function in the electrical domain [122]. Very low power (femtojoule-scale) consumption is feasible with well-designed analog electronics.

#### V. CONCLUSION

Larger DNNs in general have higher expressiveness as a classification function. Theoretical analysis has also verified that both the depth and the width of neural networks contribute to their expressive power. It has been shown that complex functions expressed by deep neural networks cannot be approximated by any shallow neural network whose size is no more than an exponential bound [123], and also that certain classes of wide neural networks cannot be realized by any narrow network whose depth is no more than polynomial bound [124]. These observations lead to the demand for the capability to efficiently process very deep or wide neural networks. The co-design approach addresses scalability (in terms of the size of neural networks) in two aspects: (1) The capability to decompose a large matrix-vector multiplication into smaller instances which significantly relaxes the requirement of photonic integrations. (2) A path to construct ultra-large scale OMMs using MRRs in the wavelength domain. This reduces the system decomposition complexity and, in turn, enables the handling of sophisticated concepts for future applications. In addition, the approach to manage the computation precision with nonnegative values can be utilized in any photonic systems, in order to reduce the implementation complexity and thus cost. This also facilitates the operation of different facets of validity in practical terms for OMMs as hardware accelerators in deep learning applications.

In summary, efficient scaling of deep learning will require dedicated hardware accelerators. We have presented an overview of silicon photonics applications for deep learning and have analyzed opportunities for scalable co-designed multiwavelength microring silicon photonic architectures.

# APPENDIX I

Theorem 1. Let  $\Phi, \Omega \subset \mathbb{R}$  such that  $\{-1, 0, 1\} \in \Phi$  and  $\Omega \subset \mathbb{R}$  $[0, +\infty]$ . Then there exists no function  $f: \Phi \to \Omega$  satisfying the followings:

For any 
$$p_1, p_2 \in \Phi$$
,  $f(p_1) + f(p_2) = f(p_1 + p_2)$  (A1)

For any 
$$p_1, p_2 \in \Phi$$
,  $f(p_1) \cdot f(p_2) = f(p_1 \cdot p_2)$  (A2)

This also holds if  $\Omega \subset (-\infty, 0)$ .

*Proof.* If such function f exists, it must satisfy the followings:

$$f(1) + f(0) = f(1)$$
 (A3)

$$f(1) \cdot f(-1) = f(-1)$$
 (A4)

$$f(1) + f(-1) = f(0)$$
 (A5)

Equation A3 implies that f(0) = 0, and Equation A4 implies that f(1) = 1. Then, Equation A5 can be re-written as

$$1 + f(-1) = 0 \tag{A6}$$

Thus, f(-1) = -1 but this value is not in the range  $\Omega$  of function f. Therefore, such f does not exist.

#### APPENDIX II

As discussed in section IV.A.1.1, the maximum sensitivity of the weight in Equation 11 occurs when

$$\tan\frac{\phi_{\text{max}}}{2} = \sqrt{\frac{r(a)}{s(a)}} \tag{A7}$$

where

$$r(a) = 3a^2 + 2 - \sqrt{9a^4 + 4a^2 + 4}$$
 (A8)

and

$$s(a) = a^2 - 2 + \sqrt{9a^4 + 4a^2 + 4}$$
. (A9)

Therefore,

$$\sin \phi_{\text{max}} = \frac{2\sqrt{r(a)s(a)}}{r(a)+s(a)} = \frac{1}{2a^2} \sqrt{r(a)s(a)}$$
 (A10)  

$$\cos \phi_{\text{max}} = \frac{s(a)-r(a)}{s(a)+r(a)} = 1 - \frac{r(a)}{2a^2}$$
 (A11)

$$\cos \phi_{\text{max}} = \frac{s(a) - r(a)}{s(a) + r(a)} = 1 - \frac{r(a)}{2a^2}$$
 (A11)

Plugging these back into the sensitivity function of Equation 12 and assuming  $T_0 \approx 0$  immediately yields:

$$\left|\frac{\partial \mu}{\partial \phi}\right|_{\text{max}} = \frac{1+a^2}{a^2} \frac{4\sqrt{r(a)s(a)}}{(r(a)+4)^2}$$
 (A12)

Assuming that  $a^2 \gg 1$  we see that

$$\frac{1+a^2}{a^2} \approx 1 \tag{A13}$$

$$\frac{1+a^2}{a^2} \approx 1$$

$$r(a) \approx 3a^2 + 2 - 3a^2 \left(1 + \frac{2}{9a^2}\right) = \frac{4}{3}$$
(A13)

$$s(a) \approx a^2 - 2 + 3a^2 \left(1 + \frac{2}{9a^2}\right) \approx 4a^2$$
 (A15)

Therefore,

$$\left| \frac{\partial \mu}{\partial \phi} \right|_{\text{max}} \approx \frac{9}{16\sqrt{3}} \ a \ .$$
 (A16)

#### REFERENCES

- 1. Simonyan, K. and A. Zisserman, Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:. 2014.
- 2. Hershey, S., et al. CNN architectures for large-scale audio classification. in 2017 ieee international conference on acoustics, speech and signal processing (icassp). 2017. IEEE.
- Bojarski, M., et al., End to end learning for self-driving cars. arXiv 3. preprint arXiv:.07316, 2016.
- 4. Esteva, A., et al., Dermatologist-level classification of skin cancer with deep neural networks. 2017. 542(7639): p. 115.
- 5 Haensch, W., T. Gokmen, and R. Puri, The Next Generation of Deep Learning Hardware: Analog Computing. Proceedings of the IEEE, 2019. 107(1): p. 108-122.
- Hu, M., et al., Dot-product engine for neuromorphic computing: programming 1T1M crossbar to accelerate matrix-vector multiplication, in Proceedings of the 53rd Annual Design Automation Conference. 2016, ACM: Austin, Texas. p. 1-6.
- 7. Athale, R.A. and W.C. Collins, Optical matrix-matrix multiplier based on outer productdecomposition. Applied Optics, 1982. 21(12): p. 2089-2090.
- Farhat, N.H., et al., Optical implementation of the Hopfield model. 8. Applied Optics, 1985. 24(10): p. 1469-1475.
- Jordan, M.I. and T.M. Mitchell, Machine learning: Trends, 9. perspectives, and prospects. Science, 2015. 349(6245): p. 255.

- Goodfellow, I., Y. Bengio, and A. Courville, *Deep learning*. 2016: MIT press.
- McCulloch, W.S. and W. Pitts, A logical calculus of the ideas immanent in nervous activity. The bulletin of mathematical biophysics, 1943. 5(4): p. 115-133.
- 12. Rumelhart, D.E., G.E. Hinton, and R.J. Williams, Learning internal representations by error propagation, in Parallel distributed processing: explorations in the microstructure of cognition, vol. 1, E.R. David, L.M. James, and C.P.R. Group, Editors. 1986, MIT Press. p. 318-362.
- Cybenko, G., Approximation by superpositions of a sigmoidal function. Mathematics of Control, Signals and Systems, 1989. 2(4): p. 303-314.
- Hornik, K., M. Stinchcombe, and H. White, Multilayer feedforward networks are universal approximators. Neural Networks, 1989.
   2(5): p. 359-366.
- LeCun, Y., et al., Backpropagation Applied to Handwritten Zip Code Recognition. Neural Computation, 1989. 1(4): p. 541-551.
- Sze, V., et al., Efficient Processing of Deep Neural Networks: A Tutorial and Survey. Proceedings of the IEEE, 2017. 105(12): p. 2295-2329.
- Bagherian, H., et al., On-Chip Optical Convolutional Neural Networks. arXiv preprint arXiv: 03303, 2018.
- Bergstra, J., et al., Theano: Deep learning on gpus with python. Journal of Machine Learning Research 2011(1): p. 1-48.
- Chetlur, S., et al., cudnn: Efficient primitives for deep learning. arXiv preprint arXiv, 2014.
- Ovtcharov, K., et al., Accelerating deep convolutional neural networks using specialized hardware. Microsoft Research Whitepaper, 2015. 2(11): p. 1-4.
- Amiri, S., et al. Multi-precision convolutional neural networks on heterogeneous hardware. in 2018 Design, Automation & Test in Europe Conference & Exhibition (DATE). 2018.
- Esser, S.K., et al., Convolutional networks for fast, energy-efficient neuromorphic computing. Proceedings of the National Academy of Sciences, 2016. 113(41): p. 11441.
- Misra, J. and I. Saha, Artificial neural networks in hardware: A survey of two decades of progress. Neurocomputing, 2010. 74(1): p. 239-255
- 24. Chen, Y., et al., Eyeriss: An Energy-Efficient Reconfigurable Accelerator for Deep Convolutional Neural Networks. IEEE Journal of Solid-State Circuits, 2017. **52**(1): p. 127-138.
- 25. Casasent, D. and A. Ghosh, *Optical linear algebra processors:* noise and error-source modeling. Optics Letters, 1985. **10**(6): p. 252-254
- Liang, Y.-Z. and H.K. Liu, Optical matrix-matrix multiplication method demonstrated by the use of a multifocus hololens. Optics Letters, 1984. 9(8): p. 322-324.
- Shen, Y., et al., Deep learning with coherent nanophotonic circuits. Nature Photonics, 2017. 11: p. 441.
- 28. Cheng, Q., et al., Scalable Microring-Based Silicon Clos Switch Fabric With Switch-and-Select Stages. IEEE Journal of Selected Topics in Quantum Electronics, 2019. 25(5): p. 1-11.
- Cheng, Q., et al., Recent advances in optical technologies for data centers: a review. Optica, 2018. 5(11): p. 1354-1370.
- Lim, A.E., et al., Review of Silicon Photonics Foundry Efforts. IEEE Journal of Selected Topics in Quantum Electronics, 2014. 20(4): p. 405-416.
- Rahim, A., et al., Open-Access Silicon Photonics: Current Status and Emerging Initiatives. Proceedings of the IEEE, 2018. 106(12): p. 2313-2330.
- Timurdogan, E., et al., An ultralow power athernal silicon modulator. Nature Communications, 2014. 5: p. 4008.
- Liu, S., et al., High-channel-count 20GHz passively mode-locked quantum dot laser directly grown on Si with 4.1Tbit/s transmission capacity. Optica, 2019. 6(2): p. 128-134.
- Chen, P., et al., High-order microring resonators with bent couplers for a box-like filter response. Optics Letters, 2014. 39(21): p. 6304-6307.
- Chrostowski, L. and M. Hochberg, Silicon photonics design: from devices to systems. 2015: Cambridge University Press.
- Sun, J., et al., Large-scale nanophotonic phased array. Nature, 2013. 493: p. 195.

- 37. Qiang, X., et al., Large-scale silicon quantum photonics implementing arbitrary two-qubit processing. Nature Photonics, 2018. 12(9): p. 534-539.
- 38. Yu, M., et al., Silicon-chip-based mid-infrared dual-comb spectroscopy. Nature Communications, 2018. 9(1): p. 1869.
- 39. Pérez, D., et al., *Multipurpose silicon photonics signal processor core*. Nature Communications, 2017. **8**(1): p. 636.
- 40. Tait, A.N., et al., Neuromorphic photonic networks using silicon photonic weight banks. Scientific Reports, 2017. 7(1): p. 7430.
- 41. Hu, H., et al., Single-source chip-based frequency comb enabling extreme parallel data transmission. Nature Photonics, 2018. 12(8): p. 469-473.
- Cheng, Q., et al., Photonic switching in high performance datacenters [Invited]. Optics Express, 2018. 26(12): p. 16022-16043
- 43. Luo, L.-W., et al., *WDM-compatible mode-division multiplexing on a silicon chip.* Nature Communications, 2014. **5**: p. 3069.
- 44. Kwon, K., et al. 128×128 Silicon Photonic MEMS Switch with Scalable Row/Column Addressing. in Conference on Lasers and Electro-Optics. 2018. San Jose, California: Optical Society of America.
- 45. Michel, J., J. Liu, and L.C. Kimerling, *High-performance Ge-on-Si photodetectors*. Nature Photonics, 2010. 4: p. 527.
- Feng, D., et al., High-Speed GeSi Electroabsorption Modulator on the SOI Waveguide Platform. IEEE Journal of Selected Topics in Quantum Electronics, 2013. 19(6): p. 64-73.
- 47. Sacher, W.D., et al., Monolithically Integrated Multilayer Silicon Nitride-on-Silicon Waveguide Platforms for 3-D Photonic Circuits and Devices. Proceedings of the IEEE, 2018. 106(12): p. 2232-2245.
- 48. Marshall, O., et al., *Heterogeneous Integration on Silicon Photonics*. Proceedings of the IEEE, 2018. **106**(12): p. 2258-2269.
- Komljenovic, T., et al., Photonic Integrated Circuits Using Heterogeneous Integration on Silicon. Proceedings of the IEEE, 2018. 106(12): p. 2246-2257.
- 50. Chen, S., et al., *Electrically pumped continuous-wave III–V quantum dot lasers on silicon*. Nature Photonics, 2016. **10**: p. 307.
- Xu, Q., et al., Micrometre-scale silicon electro-optic modulator. Nature, 2005. 435(7040): p. 325-7.
- 52. Zhang, C., et al., Highly linear heterogeneous-integrated Mach-Zehnder interferometer modulators on Si. Optics Express, 2016. 24(17): p. 19040-19047.
- 53. Ovvyan, A.P., et al., Cascaded Mach–Zehnder interferometer tunable filters. Journal of Optics, 2016. **18**(6): p. 064011.
- Bahadori, M., et al., Comprehensive Design Space Exploration of Silicon Photonic Interconnects. Journal of Lightwave Technology, 2016. 34(12): p. 2975-2987.
- 55. Luo, L.-W., et al., *High bandwidth on-chip silicon photonic interleaver*. Optics Express, 2010. **18**(22): p. 23079-23087.
- Cheng, Q., et al., Ultralow-crosstalk, strictly non-blocking microring-based optical switch. Photonics Research, 2019. 7(2): p. 155-161
- 57. DasMahapatra, P., et al., *Optical Crosspoint Matrix Using Broadband Resonant Switches*. IEEE Journal of Selected Topics in Quantum Electronics, 2014. **20**(4): p. 1-10.
- 58. Cheng, Q., et al., Demonstration of the feasibility of large-port-count optical switching using a hybrid Mach-Zehnder interferometer-semiconductor optical amplifier switch module in a recirculating loop. Opt Lett, 2014. 39(18): p. 5244-7.
- 59. Mesaritakis, C., V. Papataxiarhis, and D. Syvridis, *Micro ring resonators as building blocks for an all-optical high-speed reservoir-computing bit-pattern-recognition system.* Journal of the Optical Society of America B, 2013. **30**(11): p. 3048-3055.
- Tait, A.N., et al., Broadcast and Weight: An Integrated Network For Scalable Photonic Spike Processing. Journal of Lightwave Technology, 2014. 32(21): p. 4029-4041.
- 61. Ramaswamy, V. and R.D. Standley, *A phased, optical, coupler-pair switch*. The Bell System Technical Journal, 1976. **55**(6): p. 767-775.
- Cheng, Q., et al., Scalable, Low-Energy Hybrid Photonic Space Switch. Journal of Lightwave Technology, 2013. 31(18): p. 3077-3084
- 63. Bogaerts, W., et al., Silicon microring resonators. Laser & Photonics Reviews, 2012. 6(1): p. 47-73.
- 64. Bahadori, M., et al., Design Space Exploration of Microring Resonators in Silicon PhotonicInterconnects: Impact of the Ring

- Curvature. Journal of Lightwave Technology, 2018. 36(13): p. 2767-2782.
- Reck, M., et al., Experimental realization of any discrete unitary operator. Physical Review Letters, 1994. 73(1): p. 58-61.
- Clements, W.R., et al., Optimal design for universal multiport interferometers. Optica, 2016. 3(12): p. 1460-1465.
- Carolan, J., et al., *Universal linear optics*. 2015. 349(6249): p. 711-716.
- Miller, D.A.B., Self-configuring universal linear optical component [Invited]. Photonics Research, 2013. 1(1): p. 1-15.
- Miller, D.A.B., Meshing optics with applications. Nature Photonics, 2017. 11: p. 403.
- Harris, N.C., et al., Quantum transport simulations in a programmable nanophotonic processor. Nature Photonics, 2017.
   11: p. 447.
- 71. Ribeiro, A., et al., *Demonstration of a 4x4-port universal linear circuit.* Optica, 2016. **3**(12): p. 1348-1357.
- 72. Miller, D.A.B., Setting up meshes of interferometers reversed local light interference method. Optics Express, 2017. 25(23): p. 29233-29248.
- 73. Schmidhuber, J., *Deep learning in neural networks: An overview*. Neural Networks, 2015. **61**: p. 85-117.
- 74. Lindblad, T., J.M. Kinser, and J. Taylor, *Image processing using pulse-coupled neural networks*. 1998: Springer.
- Maass, W., Networks of spiking neurons: The third generation of neural network models. Neural Networks, 1997. 10(9): p. 1659-1671.
- Thorpe, S., A. Delorme, and R. Van Rullen, Spike-based strategies for rapid processing. Neural Networks, 2001. 14(6): p. 715-725.
- 77. Ferreira de Lima, T., et al., *Progress in neuromorphic photonics*. Nanophotonics, 2017. **6**(3): p. 577-599.
- Tait, A.N., et al., Microring Weight Banks. IEEE Journal of Selected Topics in Quantum Electronics, 2016. 22(6): p. 312-325.
- Tait, A.N., et al., Continuous Calibration of Microring Weights for Analog Optical Networks. IEEE Photonics Technology Letters, 2016. 28(8): p. 887-890.
- Tait, A.N., et al., Multi-channel control for microring weight banks.
   Optics Express, 2016. 24(8): p. 8895-8906.
- 81. Tait, A.N., et al., Feedback control for microring weight banks. Optics Express, 2018. **26**(20): p. 26422-26443.
- 82. Bahadori, M., et al., Thermal Rectification of Integrated Microheaters for Microring Resonators in Silicon Photonics Platform. Journal of Lightwave Technology, 2017: p. 1-1.
- 83. Agarap, A.F.M., On breast cancer detection: an application of machine learning algorithms on the wisconsin diagnostic dataset, in Proceedings of the 2nd International Conference on Machine Learning and Soft Computing. 2018, ACM: Phu Quoc Island, Viet Nam. p. 5-9.
- 84. Padmaraju, K., et al., Wavelength Locking and Thermally Stabilizing Microring Resonators Using Dithering Signals. Journal of Lightwave Technology, 2014. 32(3): p. 505-512.
- 85. Sun, C., et al., A 45 nm CMOS-SOI Monolithic Photonics Platform With Bit-Statistics-Based Resonant Microring Thermal Tuning. IEEE Journal of Solid-State Circuits, 2016. 51(4): p. 893-907.
- Dong, P., et al., Simultaneous wavelength locking of microring modulator array with a single monitoring signal. Optics Express, 2017. 25(14): p. 16040-16046.
- Jayatilleka, H., et al., Photoconductive heaters enable control of large-scale silicon photonic ring resonator circuits. Optica, 2019. 6(1): p. 84-91.
- Khope, A.S.P., et al., On-chip wavelength locking for photonic switches. Opt Lett, 2017. 42(23): p. 4934-4937.
- Zortman, W.A., D.C. Trotter, and M.R. Watts, Silicon photonics manufacturing. Optics Express, 2010. 18(23): p. 23598-23607.
- Xu, Q., D. Fattal, and R.G. Beausoleil, Silicon microring resonators with 1.5-µm radius. Optics Express, 2008. 16(6): p. 4309-4315.
- 91. Biberman, A., et al., *Ultralow-loss silicon ring resonators*. Optics Letters, 2012. **37**(20): p. 4236-4238.
- 92. *COMSOL 5.1*. Available from: https://www.comsol.com.
- Atabaki, A.H., et al., Optimization of metallic microheaters for high-speed reconfigurable silicon photonics. Optics Express, 2010. 18(17): p. 18312-18323.

- 94. Nishi, H., et al., *Integration of Eight-Channel Directly Modulated Membrane-Laser Array and SiN AWG Multiplexer on Si.* Journal of Lightwave Technology, 2019. **37**(2): p. 266-273.
- 95. Stern, B., et al., Battery-operated integrated frequency comb generator. Nature, 2018. **562**(7727): p. 401-405.
- 96. Witzens, J., *High-Speed Silicon Photonics Modulators*. Proceedings of the IEEE, 2018. **106**(12): p. 2158-2182.
- 97. Chen, Y.-H., J. Emer, and V. Sze, Eyeriss: a spatial architecture for energy-efficient dataflow for convolutional neural networks. SIGARCH Comput. Archit. News, 2016. 44(3): p. 367-379.
- 98. Chorowski, J. and J.M. Zurada, *Learning Understandable Neural Networks With Nonnegative Weight Constraints*. IEEE Transactions on Neural Networks and Learning Systems, 2015. **26**(1): p. 62-69.
- 99. Flenner, J. and B. Hunter, A Deep Non-Negative Matrix Factorization Neural Network. 2017.
- 100. Lecun, Y., et al., *Gradient-based learning applied to document recognition*. Proceedings of the IEEE, 1998. **86**(11): p. 2278-2324.
- 101. Kingma, D.P. and J. Ba, Adam: A method for stochastic optimization. arXiv preprint arXiv, 2014.
- 102. Pisati, M., et al. 6.3 A Sub-250mW 1-to-56Gb/s Continuous-Range PAM-4 42.5dB IL ADC/DAC-Based Transceiver in 7nm FinFET. in 2019 IEEE International Solid- State Circuits Conference (ISSCC). 2019.
- 103. Liu, D., et al., Submicron-resonator-based add-drop optical filter with an ultra-large free spectral range. Optics Express, 2019. 27(2): p. 416-422.
- 104. Eid, N., et al., FSR-free silicon-on-insulator microring resonator based filter with bent contra-directional couplers. Optics Express, 2016. 24(25): p. 29009-29021.
- Wu, X., et al., Low Power Consumption VOA Array With Air Trenches and Curved Waveguide. IEEE Photonics Journal, 2018. 10(2): p. 1-8.
- 106. Masood, A., et al. Comparison of heater architectures for thermal control of silicon photonic circuits. in 10th International Conference on Group IV Photonics. 2013.
- 107. Phare, C.T., et al., *Graphene electro-optic modulator with 30 GHz bandwidth*. Nature Photonics, 2015. **9**: p. 511.
- 108. Wright, C.D., et al. Integrated Phase-change Photonics: A Strategy for Merging Communication and Computing. in Optical Fiber Communication Conference (OFC) 2019. San Diego, California: Optical Society of America.
- 109. Wang, Q., et al., Optically reconfigurable metasurfaces and photonic devices based on phase change materials. Nature Photonics, 2015. 10: p. 60.
- 110. Li, X., et al., Fast and reliable storage using a 5 bit, nonvolatile photonic memory cell. Optica, 2019. 6(1): p. 1-6.
- 111. Teng, J., et al., Athermal Silicon-on-insulator ring resonators by overlaying a polymer cladding on narrowed waveguides. Optics Express, 2009. 17(17): p. 14627-14633.
- 112. Feng, S., et al., Athermal silicon ring resonators clad with titanium dioxide for 1.3µm wavelength operation. Optics Express, 2015. 23(20): p. 25653-25660.
- 113. Guha, B., J. Cardenas, and M. Lipson, *Athermal silicon microring resonators with titanium oxide cladding*. Optics Express, 2013. **21**(22): p. 26557-26563.
- 114. Lau, J.H., Recent Advances and New Trends in Flip Chip Technology. Journal of Electronic Packaging, 2016. 138(3): p. 030802-030802-23.
- George, J.K., et al., Neuromorphic photonics with electroabsorption modulators. Optics Express, 2019. 27(4): p. 5181-5191.
- 116. Bao, Q., et al., Monolayer graphene as a saturable absorber in a mode-locked laser. Nano Research, 2011. 4(3): p. 297-307.
- Schirmer, R.W. and A.L. Gaeta, Nonlinear mirror based on twophoton absorption. Journal of the Optical Society of America B, 1997. 14(11): p. 2865-2868.
- 118. Soljačić, M., et al., *Optimal bistable switching in nonlinear photonic crystals.* Physical Review E, 2002. **66**(5): p. 055601.
- Xu, B. and N.-B. Ming, Experimental observations of bistability and instability in a two-dimensional nonlinear optical superlattice.
   Physical Review Letters, 1993. 71(24): p. 3959-3962.
- 120. Coarer, F.D., et al., *All-Optical Reservoir Computing on a Photonic Chip Using Silicon-Based Ring Resonators*. IEEE Journal of Selected Topics in Quantum Electronics, 2018. **24**(6): p. 1-8.

- Miscuglio, M., et al., All-optical nonlinear activation function for photonic neural networks [Invited]. Optical Materials Express, 2018. 8(12): p. 3851-3863.
- Hamerly, R., et al., Large-Scale Optical Neural Networks Based on Photoelectric Multiplication. Physical Review X, 2019. 9(2): p. 021032.
- Cohen, N., O. Sharir, and A. Shashua. On the expressive power of deep learning: A tensor analysis. in Conference on Learning Theory. 2016.
- 124. Lu, Z., et al. The expressive power of neural networks: A view from the width. in Advances in Neural Information Processing Systems. 2017.

**Qixiang Cheng** received his B.S from Huazhong University of Sci. & Tech., China in 2010 and Ph.D. from the University of Cambridge, UK in 2014. In March 2015, he joined Shannon Lab., Huawei researching future optical computing systems.

He is now a research scientist at the Lightwave Research Lab, Columbia University, New York, USA. His current research interests include design, simulation, and characterization of large-scale optical integrated devices for data center and optical computing applications.

**Jihye Kwon** received the B.S. degree in mathematical sciences and the M.S. degree in computer science and engineering from Seoul National University, Seoul, Korea, in 2012 and 2014, respectively. She is currently pursuing the Ph.D. degree in computer science at Columbia University, New York, NY, USA.

Her research interest includes system-level design methodology enhanced via learning and interaction, high-level synthesis for designing hardware accelerators, real-time scheduling theory, and solving optimization problems. She has interned at IBM T. J. Watson Research Center during the summer in 2015, 2016, and 2017. During her PhD studies, she has received Presidential Fellowship from Columbia University.

Madeleine Glick received her Ph.D. in physics at Columbia University for research on electro-optic effects GaAs/AlGaAs quantum wells. After receiving her degree, she joined the Department of Physics, Ecole Polytechnique Federale de Lausanne (EPFL) Lausanne, Switzerland, where she continued her research in electro-optic effects in GaAs and InP-based materials. From 1992 to 1996, she was a Research Associate with CERN, Geneva, Switzerland, as part of the Lightwave Links for Analogue Signal Transfer Project for the Large Hadron Collider. From 2002-2011, Madeleine was Principal Engineer at Intel (Intel Research Cambridge UK, Intel Research Pittsburgh) leading research on optical interconnects for computer systems. Her research interests are in applying photonic devices and interconnects to computing systems. Madeleine is a Fellow of the Institute of Physics, and a Senior Member of IEEE and OSA.

**Meisam Bahadori** received his B.Sc. degree in electrical engineering, majoring in Communication Systems, with honors from Sharif University of Technology in 2011. After that, he worked toward M.Sc. degree in electrical engineering, majoring in Microwaves and Optics, at the same school and graduated with the highest honors in June 2013.

From fall 2011 to spring 2014, he worked as a research assistant at the Integrated Photonics Laboratory at Sharif University of Technology. He joined the Lightwave Research Laboratory at Columbia University in fall 2014 where he obtained the PhD degree in Electrical Engineering in 2018 with a focus on Silicon Photonics. His current research interests include silicon photonic devices, thin-film lithium niobate photonics, and nano-photonics.

Luca P. Carloni (S'95-M'04-SM'09-F'17) received the Laurea degree (summa cum laude) in electrical engineering from the Universita' di Bologna, Bologna, Italy, in 1995 and the M.S. and Ph.D. degrees in electrical engineering and computer sciences from the University of California at Berkeley, Berkeley, CA, USA, in 1997 and 2004, respectively. He is a Professor of computer science with Columbia University, New York, NY, USA. He has authored over 130 publications and holds two patents. His current research interests include systemon-chip platforms, system-level design, distributed embedded systems, and high-performance computer systems. Dr. Carloni was a recipient of the Demetri Angelakos Memorial Achievement Award in 2002, the Faculty Early Career Development (CAREER) Award from the National Science Foundation in 2006, the ONR Young Investigator Award in 2010, and the IEEE CEDA Early Career Award in 2012. He was selected as an Alfred P. Sloan Research Fellow in 2008. His paper on the latency-insensitive design methodology was selected for the Best of ICCAD in 1999, a collection of the best papers published in the first 20 years of the IEEE International Conference on Computer-Aided Design. In 2013, he served as the General Chair of Embedded Systems Week, the premier event covering all aspects of embedded systems and software.

Keren Bergman (S'87–M'93–SM'07–F'09) received the B.S. degree from Bucknell University, Lewisburg, PA, in 1988, and the M.S. and Ph.D. degrees from the Massachusetts Institute of Technology, Cambridge, in 1991 and 1994, respectively, all in electrical engineering. Dr. Bergman is currently a Charles Batchelor Professor at Columbia University, New York, NY, where she also directs the Lightwave Research Laboratory. She leads multiple research programs on optical interconnection networks for advanced computing systems, data centers, optical packet switched routers, and chip multiprocessor nanophotonic networks-on-chip. Dr. Bergman is a Fellow of the IEEE and OSA.