Analysis of Optical Loss and Crosstalk Noise in MZI-based Coherent Photonic Neural Networks

Amin Shafiee¹, Student Member, IEEE, Sanmitra Banerjee², Student Member, IEEE, Krishnendu Chakrabarty³, Fellow, IEEE, Sudeep Pasricha¹, Fellow, IEEE, and Mahdi Nikdast¹, Senior Member, IEEE

¹Department of Electrical and Computer Engineering, Colorado State University, Fort Collins, CO 80523, USA

²NVIDIA, Santa Clara, CA 95134, USA

Abstract-With the continuous increase in the size and complexity of machine learning models, the need for specialized hardware to efficiently run such models is rapidly growing. To address such a need, silicon-photonic-based neural network (SP-NN) accelerators have recently emerged as a promising alternative to electronic accelerators due to their lower latency and higher energy efficiency. Not only can SP-NNs alleviate the fan-in and fan-out problem with linear algebra processors, their operational bandwidth can match that of the photodetection rate (typically ≈100 GHz), which is at least over an order of magnitude faster than electronic counterparts that are restricted to a clock rate of a few GHz. Unfortunately, the underlying silicon photonic devices in SP-NNs suffer from inherent optical losses and crosstalk noise originating from fabrication imperfections and undesired optical couplings, the impact of which accumulates as the network scales up. Consequently, the inferencing accuracy in an SP-NN can be affected by such inefficiencies-e.g., can drop to below 10%the impact of which is yet to be fully studied. In this paper, we comprehensively model the optical loss and crosstalk noise using a bottom-up approach, from the device to the system level, in coherent SP-NNs built using Mach-Zehnder interferometer (MZI) devices. The proposed models can be applied to any SP-NN architecture with different configurations to analyze the effect of loss and crosstalk. Such an analysis is important where there are inferencing accuracy and scalability requirements to meet when designing an SP-NN. Using the proposed analytical framework, we show a catastrophic inferencing accuracy drop of up to 84% for SP-NNs of different scales with three known MZI mesh configurations (i.e., Reck, Clements, and Diamond) due to accumulated optical loss and crosstalk noise.

Index Terms—Silicon photonic integrated circuits, deep learning, optical neural networks, optical loss, optical crosstalk noise.

I. INTRODUCTION

With the rising demand for larger neural networks to address complex and computationally expensive problems, artificial intelligence (AI) accelerators need to consistently deliver better performance and improved accuracy while being energy-efficient. In this context, deep neural networks have attracted substantial interest for various applications ranging from image recognition to network anomaly detection, decision-making problems, self-driving cars, pandemic rate prediction, and early-stage cancer detection [1]. Given the significant growth in the demand for data-driven and computationally expensive applications, the energy efficiency of electronic-based deep-learning inference accelerators has been relatively low, and

This work was supported in part by the National Science Foundation (NSF) under grant numbers CCF-2006788 and CNS-2046226.

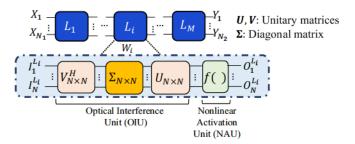


Fig. 1: Overview of a coherent SP-NN with N_1 inputs, N_2 outputs, and M layers.

they have been unable to keep up with the performance requirements of emerging deep-learning applications [2], [3].

Silicon photonics (SiPh) enables the deployment of integrated photonics across a wide range of applications, from realizing ultra-fast communication for Datacom applications [4]-[7] to energy-efficient optical computation in emerging hardware accelerators for deep learning [8]-[10]. To alleviate the limitations of conventional CMOS-based electronic accelerators in terms of energy consumption and latency, new SiPh-based hardware accelerators optimized for deep learning applications are on the rise. By leveraging optical interconnects for communication and photonic devices for computation, silicon-photonic-based neural network (SP-NN) accelerators offer the promise of up to 1000 times higher energy efficiency for performing computationally expensive multiply-and-accumulate operations [2], which are the most power-hungry and common operations in deep learning applications [2].

Among different SP-NN implementations, coherent SP-NNs, which operate on a single wavelength, have an inherent advantage over noncoherent SP-NNs that require power-hungry wavelength-conversion steps and multiple wavelength sources [11]. Fig. 1 presents an overview of a multi-layer coherent SP-NN with N_1 inputs, N_2 outputs, and M layers. Each layer comprises an optical-interference unit (OIU) implemented using an array of Mach–Zehnder interferometers (MZIs) with a specific architecture, connected to a nonlinear-activation unit (NAU). Within an OIU, MZIs can be used to realize matrix-vector multiplication as shown by [2], [12]. Accordingly, several coherent SP-NN architectures have been proposed by cascading arrays of MZIs to perform large-scale

³ Department of Electrical, Computer and Energy Engineering, Arizona State University, Tempe, AZ 85281, USA

linear multiplication in the optical domain [13]–[15].

While SP-NNs are promising alternatives to electronicbased deep learning hardware accelerators, several factors limit their performance and scalability. The underlying devices in SP-NNs (e.g., MZIs in coherent SP-NNs) suffer from optical loss and crosstalk noise due to device fabrication imperfections (e.g., sidewall roughness) and undesired mode couplings [8], [16]. For example, prior work has shown up to 1.5 dB insertion loss and 18 dB crosstalk in a 2 2 MZI [17]. While optical loss and crosstalk noise are small and seem to be negligible at the device level, they accumulate as the network scales up, hence leading to severe performance degradation at the network and system level (e.g., drop in inferencing accuracy). Moreover, crosstalk cannot be filtered in coherent SP-NNsour focus in this paper—due to the coherence between the noise and victim signals. Therefore, there is a critical need for careful analysis of optical loss and crosstalk noise in coherent SP-NNs and exploring their impact on SP-NN performance.

The novel contribution of this paper is to comprehensively analyze optical loss and crosstalk noise and their impact on coherent SP-NN performance from the device level to the system level. We develop a realistic device-level MZI compact model to analyze the optical loss from different sources (i.e., propagation loss, directional coupler loss, and metal absorption loss) and the coherent crosstalk noise in the MZI. This model is also able to capture the impact of optical phase settings, which represent weight parameters in coherent SP-NNs, on the MZI's optical loss and crosstalk noise. Leveraging our accurate device-level models, we present layer- and networklevel optical loss and coherent crosstalk models that scale with the number of inputs and layers in coherent SP-NNs. In addition, we propose a detailed analysis of the effect of optical loss and crosstalk in SP-NNs when the optoelectronic NAU units are used. The proposed framework enables an accurate exploration of the inferencing accuracy drop in SP-NNs with different mesh configurations under the effect of optical loss and crosstalk noise. Leveraging our proposed framework, we also quantify the maximum optical loss acceptable in the underlying devices when specific inferencing accuracy goals must be met within an SP-NN.

The proposed analytical framework can be applied to any coherent SP-NN architecture of any size to analyze the optical loss and crosstalk noise in the network. In this paper, we consider three well-known MZI-based coherent SP-NN architectures, namely Clements [14], Reck [15], and Diamond [13]. For example, for the case study of the Clements SP-NN with 16 inputs and 2 hidden layers, we show that the optical loss can be as high as 38 dB and the signal-to-noise ratio (SNR) can be as low as 7.3 dB. Furthermore, considering the MNIST classification task as an example, we show that the network inferencing accuracy can drop by about 84.6% due to optical loss and crosstalk noise. We also show that by increasing the number of inputs from 16 to 64 in the same network, the resulting optical loss increases significantly to as high as 140 dB and the SNR decreases to as low as 1 dB. Moreover, we also demonstrate that increasing the input optical power to compensate for the effect of accumulated optical loss and crosstalk will be an impractical solution as it also increases the noise at the outputs of the SP-NNs while maintaining the same SNR. The proposed analyses in this paper extend our prior work in [12] by performing the loss and crosstalk analysis for two more SP-NN configurations, analyzing the effect of optical loss and crosstalk in SP-NNs when optoelectronic nonlinear activation units are used, presenting SNR and optical power penalty models for SP-NNs, and analysing the scalability constraints due to the optical loss and crosstalk in the MZI-based SP-NNs when being used as a photonic processing unit for high-performance computation.

The rest of the paper is organized as follows. Section II presents an overview of the building blocks in SP-NNs, SP-NN design and working mechanism, and prior related work. Section III presents analytical models to analyze the impact of optical loss and crosstalk from the device level to the system level in SP-NNs. The impact of the optical loss and crosstalk in optoelectronic NAU units is modeled in this section. Section IV presents the simulation results to show the impact of loss and crosstalk on the performance of SP-NNs with the three MZI mesh configurations of Clements, Reck, and Diamond. Section V presents the discussion on the effect of optical loss and crosstalk noise on SP-NN power consumption (i.e., laser power penalty) as well as scalability constraints in SP-NNs. Finally, Section VI concludes this work.

II. BACKGROUND AND PRIOR RELATED WORK

In this section, we present an overview of the MZI building block in coherent SP-NNs as the primary vector-matrix multiplier unit and some fundamentals of MZI-based coherent SP-NNs. We also discuss different sources of optical loss and crosstalk in MZIs. Moreover, we review prior work on studying the effect of loss and crosstalk in SP-NNs.

A. Mach-Zehnder Interferometer (MZI)

MZIs can be used to realize linear multiplication between a 2 1 vector (signals applied to the two inputs) and a 2 2 matrix (defined based on the phase settings in the MZI). Such an MZI-based multiplier unit can be constructed using two 3dB directional couplers (DCs) with an ideal splitting ratio of 50:50 and two integrated phase shifters (and), as shown in Fig. 2. Phase shifters in this design can be implemented, for example, using microheaters on top of the underlying waveguide [17]. By introducing a temperature change using microheaters, the refractive index of the underlying silicon waveguide will change due to the thermo-optic effect, leading to a change in the phase of the electric field of the propagating optical signal. Therefore, by controlling the phase shift between the two arms in an MZI, we can control the interference in the output. Note that in the MZI in Fig. 2, . The transfer matrix of an MZI-based multiplier

. The transfer matrix of an MZI-based multiplier unit can be realized by multiplying the transfer matrices of the two 3-dB DCs () and the transfer matrices of the two phase shifters (). Accordingly, the ideal transfer

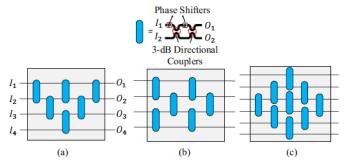


Fig. 2: Schematic of a 4×4 single MZI-based optical interface unit (top) with different mesh configurations. (a) Reck, (b) Clements, and (c) Diamond.

matrix of an MZI-based multiplier unit (i.e., without optical loss and crosstalk noise) can be defined as [13], [18]:

$$T_{MZI}(\theta, \phi) = T_{DC2} \cdot T_{\theta} \cdot T_{DC1} \cdot T_{\phi}$$

$$= \begin{pmatrix} \frac{e^{i\phi}}{2} (e^{i\theta} - 1) & \frac{i}{2} (e^{i\theta} + 1) \\ \frac{ie^{i\phi}}{2} (e^{i\theta} + 1) & -\frac{1}{2} (e^{i\theta} - 1) \end{pmatrix}. \quad (1$$

While MZIs can help perform matrix-vector multiplication in the optical domain, they have a large footprint. For example, state-of-the-art MZI multipliers can be up to about 300 μ m long and 30 μ m wide [8], limiting the scalability of MZI-based coherent SP-NNs [17], [19]. In addition to a large footprint, they also suffer from high optical losses and crosstalk noise [20]. The optical loss in an MZI originates from the absorption in the metallic contacts in proximity to the waveguides when using microheaters for applying the required phase shifts on the MZI arms, the directional coupler imperfections, and the propagation loss of the waveguides, which is mainly due to their sidewall roughness. The optical loss and crosstalk noise in MZIs will be discussed in detail in Section III.

B. MZI-based Coherent SP-NNs

MZI-based SP-NNs rely on the manipulation of the electrical field's phase of a single optical wavelength to perform matrix-vector multiplication. MZIs in such SP-NNs are responsible for phase manipulations and interference to carry out the computations [11], [13]-[15], [17]. MZIs can be cascaded in the form of an array following a specific configuration to implement an OIU for performing large matrix-vector multiplication in the optical domain. Fig. 1 shows an overview of a coherent SP-NN composed of an optical interface unit (OIU) and nonlinear activation unit (NAU). A fully connected layer (L_i) with n inputs performs matrix-vector multiplication between the input vector (I_i) and a weight matrix (W_i) . Note that to realize unitary matrices with arbitrary magnitudes, we can scale the input optical power when inferencing [2], [13], [21]. The output vector of the OIU eventually will be passed into the NAU to apply a nonlinear activation function (f_i) , the result of which will be the input to the next layer (L_{i+1}) . The output of L_i can be mathematically modeled as $O_i^{n_i \times 1} = f_i(W_i^{n_i \times n_{i-1}}, O_{i-1}^{n_{i-1} \times 1})$, in which W_i is the layer's corresponding weight matrix. The weight matrix (W_i) can be

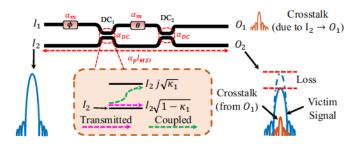


Fig. 3: Schematic of a 2×2 MZI multiplier unit with different sources of optical loss and crosstalk noise. Here, $I_2\to O_2$ is shown as an example with $\theta=\pi$ (l_{MZI} : MZI length, α_m : loss due to metallic absorption, α_{DC} : DC's insertion loss, α_p : propagation loss).

obtained by training the network. Using singular value decomposition (SVD), the obtained weight matrix can be mathematically modeled as $W_i = U_i^{n_i \times n_i} \Sigma_i^{n_i \times n_{i-1}} V_i^{H,n_{i-1} \times n_{i-1}}$. In this formulation, U_i and V_i are unitary matrices. Moreover, V_i^H denotes the Hermitian transpose of V_i and Σ_i is a diagonal matrix consisting of the eigenvalues of W_i . A unitary matrix can be implemented by an array of cascaded 2×2 MZIs in a specific configuration according to:

$$U_i^{n_i \times n_i} = D\left(\prod_{(m,n) \in S} T_{MZI_{m,n}}\right). \tag{2}$$

Using this scheme, each unitary matrix U can be decomposed into the products of several MZIs' transfer matrices. The order of the multiplication of MZIs' transfer matrices plays an important role in SP-NNs, determining the configuration of the MZIs in the SP-NNs (i.e. Clements, Reck, or Diamond see Fig. 2). In (2), D is a diagonal matrix with complex elements with a unity modulus [14], and S denotes the order of the multiplication of the MZIs' transfer matrices. S will be determined based on the configuration of the array of cascaded MZIs used to map the weight matrices in order to perform matrix-vector multiplication in the optical domain. Moreover, m and n denote the input ports (i.e., $I_1 - I_4$ in Fig.2 (a)) which require transformation using each MZI in the network. For example, m = 1 and n = 2 refer to the MZI in between input ports 1 and 2 in the network. Note that n = m+1 always applies [14]. The location of each MZI in the network can be determined during the mapping of the weights to the array of cascaded MZIs which itself depends on the configuration of the network.

Several configurations have been proposed for the network topology of cascaded MZIs and to carry out the unitary transformation of the input optical signals. Three well-known mesh configurations: Reck, Clements, and Diamond [13]–[15] are considered in this paper to realize MZI-based unitary multipliers. A 4×4 Reck mesh configuration is shown in Fig. 2(a) organizing an array of MZIs in a triangular shape. In general, any $N \times N$ unitary multiplier based on the Reck design consists of $\frac{N(N-1)}{2}$ MZIs, where N is the number of the input ports. The same number of MZIs can be configured in a rectangular shape as is shown in Fig. 2(b), and this

configuration is called the Clements mesh. The advantage of the Clements design over the Reck is that the network is more symmetric, hence making the unitary multiplier more resilient to propagation loss due to more symmetric and, on average, shorter optical paths compared to the Reck design [14]. We can also increase the number of MZIs to make the Reck design more symmetric. The work in [13] proposed this symmetric design by adding additional - MZIs to the Reck configuration to design the Diamond mesh configuration, shown in Fig. 2(c). For an array of cascaded MZIs with MZIs will be used in a diamond shape. Although the network topology includes a higher number of input and output ports, only the last inputs are used to perform matrix-vector multiplication, while the rest of the inputs can be used for characterization and calibration of the MZIs in the network [13].

C. Optical Loss and Crosstalk Noise in MZIs

MZIs intrinsically suffer from optical loss and crosstalk noise. The schematic of a single 2 2 MZI multiplier is shown in Fig. 3. An optical signal traversing an MZI can undergo different losses based on the MZI phase settings. The main sources of optical loss in an MZI are the DCs, metal absorption due to the proximity of microheater's metallic contacts to waveguides, and propagation loss mainly due to sidewall roughness in waveguides [22]. Note that the metallic absorption loss varies slightly with the adjusted phase settings, corresponding to different temperatures, in the MZI. The reason for this slight change (i.e., 1 - 3%) in the loss is due to slight changes in the propagation constant at different temperatures due to the thermo-optic effect, which also affect the optical mode confinement and light-matter interaction in the waveguide. In addition, crosstalk noise can originate from the undesired coupling of light in the DCs in an MZI due to an imbalance in the DC's splitting ratio. In coherent SP-NNs, we deal with coherent (i.e., intra-channel or in-band) crosstalk. As shown in Fig. 3, coherent crosstalk noise can interfere with the main optical signal (victim signal) based on their phase difference, imposing power fluctuations on the victim signal. Unlike incoherent networks, in coherent networks like MZIbased SP-NNs, the in-band coherent crosstalk noise cannot be easily filtered in the output due to the coherence between the crosstalk and victim signals (i.e., on the same wavelength).

Unlike in SP-NNs, optical loss and crosstalk noise have been widely studied in chip-scale Datacom photonic networks (e.g., [23] and [24]), showing signal integrity degradation and scalability constraints in these networks due to optical loss and crosstalk noise. Unfortunately, the existing work on optical loss and crosstalk analysis in such networks cannot be applied to SP-NNs because optical loss and crosstalk noise characteristics of silicon photonic devices for photonic computation in SP-NNs are different. For example, a 2–2 MZI switching cell, whose structure is similar to the one in Fig. 3 but without , in an optical switch fabric can only take two functional states based on for optical loss and crosstalk analysis: the Cross-state, where

, and the Bar-state, where and

and . As a result, crosstalk noise can be easily characterized in such a device. However, in the coherent SP-NNs, , which determines the MZI state, can take any value between 0 and (0) to perform computation via interference between the inputs. The analysis of optical loss and crosstalk in SP-NNs should therefore account for various phase settings in the underlying MZI devices.

D. Prior Related Work

While several coherent SP-NNs have been recently proposed [13], [25], the work in [8], [26] showed that the inferencing accuracy of such networks can drop by up to 90% due to fabrication-process variations and thermal crosstalk. In addition to such variations, the work in [13] explored the impact of optical loss and phase noise in MZIs for the Diamond SP-NN configuration and showed that the SP-NN's inferencing accuracy can drop to below 20% when scaling the network in the presence of phase errors and MZI losses. However, the work in [13] simply assumed that the loss in different devices in a network to be normally distributed (irrespective of their phase settings).

The work in [27] systematically analyzed the impact of loss and crosstalk in SP-NNs. However, the loss and crosstalk models presented in [27] are not a function of the phase settings of the MZIs. Hence, only the maximum and minimum loss and crosstalk are calculated. In addition, the models illustrated in [27] cannot be extended to OIUs of any arbitrary configuration. As the size and complexity of emerging SP-NNs increase to handle more complex tasks, the total insertion loss accumulated in the network increases as well. This necessitates the use of power-hungry optical amplification devices [28] and higher laser power at the input. Uncertainties due to fabrication-process variations—the analysis of which is beyond the scope of this paper—in the two DCs in an MZI can degrade the extinction ratio (ER) of the device which, in turn, will increase the loss and crosstalk in the output [29]. Yet, no prior work comprehensively analyzes the impact of optical loss and crosstalk noise in coherent SP-NNs from the device to the system level. Although the work in [29] suggests that using silicon nitride instead of silicon to implement MZI-based SP-NNs as a possible solution to reduce losses, the performance degradation due to coherent crosstalk in SP-NNs still remains unaddressed.

Different from the aforementioned work, this paper presents a comprehensive modeling framework for the optical loss and coherent crosstalk noise in coherent SP-NNs. The proposed loss and crosstalk analysis at the MZI device level takes into account the phase settings on the device. In addition, the models developed at the network level are adaptable meaning that the number of inputs/outputs and layers in SP-NNs can vary to explore SP-NN optical power penalty and scalability constraints while evaluating average and worst-case optical loss and crosstalk. Compared to our prior work in [12], we extended our analysis for two more well-known SP-NN configurations (Reck and Diamond) in addition to the Clements configurations. We show that our proposed models can be applied to any SP-NN with any configuration, proving

their versatility. We present a comprehensive analysis of the laser power penalty and inferencing accuracy loss in SP-NNs due to optical loss and crosstalk noise. In addition, we present an analysis of the effect of optical loss and crosstalk when optoelectronic NAUs are used in the SP-NNs. We also present a detailed analysis of scalability constraints of SP-NNs due to optical loss and crosstalk when a single OIU is being used as a processing unit to carry out the computations in the optical domain.

III. OPTICAL LOSS AND CROSSTALK NOISE ANALYSIS

This section presents the compact models developed to analyze optical loss and crosstalk noise in MZI-based coherent SP-NNs from the device level to the network level. We also model the loss and crosstalk noise in the optoelectronic nonlinear activation unit in coherent SP-NNs. All these models will be used to explore the power penalty, performance (e.g., inferencing accuracy), and scalability constraints in SP-NNs of different sizes and architectures under optical loss and crosstalk noise, as will be discussed in Section IV.

A. Modeling Optical Loss and Crosstalk at Device Level

The schematic of a single 2 2 MZI-based multiplier unit is shown in Fig. 3. It comprises two 3-dB DCs and two integrated phase shifters (and) on the upper arm. DCs in the MZI structure are responsible for splitting and combining the optical signals traversing the MZI. Optical crosstalk noise stems from undesired mode coupling in these DCs in the MZI structure. The splitting ratio for an optical signal at the input of a DC can be determined by the cross-over coupling coefficient () and transmission coefficient () in the DC. Considering the DC's total transmission at its outputs, encapsulates the DC's optical insertion loss as we can write . Note that for a lossless DC,

and for a DC including its optical insertion loss, To be used in an MZI-based multiplier unit for coherent SP-NNs, the DCs should perform exact 50:50 splitting/combining —). Both and are wavelength dependent and they also depend on the waveguide width and thickness and the gap between the waveguides in DCs. As discussed in Section II, the main sources of optical loss in an MZI are the DC's loss (), propagation loss (), and the metal) due to interaction of the propagating light absorption loss (with the metallic contacts integrated on top of the waveguide when using microheaters to implement the phase shifters. The propagation loss originates from sidewall roughness and scattering loss in the waveguides [16]. The metal absorption loss depends on the metal used to implement the heater's contacts and their geometry such as thickness, width, length, and the longitudinal distance between the waveguide and the contacts [30].

Considering (1) and the MZI's optical loss parameters defined (i.e., , , and , as shown in Fig. 3), an optical-loss-aware transfer matrix model for the 2 2 MZI multiplier unit can be defined as:

(3.e)

In this model, and are the cross-over coupling coefficients in the two DCs and represents the total propagation loss related to the MZI.

Optical crosstalk noise in a 2 2 MZI switching element, whose structure is similar to the one shown in Fig. 3 but without phase shifter at the input of the MZI, can be analyzed by considering the MZI in the Bar (in the Cross-state (0). By injecting an optical signal into one of the input ports of the MZI, the crosstalk noise can be determined by capturing the undesired optical signal transmission on the opposite output when in Cross- or Barstates. For example, considering the MZI schematic depicted in Fig. 3, when injecting light into , the , by setting crosstalk noise in the Bar-state () can be captured on Similarly, it is possible to capture the Cross-state crosstalk by setting 0. Different from the MZIs used in switching networks, MZIs in coherent SP-NNs can take more than the two Bar-state and Cross-state, depending on the value of . In a 2 2 MZI-based multiplier unit, can assume any value between 0 and) to perform a 2 2 unitary multiplication (does not change the MZI state). Therefore, the exact analysis of the crosstalk noise per device's output ports cannot be easily performed in coherent SP-NNs.

To address this limitation, we can define a statistical model to estimate the crosstalk noise on the output ports of each MZI in coherent SP-NNs depending on the MZI's setting. Considering the crosstalk noise in the two known Cross-state and Bar- state (and , where typically [31]), we can statistically model the crosstalk noise in the MZI in any intermediate state () as the function of using a Gaussian distribution whose mean can be calculated according to and standard deviation of 0.05 , considered here as an example. Note that the Gaussian distribution for the crosstalk noise model in the MZI in intermediate states is assumed as an example, and it is similar to the assumptions made for the noise distribution in optical communication and computation systems [32]. Using the crosstalk noise in intermediate states , the transfer matrix of a 2 2 MZI-based multiplier unit in (3) can be written as:

$$(3) (4)$$

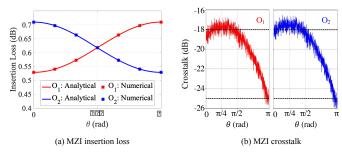


Fig. 4: (a) Insertion loss and (b) crosstalk noise at the outputs of the 2 2 MZI in Fig. 3 simulated using the parameters listed in Table I.

Models presented in (3) and (4) can be used to analyze the optical loss and crosstalk noise in any MZI with different phase settings. Considering the parameters listed in Table I, Fig. 4 shows the insertion loss and crosstalk for a single 2 2 MZI multiplier unit. In this figure, the x-axis shows the state of the MZI based on where

It should be noted that does not change the MZI state, but its loss is included in the results shown. Observe that the insertion loss on each output port is 0.52–0.72 dB. we used Ansys Lumerical [33] to validate the results in Fig. 4(a). Note that commercial tools (e.g., Lumerical or Synopsys) cannot analyze the crosstalk in intermediate states in the MZI, hence their results are not considered in Fig. 4(b). We can see from Fig. 4 that by changing the state of the MZI, the device-level loss and crosstalk noise also change. In Fig. 4(a), compared to input, the optical signal on experiences higher insertion (see Fig. 3). Therefore, for loss because of through example, the insertion loss is higher on for the Cross-state when and it is higher on for the Bar-state when , the MZI in fact acts like a 50:50 . As for splitter where the loss at the outputs is the same.

It should be mentioned that the fluctuations in the crosstalk in Fig. 4(b) are due to the Gaussian noise model defined for the MZI. The coherent crosstalk in the MZI output changes between 18 dBm and 25 dBm, when the input power is 0 dBm.

B. Modeling Optical Loss and Crosstalk at Network Level

As discussed in Section II-B, the weight matrix of each hidden layer in an SP-NN architecture with inputs and outputs can be decomposed into a multiplication of two unitary and one diagonal matrix implemented by OIUs followed by an NAU, as shown in Fig. 1. The number of MZIs in an OIU depends on the mesh configuration (Clements, Reck, or Diamond) and the number of inputs () and outputs (). Moreover, as discussed in Section III-A, the optical loss and crosstalk noise at the outputs of an MZI in an OIU is a function of and in the MZI, which can be obtained during the training of the network [12].

The optical loss in the output of an SP-NN can be systematically modeled as:

In (5), is the insertion loss in the OIUs that can be calculated using (3) for each MZI, and it depends on the OIU's mesh configuration (), its dimension (and), and the phase settings of MZIs in the OIU. Moreover the is the insertion loss due to the NAU. Also, the insertion of the optoelectronic NAU is considered to be 0–1 dB, depending on its nonlinear response [25].

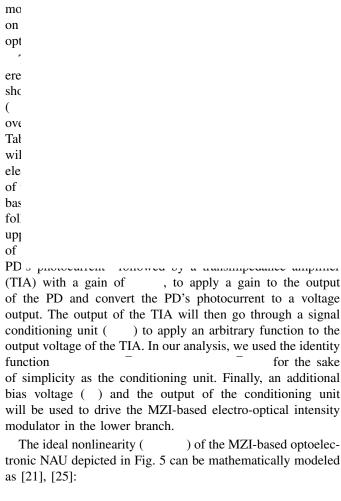
Optical crosstalk noise from each MZI will be accumulated at the outputs of the OIU as the optical signal propagates through the SP-NN. Therefore, following the same approach for optical loss calculation, the accumulated optical crosstalk noise power at the end of the SP-NN can be modeled as:

(6)

In (6), is the total number of MZIs in the network and it depends on the configuration of the OIUs as well as their dimension. Also, is the optical power at input . Moreover, can be calculated using (4) and is the coherent crosstalk on the output of the OIU originating from MZI in the network. Also, is the optical phase of the crosstalk signal. Similarly, is the insertion loss in the OIU, which can be calculated using (3) for each MZI, experienced by as it traverses the network towards output . Note that using an optical gain unit such as semiconductor optical amplifiers (SOAs) to compensate for the accumulated optical losses or scaling the outputs of the OIUs can be problematic. This is due to the fact that SOAs will also amplify the coherent crosstalk noise at the outputs of the OIUs. In addition, at low input powers due to accumulated optical losses, SOAs suffer from amplified spontaneous emission noise (ASE), resulting in high noise figures in the output, where the SOA's noise can dominate the main output signal and coherent crosstalk noise [34]–[36], hence further degrading the network performance. Note that to realize unitary matrices with arbitrary magnitudes without using SOAs, we can scale the input optical power as shown by [2], [13], [21].

C. Modeling Optical Loss and Crosstalk in Optoelectronic NAUs

Nonlinear activation functions are an integral part of deep neural networks due to the essential need for the realization of the complex nonlinear relationship between the inputs and outputs of the SP-NNs [25]. NAUs are responsible to trigger a single activation at the end of each layer's output and pass the output to the input of the next layer. NAUs can be implemented electronically [2], optoelectronically [25], or optically [37], each with different costs. High power consumption and latency as well as the need for lasers because of the need for multiple E-O and O-E conversions can be named as limitations related to electronically implemented NAUs. Moreover, very large waveguide lengths and high optical power must be used for optical NAUs due to the weak nonlinearity of photonic platforms [25], [38]. Optoelectronic NAUs presented in [21], [25] show great promise as alternatives to electronic ones due to the ability to implement an arbitrary nonlinear response via self-intensity modulation (e.g., MZI-based electro-optical



is the voltage that is required to impose a shift in the MZI-based intensity modulator in the optoelectronic NAU. A ReLU-like nonlinear activation response can be realized by setting in the formulation proposed in (7). Note that in the activation function modeled in (7) and [21], [25], the effect of optical loss and crosstalk noise from the OIU unit, PD's sensitivity, shot noise, dark current (output current of the PD in the absence of any optical signal), and optical insertion loss of the DC and MZI in the NAU architecture is not considered.

To update (7) for analyzing the effect of optical loss and crosstalk noise in SP-NNs when optoelectronic NAUs are used, we can replace in (7) with

. Note that in these formulations, the loss values are considered in terms of power amplification and attenuation coefficient. Moreover, takes into account the phase of the crosstalk noise at the output of the OIU unit. Using this approach, we can write the input of the optoelectronic NAU unit including the optical loss and crosstalk from OIU (

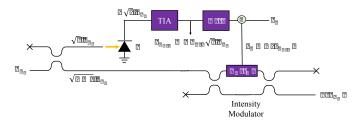


Fig. 5: Schematic block diagram of the optoelectronic nonlinear activation unit used to realize ReLU-like nonlinear activation function [21], [25].

as [39]: (8)

The output photocurrent of the PD with a responsivity of shot noise of and dark current of can be modeled as:

in which is the PD's bandwidth where C is the electronic charge [39]. The and output photocurrent then enters the TIA with the gain of and will be converted to an amplified voltage according to:

(10)

(9)

Keeping the assumption of we can systematically model the activation function response using an MZI-based optoelectronic intensity modulator mentioned in [21], [25] under optical loss and crosstalk noise from the OIU unit as:

(11)

where (12)

Here, we also assumed that to realize ReLU-like activation response and the output photocurrent of the PD is considered to be zero () for , where stands for PD sensitivity (the minimum detectable optical power by the PD). In (11), is the optical crosstalk noise related to the MZI-based electro-optical modulator, is its insertion loss, and is the phase of the optical crosstalk noise from the MZI in the intensity modulator in the optoelectronic NAU [21].

IV. SIMULATION RESULTS AND DISCUSSIONS

Optical loss and crosstalk noise lead to the deterioration of the performance of SP-NNs. We developed a framework to analyze the effect of optical loss and crosstalk noise in MZIbased coherent SP-NNs on top of Neuroptica [13], [42]. Neuroptica is a flexible chip-level simulation platform for nanophotonic neural networks written in Python/NumPy. It provides a wide range of abstraction levels for simulating

TABLE I: Device-level loss, crosstalk, power, gain, and NAU parameters considered in this paper (PhS: Phase shifter, PD: Photodetector).

| Par. | Definition | Value | Ref. |
|------|-----------------------------|-----------|------------------|
| | Crosstalk in Bar-state | -25 dB | [31] |
| | Crosstalk in Cross-state | -18 dB | [31] |
| | MZI length | 300 m | [17] |
| | PhS (metal) absorption loss | 0.23 dB | [30] |
| | Propagation loss | 2 dB/cm | [16] |
| | Insertion loss of DC | 0.13 dB | [16] |
| | Input optical power | 0 dBm | - |
| | TIA gain | | [25] |
| | Responsivity | 1 A/W | [21], [25], [40] |
| | DC's splitting ratio in NAU | 0.1 | [21], [25], [40] |
| | Intensity modulator voltage | 10 V | [21], [25], [40] |
| | NAU's bias voltage | 10 V | [21], [25], [40] |
| | PD's bandwidth | 42.5 GHz | [40] |
| | PD's dark current | 3.5 A | [40] |
| | Intensity modulator's loss | 1 dB | [25] |
| | PD's sensitivity | -11.7 dBm | [40], [41] |

optical neural networks [42]. We expanded the analysis of loss and crosstalk noise for a single MZI using the defined mathematical models in Section III and parameters listed in Table I to perform layer-level (i.e., OIU), network-level (i.e., multi-layer SP-NN), and system-level (i.e., network accuracy) analyses using Neuroptica [42].

For layer- and network-level analyses, we consider random phase settings for MZIs in OIUs of different dimensions and configurations (8, 16, 32, and 64, 1, 2, and 3 layers). SVD is used to obtain the corresponding weight matrix in an SP-NN (see Fig. 1) with three mesh configurations of Reck, Diamond, and Clements (see Fig. 2). Note that the random phases are only used for layer- and network-level analyses. As for the system-level analysis, we use shifted fast Fourier transform (shifted-FFT) on the MNIST handwritten digit dataset. The reason for using shifter-FFT is to reduce the number of inputs which leads to the size of the mesh configurations being smaller and more manageable so we can carry out the training of the SP-NNs of different configurations (i.e., Reck, Diamond, and Clements). Note that the training of the network is performed on complex inputs which leads to complex values of the weights when the training of the SP-NN is finished [43].

We use inferencing accuracy as a figure of merit to analyze the effect of optical loss and crosstalk noise at the system level in the SP-NNs using the ideal *ReLU* activation function. Moreover, we use relative-variation distance (RVD) (a measure of the deviation of two matrices [26]) for the scalability analysis of OIUs of different scales with different mesh configurations. Such an analysis is helpful when using standalone OIUs as a photonic multiplication unit to perform matrix-vector multiplication. We also present a comprehensive analysis of the effect of optical loss and crosstalk noise on the SP-NNs' performance when optoelectronic NAUs are used instead of the ideal *ReLU* activation response. In addition, we analyze the laser power penalty in SP-NNs with different scales and configurations to compensate for optical loss and crosstalk noise.

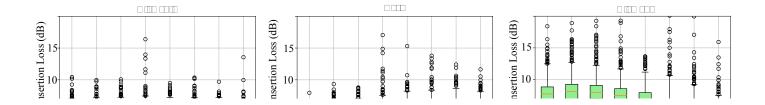
A. Optical Loss and Crosstalk in OIUs

Using the analytical models in (5) and (6) proposed in Section III, a single layer () with 8 (considered as an example), three different mesh configurations of Clements, Reck, and Diamond are simulated to capture the impact of optical loss and crosstalk noise at the SP-NN outputs. In these simulations, the optical insertion loss is calculated for 1000 random weight matrices and the results are shown in the form of conventional box plots in Fig. 6. We can see from Fig. 6(a) that for the Clements configuration, the average insertion loss among all the output ports is 6.5 dB, while the worst-case insertion loss can be as high as 16 dB. The average insertion loss for Reck and Diamond configurations is 6 dB and 6.5 dB, respectively. As for the worst-case loss, the Reck configuration experiences 16 dB, while the Diamond configuration undergoes 18 dB of optical loss. In addition, observe that the insertion loss for the Clements and Reck configuration is almost similar on all the outputs due to their more symmetric layer configuration compared to Diamond [2].

We also analyze optical crosstalk noise power for a single layer OIU with 8 across different mesh configurations. We used 1000 random weight matrices to perform statistical analysis of accumulated coherent crosstalk noise at each output in the OIU with Clements, Reck, and Diamond configurations. In each iteration, a random optical phase (), where , is assigned to the crosstalk noise signal from each MZI in the OIU structure to emulate the crosstalk noise signal phase and behavior throughout the network. Note that the crosstalk signal from each MZI will interact with each other and the victim signal at the outputs of the OIUs. This approach is acceptable when optical signals traverse a large network of devices (e.g., in OIUs), and hence experience random phase shifts. The crosstalk noise including the insertion loss is shown in Fig. 7 for each network. Note that 0 dBm input optical power at the input of the OIUs has been considered in these simulations. Observe that the average crosstalk noise for the Clements, Reck, and Diamond meshes can be as high as 21.3, 22.1, and 22.4 dB, respectively. As for the worst-case crosstalk noise, the three mesh configurations of Clements, Reck, and Diamond exhibit 15.3, 15.1, and 12.2 dB, respectively. Note that the reason for the higher worst-case crosstalk noise in the Diamond mesh is due to the larger number of MZIs in this network where each MZI generates some crosstalk that is accumulated on the outputs. In addition, note that although the Diamond mesh has a higher number of MZIs compared to Reck and Clements when being used in the OIU units, some of the output ports (are not used during inferencing and are reserved only for the characterization and calibration of the MZIs in the Diamond mesh [13].

B. Optical Loss and Crosstalk in Multi-Layer SP-NNs

The optical loss and crosstalk power can be analyzed in an SP-NN comprising multiple hidden layers and with different dimensions and configurations using the models developed in (5) and (6). We extended the layer-level insertion loss and crosstalk models in (5) and (6) for the full-network to



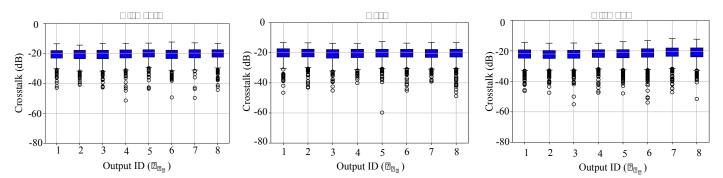


Fig. 7: Optical crosstalk analysis for 1000 random weight matrices mapped to a fully connected 8 8 OIU with different mesh configurations: (a) Clements, (b) Reck, and (c) Diamond. Output IDs are numbered from top to bottom.

demonstrate how the crosstalk and optical loss impact the integrity of optical signal (measured by SNR) as we scale up SP-NNs. We considered 1000 random weight matrices for SP-NNs with different dimensions (8, 16, 32, and 64) and numbers of layers (1, 2, and 3). For each matrix, we performed three different experiments on the same network. The first experiment included the ideal network with ideal output powers (i.e., without optical loss and crosstalk noise). The second experiment included the case where optical losses of the MZIs in the network were included. The optical loss reported in Figs. 8(a)-8(c) is calculated by comparing the network outputs from the first and the second experiments. As for the third experiment, we included both the optical loss of the components and accumulated crosstalk noise from them on the outputs. The outputs of the second and the third experiments were used to calculate the crosstalk at the outputs, and also the SNR.

The optical loss, crosstalk, output power, and SNR distribution among the outputs of MZI-based SP-NNs with different dimensions and configurations are depicted in the form of box plots in Figs. 8. Note that the average values over 1000 random weight matrices for each output are used for the results in Fig. 8. In addition, the insertion loss of the optoelectronic NAU is included in the results. As can be seen from Figs. 8(a)–8(c), insertion loss for SP-NNs increases drastically as and increase for all the three mesh configurations. Observe that the dynamic range in the output of the Diamond's network is larger compared to that in Clements and Reck, because of the asymmetric design of the network (i.e., asymmetric distribution of MZIs) in the Diamond network. Furthermore,

the optical loss (in dB) scales almost linearly (see Figs. 8(a)–8(c)) when scaling the SP-NNs.

Following the same approach, the crosstalk noise and output power distribution have been captured over the SP-NNs outputs of different dimensions and configurations. The results are depicted in Figs. 8(d)–8(i) using box plots. Observe that for all of the cases, the crosstalk and the output power decrease as we increase the number of inputs and/or the number of layers. The reason for this is the increase in the accumulated insertion loss as we scale up the SP-NNs. Similar to the optical loss, for the crosstalk noise and output power, compared to Reck and Clements, Diamond mesh suffers from significantly larger dynamic range in the output due to its asymmetric design.

Using the crosstalk and output power for SP-NNs of different scales and configurations under the impact of optical loss and crosstalk noise, the signal-to-noise ratio (SNR) at each output can be calculated as:

(13)

Here, and are the optical powers at the output and input , respectively. Also, is the crosstalk power at output and can be calculated based on (6). It should be noted that the output powers in SP-NNs are a function of network input powers, optical losses in each MZI, and phase configurations in each MZI in the network. Therefore, the output optical power values cannot be precisely estimated by having only the optical loss of the network. The results for SNR related to SP-NNs of different scales and mesh configurations are depicted in Figs. 8(j)–8(l). Observe that

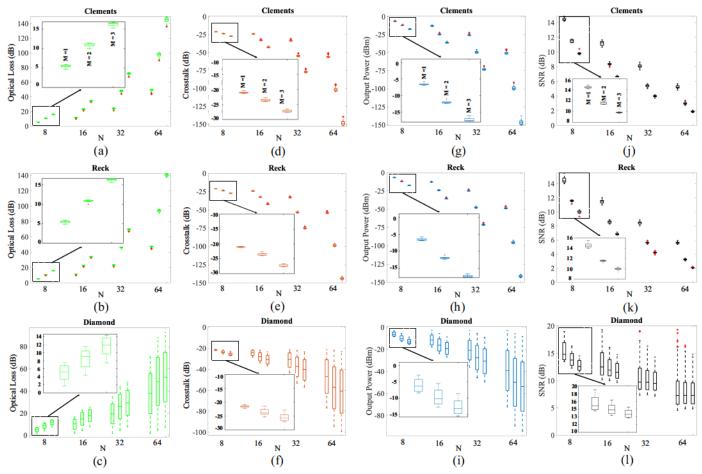


Fig. 8: Box plot presentation of the average optical loss, crosstalk, output power, and SNR distribution over the SP-NNs' outputs for 1000 random weight matrices, based on the network in Fig. 1 and with different numbers of inputs (N), layers (M), and mesh configurations. In (g), (h), and (i), the optical input power at layer one is 0 dBm. The SNR is calculated considering the output powers, which consider optical losses, and the crosstalk noise on the outputs. For better visibility, the insets show the results for N=8.

for all of the cases, SNR can be as high as 15 dB for Reck and Clements and about 19 dB for Diamond, and decreases significantly as we scale up the network which shows the domination of accumulated crosstalk over the main optical signal at the SP-NNs' outputs. Note that also in this case, the SNR distribution on the Diamond mesh's outputs is spread more compared to that in Reck and Clements, because of the asymmetric mesh design in Diamond configuration.

C. Impact of Optical Loss and Crosstalk Noise on SP-NN Inferencing Accuracy With Ideal NAUs

To understand the system-level performance degradation in SP-NNs due to the impact of optical loss and crosstalk noise, we consider three case studies of an SP-NN with two hidden layers (M=3) of 16 neurons each (N=16) with Clements, Reck, and Diamond mesh configurations, trained on the MNIST handwritten digit classification task using the same hyperparameters in all the three case studies. Each image in the MNIST dataset is converted to a complex feature vector of length 16 using the shifted fast Fourier transform (shifted-FFT), discussed in [26]. The nominal test accuracy after training is 91.5%, 92.4, and 90.6 for Reck, Clements,

and Diamond mesh configurations, respectively. Note that the nominal test accuracy is slightly different for the three case studies due to structural differences in the mesh configurations which leads to different phase settings and distribution over the MZIs in the SP-NNs. To analyze the effect of optical loss and crosstalk noise during inferencing, we integrated the proposed MZI models in (3) and (4) into our SP-NN model implementation. Note that in all of the simulations in this section, an ideal *ReLU* nonlinear activation response was considered focusing on the impact of optical loss and crosstalk from OIUs on the SP-NNs inferencing accuracy.

We consider the expected values of α_{DC} , α_m , and α_p within the ranges 0.1–0.4 dB [16], 0.1–0.3 dB [30], and 1–4 dB/cm [16], respectively. Considering an MZI with a length of $l_{MZI}=300~\mu m$ from [17], the propagation loss per MZI $(\alpha_p \cdot l_{MZI})$ is 0.03–0.12 dB. Our analyses showed that the impact of the DC insertion loss (α_{DC}) on the SP-NN's inferencing accuracy is significantly higher compared to that of α_m and α_p for Clements: the accuracy dropped to $\approx 10\%$ when only the impact of α_{DC} was considered in the network (see Fig. 9(a)-(c)). Note that the simulation results in Fig. 9(a) showed the same impact of α_{DC} on inferencing accuracy for

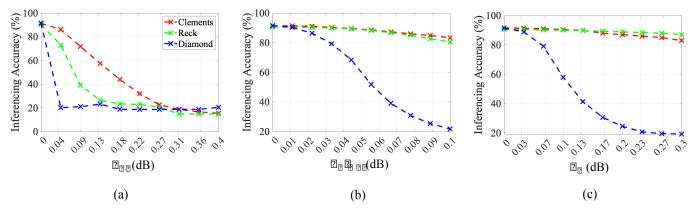


Fig. 9: SP-NN inferencing accuracy in the presence of (a) DC insertion loss (), (b) propagation loss in the MZI (), and (c) absorption loss through phase shifter metal planes ().

SP-NNs with Reck and Diamond mesh configurations when only one source of optical loss was considered at a time. Moreover, the Diamond mesh configuration showed the most susceptibility to the optical insertion loss when one source of optical loss in the MZI was considered at a time (see Fig. 9(a)-(c)).

To understand the simultaneous effect of different optical loss sources (i.e., DCs) on the SP-NN inferencing accuracy, for each case study (i.e., SP-NNs with Clements, Reck, and Diamond mesh configurations), two experiments are defined in which all of the three device-level loss values vary simultaneously from a half-normal distribution. Note that for each experiment, the simulations were repeated 1000 times to avoid loss of generality. The statistical characterization of the two experiments is as the following:

(EXPT1) considers a mean of minimum expected loss value and standard deviation, , such that their maximum expected loss value. Results for EXPT1 are shown in Fig. 10(a)–(c).

(EXPT2) considers loss values with a mean of 0 and the same standard deviation as EXPT1. Results for EXPT2 are shown in Fig. 10(d)–(f).

As can be seen from EXPT1 results in Fig. 10(a)–(c), out of the three mesh configurations, Clements shows the most resilience to optical loss. Considering Fig. 10(a), almost 60% of the 1000 scenarios show inferencing accuracy less than 25%. This is much lower compared to the Reck and Diamond which show about 90% and 100% of the cases with an accuracy below 25%, respectively (see Figs. 10(b) and (c)). Note that in the Clements mesh configuration, only 14% of 1000 scenarios show inferencing accuracy higher than 50% and none of the scenarios in Reck and Diamond mesh configurations show accuracy higher than 50%. Considering the results for EXPT2 shown in Figs. 10(d)–(c), we can obtain the maximum tolerable device-level optical loss (for each source of optical loss) while considering a threshold for maximum acceptable drop in the inferencing accuracy. In this paper, we consider 5% accuracy drop in the presence of the simultaneous effect of all sources of optical loss. Accordingly, the Clements mesh configuration can tolerate up to 0.22 dB of metallic loss, 0.04 dB of propagation loss, and 0.08 dB of DC's loss in

the MZIs, while these values are lower for the other two configurations (see Figs. 10(d)-(c)).

Considering (6), to understand the impact of optical crosstalk noise on SP-NN inferencing accuracy, we model using a linear interpolation between the crosstalk noise the worst-case (MZI in Cross-bar, 18 dB) and the best-case (MZI in Bar-state, 25 dB) crosstalk; see Section III. Fig. 11 shows the inferencing accuracy in the three SP-NN mesh configurations in the presence of both optical loss and crosstalk when and for different and values and with optical losses set to their corresponding minimum expected values. As shown in Fig. 11, when 18 dB and 25 dB (considered as an example based on the work in [31]), the accuracy drops to 11.5% for Clements, 10.1 % for Reck, and 8.3% for Diamond. We also found that under expected values of optical crosstalk and loss, the accuracy remains at 10% in all three mesh configurations. Note that when decreases to below -50 dBm (lower-left corner in Figs. 11(a)-(c)), the accuracy saturates at about 68% for Clements which is significantly higher than Reck at 43.3% and Diamond at 10.1%. In all three case studies presented in this paper, only considering crosstalk noise and neglecting the loss parameters led to an accuracy lower than 25%. This shows the more critical effect of the crosstalk noise than the loss parameters on SP-NNs' accuracy.

The results presented in this section motivate the need for SP-NN design exploration and optimization from the device to the system level to alleviate the impact of optical loss and crosstalk. Moreover, our proposed approach can be used to determine the suitable mesh configuration based on design requirements, and also the maximum tolerable crosstalk and component optical losses to guarantee certain inferencing accuracy.

D. Impact of Optical Loss and Crosstalk Noise on SP-NN Inferencing Accuracy with Non-Ideal NAUs

To understand the impact of optical loss and crosstalk noise on SP-NNs' system-level performance when optoelectronic NAUs are used, we use the analytical model in (11) and parameters reported in Table I to emulate the *ReLU* as the nonlinear

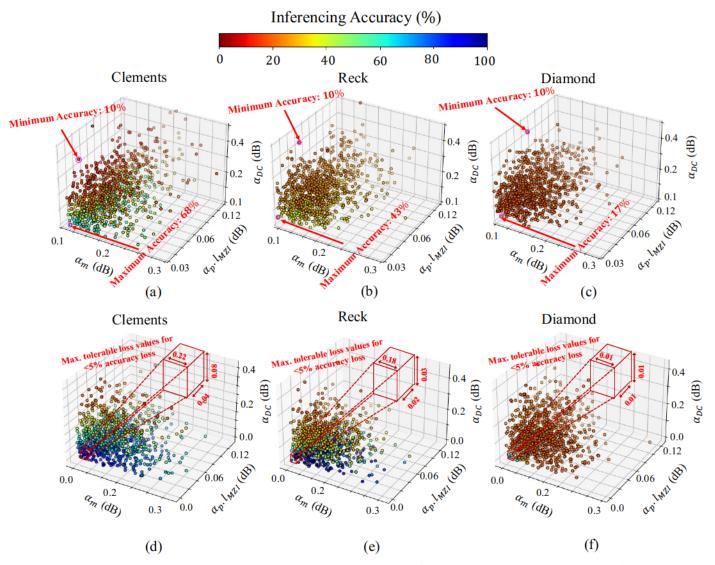


Fig. 10: Inferencing accuracy when the loss parameters $(\alpha_{DC}, \alpha_m, \text{ and } \alpha_p \cdot l_{MZI})$ are simultaneously varied for three different OIU mesh configurations. Two 16×16 hidden layers (M=3) have been considered for each case. (a)–(c) show the case where each of the 1000 points in the scatter plot represents an instance of the SP-NN where the α 's are sampled from a half-normal distribution with mean, μ = their minimum expected value and standard deviation, σ , such that 3σ = their maximum expected value. (d)–(f) show the case where α 's are sampled from a half-normal distribution with mean, μ = 0 and standard deviation, σ , such that 3σ = their maximum expected value to show the maximum tolerable α values for each configuration. Note that the effect of the optical crosstalk noise is neglected in these simulations focusing on the effect of optical loss on the SP-NNs inferencing accuracy.

activation function using the NAU depicted in Fig. 5. The mean-square error (MSE) of the NAU's response (compared to the ideal case) with parameters listed in Table I is used as a metric to understand how optical loss and crosstalk noise from OIUs deteriorates the optoelectronic NAU's performance. Note that for each iteration, 1000 random crosstalk noise phases (θ_{err} , see (11)) are considered and the mean of the resulting MSE is reported in Fig. 12. In this simulation, the PD's sensitivity, dark current, and bandwidth are considered according to Table I in the optoelectronic NAU [30]. Moreover, we also consider the optical loss and crosstalk noise related to the input DC and the MZI in the optoelectronic NAU (see Fig. 5). Considering the crosstalk noise and loss at the output of the OIU unit, we can see that the MSE for a single optoelectronic

NAU used in the SP-NN can be as high as 180 meaning that the actual response of the NAU is significantly deviated from the ideal ReLU nonlinear activation response. To analyze the system-level effect of optical loss and crosstalk noise in SP-NNs implemented using optoelectronic NAUs, three different SP-NN networks with 2 hidden layers (M=3) and 16 inputs with Clements, Reck, and Diamond mesh configuration for the OIU and the optoelectronic activation as the NAU (see (11)) were trained on shifted-FFT MNIST handwritten digit dataset using Neuroptica with nominal test set accuracy of 90%, 93.8, and 94.1% for the Diamond, Reck, and Clements, respectively. We found out that for all three cases when considering the minimum expected values of the device-level loss parameters reported in Table I for MZIs in

and
) have been simulations

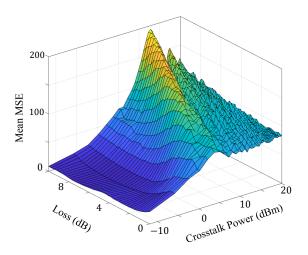


Fig. 12: Average MSE for optoelectronic NAU under the effect of the optical loss and crosstalk noise from OIUs. Each point in the figure has been averaged over 1000 random phases for input crosstalk noise interfering with the victim optical signal.

the OIU and optoelectronic NAU, even a negligible crosstalk noise related to the MZIs in the OIUs (60 dBm) can lead to a significant drop in the inferencing accuracy to below 15%, which is drastically lower than the accuracy reported in Fig. 11 when ideal *ReLU* activation function is used (68% for Clements and 43.3% for the Reck). Therefore, the SP-NNs which are implemented using optoelectronic NAUs are significantly sensitive to the optical loss and crosstalk noise from OIUs. Note that only considering standalone optical loss and crosstalk from optoelectronic NAUs while neglecting the optical loss and crosstalk from OIUs leads to less than 2% drop in the SP-NNs inferencing accuracy.

V. POWER PENALTY AND SCALABILITY CONSTRAINT

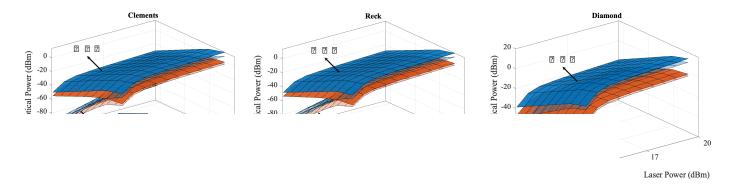
Leveraging the results from the previous section, here we analyze the impact of optical loss and crosstalk noise on SP-NN power consumption (i.e., laser power penalty) as well as scalability constraints in SP-NNs.

A. Power Penalty due to Optical Loss and Crosstalk

The optical loss and coherent crosstalk necessitate an increase in the laser power at the input to compensate for optical loss and crosstalk in SP-NNs. We study this power penalty by considering the input optical laser power () required at the SP-NN input to compensate for the impact of optical loss and coherent crosstalk at the output (see Fig. 1). For the network output in a coherent SP-NN, the input optical laser power should satisfy:

(13)

Here, are the insertion loss (in dB) and and coherent crosstalk power (in dBm), respectively, at the network output . They can be calculated using (5) and (6). Note that the total insertion loss for a signal at output is determined by both and the interference between the victim signal and the coherent crosstalk signal (determined by crosstalk signal phase) at the same output, where the coherent crosstalk power also depends on . Also, is the sensitivity of the photodetector (in dBm) in electronic or optoelectronic NAUs [25], taken to be 11.7 dBm in this paper [41]. The output power and crosstalk power, averaged across all the network outputs, for SP-NNs of different scales and configurations as a function of the input laser power is depicted in Figs. 13(a)-13(c). Observe that as the input laser power increases to compensate for the effect of optical loss and crosstalk, both the output power and the crosstalk power increase at almost similar rates. This indicates that increasing the input laser power will not enhance the SNR in the output due to the increased amount of coherent optical crosstalk power, which increases with the input power as shown in (13), on the outputs of the SP-NNs. Moreover, from Fig. 13 we can see that as we scale up the network, higher input power is required for the output power to be detected by the PDs. In addition, observe the significant decrease in output power due to accumulated optical loss and crosstalk as we scale up SP-NNs, making experimental detection of the output power extremely challenging (see Fig. 13 (a)-(c)).



1 configurations, out laser power.

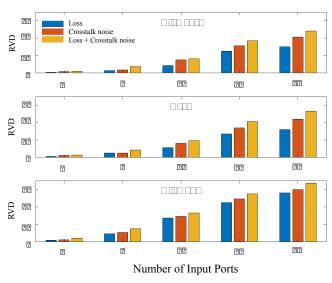


Fig. 14: Average RVD values for 1000 random weight matrices for SP-NNs with different sizes and mesh configurations. Three different scenarios of loss only, crosstalk only, and loss and crosstalk values are considered.

B. Scalability Constraints due to Optical Loss and Crosstalk

As it was shown in the previous sections, optical loss and crosstalk significantly limit the scalability of the MZI-based SP-NNs. The work in [27] showed that the matrix-vector multiplications can be done in multiple steps using a single MZI-based OIU. The trained weight matrix can be broken down into multiple sub-matrices and by loading the sub-matrices into memory and repeatedly updating the phase settings of the MZIs in the network, one large matrix-vector multiplication can be carried out using a smaller mesh of OIU to limit the effect of optical loss and crosstalk on the network's performance. Another example of this application is when MZI-based OIUs are used in a crossbar architecture to carry out optical multiplication, which was presented in [44].

To analyze the scalability constraints due to optical loss and crosstalk noise in SP-NNs for the aforementioned scenarios, we use RVD as a parameter to measure the deviation between

an intended transfer matrix and a deviated transfer matrix (due to loss and crosstalk) in an OIU. As a result, this metric can be used to assess how the ideal transfer matrix deviates when optical loss and crosstalk noise are included in OIUs. RVD can be written as:

In this formulation, and are the transfer matrices of the OIU with different mesh configurations with and without the loss and crosstalk noise, respectively. When the RVD is closer to 0, the deviated transfer matrix is closer to the ideal one, hence the inferencing accuracy is higher, as shown by [8]. Fig. 14 reports the RVD values for Clements, Reck, and Diamond mesh configurations with different sizes under optical loss and crosstalk reported in Table I. For these simulations, 1000 random weight matrices were tested on a single MZI-based OIU with Clements, Reck, and Diamond configuration with 4, 8, 16, 32, and 64. The reason for using 1000 random weight matrices is to avoid loss of generality in our analysis. As shown by Fig. 14, in all cases the impact of optical crosstalk is more critical for scalability than the optical loss (see the results for crosstalk noise in the figure). One reason for this is due to using SOAs to compensate for the total insertion loss of the network that also amplifies the coherent optical crosstalk noise. Moreover, the RVD increases as the network scales up. Observe that the RVD increase is even worse for the Diamond mesh due to using a larger number of MZIs when scaling the network.

To better understand the relationship between RVD and inferencing accuracy of SP-NNs, a single OIU of different mesh configurations with 4, 8, 16, 32, and 64 followed by an ideal *ReLU* NAU were trained on a linearly separable Gaussian dataset presented in [13]. The test set accuracy is 100% for all the case studies. Note that the proposed Gaussian dataset is considered because it is simpler than MNIST dataset and requires a single layer only. As a result, we use it as an example to show the relationship between RVD, network accuracy, and the scalability constraints in OIUs across different OIU mesh configurations. The RVD values against the inferencing accuracy of a single OIU of

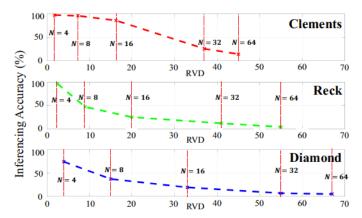


Fig. 15: RVD values for weight matrices for SP-NNs with different sizes and mesh configurations trained on a linearly separable Gaussian dataset against the inferencing accuracy

different sizes and configurations under the impact of both optical loss and crosstalk noise (using the parameters listed in Table I) trained on a linearly separable Gaussian dataset is reported in Fig. 15. Observe that as the RVD increases, the network's accuracy decreases. Moreover, out of three mesh configurations, Clements showed less than 4% and 15% drop in the inferencing accuracy when N=8 and N=16, respectively, which is significantly lower than Reck with 53% and 75% and than Diamond with 63% and 81% accuracy drops, respectively. The Reck configuration shows about a 4% drop in accuracy when N=4. Furthermore, the Diamond configuration shows a catastrophic 30% drop in accuracy even when a deficient number of inputs N=4 is used.

Although the Diamond structure shows the least tolerance to optical loss and crosstalk, the work in [13] suggests that this mesh configuration has the most tolerance to fabrication-process variations compared to Clements [8]. However, our results presented in this paper show that Diamond mesh cannot be scaled up in the presence of optical loss and crosstalk noise. Note that among the three mesh configurations studied in this paper, Clements architecture showed the highest resilience to optical loss and coherent crosstalk, making it the preferred configuration to implement SP-NNs.

System architects can benefit from the proposed analyses in this paper to better understand the impact of optical loss and crosstalk noise in SP-NNs, and how such an impact change among different SP-NN architecture choices. Also, our analyses can help device designers to better understand device-level performance requirements (e.g., maximum optical loss at the device level) to achieve certain performance and accuracy in SP-NNs. Loss- and crosstalk-aware training of the MZI-based SP-NNs can be considered as a possible solution to alleviate the effect of optical loss and crosstalk in MZI-based photonic computing systems, not studied in this paper.

VI. CONCLUSION

The performance and scalability of SP-NNs are limited by optical loss and crosstalk noise in silicon photonic devices. This paper presents a framework for modeling optical loss and crosstalk noise for SP-NNs of different scales with different

mesh configurations. We presented a detailed and comprehensive analysis of optical loss and crosstalk noise and the corresponding SNR in SP-NNs with three different mesh configurations of Clements, Reck, and Diamond while exploring the drops in inferencing accuracy under different scenarios. In particular, the results showed a significant accuracy loss of at least 84% and dramatic SNR degradation for all case studies due to optical loss and crosstalk noise as we scale up the SP-NNs. Additionally, we conducted an extensive analysis of optical loss and crosstalk in optoelectronic NAUs. Moreover, we thoroughly analyzed the scalability limitations of SP-NNs arising from optical loss and crosstalk. The valuable insights presented in this study can be leveraged by SP-NN device and system architects to explore and optimize different challenges in the development and evaluation of SP-NNs in the presence of inevitable optical loss and crosstalk noise.

REFERENCES

- B. J. Shastri, C. Huang, A. N. Tait, T. F. de Lima, and P. R. Prucnal, "Silicon photonic neural network applications and prospects," in AI and Optical Data Sciences III, vol. 12019, 2022.
- [2] Q. Cheng, J. Kwon, M. Glick, M. Bahadori, L. P. Carloni, and K. Bergman, "Silicon photonics codesign for deep learning," *Proc. IEEE*, vol. 108, no. 8, pp. 1261–1282, 2020.
- [3] S. Pasricha and M. Nikdast, "A survey of silicon photonics for energy-efficient manycore computing," *IEEE Des Test.*, vol. 37, no. 4, pp. 60–81, 2020.
- [4] A. Palmieri, A. Shafiee, M. Alasio, A. Tibaldi, G. Ghione, F. Bertazzi, M. Goano, and M. Vallone, "Enhanced dynamic properties of ge-on-si mode-evolution waveguide photodetectors," in 2020 International Conference on Numerical Simulation of Optoelectronic Devices (NUSOD). IEEE, 2020, pp. 1–2.
- [5] S. Kumari, Y. K. Verma, and S. M. Tripathi, "Plasmonic ring resonator sensor with high sensitivity and enhanced figure of merit using an ag-si-ag bus waveguide," *IEEE Transactions on Nanotechnology*, vol. 22, pp. 200-205, 2023.
- [6] M. A. Hamdani and G. Qazi, "Highly efficient and compact silicon based novel michelson interferometer modulator," in 2021 Devices for Integrated Circuit (DevIC). IEEE, 2021, pp. 581–585.
- [7] Y. K. Verma, S. Kumari, G. Bawa, and S. M. Tripathi, "Temperature insensitive large free spectral range micro-ring resonator," *Optical and Quantum Electronics*, vol. 54, no. 12, p. 839, 2022.
- [8] A. Mirza, A. Shafiee, S. Banerjee, K. Chakrabarty, S. Pasricha, and M. Nikdast, "Characterization and optimization of integrated siliconphotonic neural networks under fabrication-process variations," *IEEE Transactions on Nanotechnology*, vol. 21, pp. 763–771, 2022.
- [9] H. Zhang, M. Gu, X. Jiang, J. Thompson, H. Cai, S. Paesani, R. Santagati, A. Laing, Y. Zhang, M. Yung et al., "An optical neural chip for implementing complex-valued neural network," *Nature communications*, vol. 12, no. 1, p. 457, 2021.
- [10] P. R. Wiecha, A. Arbouet, C. Girard, and O. L. Muskens, "Deep learning in nano-photonics: inverse design and beyond," *Photonics Research*, vol. 9, no. 5, pp. B182–B200, 2021.
- [11] F. P. Sunny, E. Taheri, M. Nikdast, and S. Pasricha, "A survey on silicon photonics for deep learning," ACM JETC, vol. 17, no. 4, pp. 1–57, 2021.
- [12] A. Shafiee, S. Banerjee, K. Chakrabarty, S. Pasricha, and M. Nikdast, "LoCI: An analysis of the impact of optical loss and crosstalk noise in integrated silicon-photonic neural networks," in *Proceedings of the Great Lakes Symposium on VLSI 2022*, 2022, pp. 351–355.
- [13] F. Shokraneh, S. Geoffroy-Gagnon, and O. Liboiron-Ladouceur, "The diamond mesh, a phase-error-and loss-tolerant field-programmable mzibased optical processor for optical neural networks," *Optics Express*, vol. 28, no. 16, pp. 23495–23508, 2020.
- [14] W. R. Clements, P. C. Humphreys, B. J. Metcalf, W. S. Kolthammer, and I. A. Walmsley, "Optimal design for universal multiport interferometers," *Optica*, vol. 3, no. 12, pp. 1460–1465, 2016.
- [15] M. Reck, A. Zeilinger, H. J. Bernstein, and P. Bertani, "Experimental realization of any discrete unitary operator," *Physical Review Letters*, vol. 73, no. 1, p. 58, 1994.

- [16] M. Bahadori, S. Rumley, D. Nikolova, and K. Bergman, "Comprehensive design space exploration of silicon photonic interconnects," *IEEE JLT*, vol. 34, no. 12, pp. 2975–2987, 2016.
- [17] F. Shokraneh, M. S. Nezami, and O. Liboiron-Ladouceur, "Theoretical and experimental analysis of a 4 4 reconfigurable MZI-based linear optical processor," *IEEE JLT*, vol. 38, no. 6, pp. 1258–1267, 2020.
- [18] M. A. Hamdani and G. Qazi, "Modelling and theoretical analysis of a novel common mirror based silicon photonic michelson modulator," *Optical and Quantum Electronics*, vol. 55, no. 1, pp. 1–23, 2023.
- [19] A. Shafiee, S. Pasricha, and M. Nikdast, Silicon Photonics for Future Computing Systems. John Wiley & Sons, Ltd, 2022, pp. 1–26. [Online]. Available: https://onlinelibrary.wiley.com/doi/abs/10.1002/047134608X. W8435
- [20] M. Bahadori, S. Rumley, R. Polster, and K. Bergman, "Loss and crosstalk of scalable mzi-based switch topologies in silicon photonic platform," in 2016 IEEE Photonics Conference (IPC). IEEE, 2016, pp. 615–616.
- [21] I. A. Williamson, T. W. Hughes, M. Minkov, B. Bartlett, S. Pai, and S. Fan, "Reprogrammable electro-optic nonlinear activation functions for optical neural networks," *IEEE JSTQE*, vol. 26, no. 1, pp. 1–12, 2019.
- [22] H. Shang, D. Sun, P. Yu, B. Wang, T. Yu, T. Li, and H. Jiang, "Investigation for sidewall roughness caused optical scattering loss of silicon-on-insulator waveguides with confocal laser scanning microscopy," *Coatings*, vol. 10, no. 3, p. 236, 2020.
- [23] M. Nikdast, J. Xu, L. H. K. Duong, X. Wu, X. Wang, Z. Wang, Z. Wang, P. Yang, Y. Ye, and Q. Hao, "Crosstalk noise in wdm-based optical networks-on-chip: A formal study and comparison," *IEEE TVLSI*, vol. 23, no. 11, pp. 2552–2565, 2015.
- [24] M. Kouissi, B. Charbonnier, and C. Algani, "Comprehensive model for evaluating the performance of mach-zehnder-based silicon photonic switch fabrics in large scale," *Appl. Sci.*, vol. 10, no. 23, 2020.
- [25] M. M. P. Fard, I. A. Williamson, M. Edwards, K. Liu, S. Pai, B. Bartlett, M. Minkov, T. W. Hughes, S. Fan, and T.-A. Nguyen, "Experimental realization of arbitrary activation functions for optical neural networks," *Opt. Express*, vol. 28, no. 8, pp. 12138–12148, 2020.
- [26] S. Banerjee, M. Nikdast, and K. Chakrabarty, "Modeling siliconphotonic neural networks under uncertainties," in *IEEE/ACM DATE*, 2021, pp. 98–101.
- [27] Y. Liu, J. Zhang, J. Feng, S. Chen, and J. Xu, "A reliability concern on photonic neural networks," in 2022 Design, Automation & Test in Europe Conference & Exhibition (DATE). IEEE, 2022, pp. 1059–1064.
- [28] B. Haq, S. Kumari, K. Van Gasse, J. Zhang, A. Gocalinska, E. Pelucchi, B. Corbett, and G. Roelkens, "Micro-Transfer-Printed III-V-on-Silicon C-Band Semiconductor Optical Amplifiers," *Laser Photonics Rev*, vol. 14, no. 7, p. 1900364, 2020.
- [29] L. De Marinis, M. Cococcioni, O. Liboiron-Ladouceur, G. Contestabile, P. Castoldi, and N. Andriolli, "Photonic integrated reconfigurable linear processors as neural network accelerators," *Appl. Sci.*, vol. 11, no. 13, 2021.
- [30] F. Ding, J. Dai, Y. Chen, J. Zhu, Y. Jin, and S. I. Bozhevolnyi, "Broadband near-infrared metamaterial absorbers utilizing highly lossy metals," *Scientific Reports*, vol. 6, no. 1, pp. 1–9, 2016.
- [31] Y. Shoji, K. Kintaka, S. Suda, H. Kawashima, T. Hasama, and H. Ishikawa, "Low-crosstalk 2 2 thermo-optic switch with silicon wire waveguides," *Opt. Express*, vol. 18, no. 9, pp. 9071–9075, 2010.
- [32] S. Okamoto, K. Minoguchi, F. Hamaoka, K. Horikoshi, A. Matsushita, M. Nakamura, E. Yamazaki, and Y. Kisaka, "A study on the effect of ultra-wide band wdm on optical transmission systems," *Journal of Lightwave Technology*, vol. 38, no. 5, pp. 1061–1070, 2020.
- [33] Ansys Lumerical. [Online]. Available: https://www.lumerical.com/ products/
- [34] D. M. Baney, P. Gallion, and R. S. Tucker, "Theory and measurement techniques for the noise figure of optical amplifiers," *Optical fiber* technology, vol. 6, no. 2, pp. 122–154, 2000.
- [35] S. Yu, A. Gallet, N. El Dahdah, H. Elfaiki, I. Demirtzioglou, L. Godard, and R. Brenot, "Flat noise figure semiconductor optical amplifiers," in 2021 European Conference on Optical Communication (ECOC), 2021, pp. 1–3.
- [36] P. W. Juodawlkis, J. J. Plant, W. Loh, L. J. Missaggia, K. E. Jensen, and F. J. O'Donnell, "Packaged 1.5- m quantum-well soa with 0.8-w output power and 5.5-db noise figure," *IEEE Photonics Technology Letters*, vol. 21, no. 17, pp. 1208–1210, 2009.
- [37] Y. Shen, N. C. Harris, S. Skirlo, M. Prabhu, T. Baehr-Jones, M. Hochberg, X. Sun, S. Zhao, H. Larochelle, D. Englund *et al.*, "Deep learning with coherent nanophotonic circuits," *Nature Photonics*, vol. 11, no. 7, pp. 441–446, 2017.

- [38] O. Destras, S. Le Beux, F. G. de Magalhães, and G. Nicolescu, "Survey on activation functions for optical neural networks," ACM Computing Surveys, 2023.
- [39] G. Ghione, Semiconductor devices for high-speed optoelectronics, vol. 116.
- [40] M. M. P. Fard, G. Cowan, and O. Liboiron-Ladouceur, "Responsivity optimization of a high-speed germanium-on-silicon photodetector," *Optics Express*, vol. 24, no. 24, pp. 27738–27752, 2016.
- [41] X. Bi, J. Li, Z. Gu, B. Tang, C. Sheng, Y. Yang, and Q. Xu, "High sensitivity and dynamic-range 25 Gbaud silicon receiver chipset with current-controlled DC adjustment path and cube-shape Ge-on-Si PD," *IEEE TCSI: Regular Papers*, vol. 67, no. 11, pp. 3991–4001, 2020.
- [42] S. Geoffroy-Gagnon, "Neuroptica," https://gitlab.com/simongg/ neuroptica, 2020.
- [43] S. Banerjee, M. Nikdast, and K. Chakrabarty, "Characterizing coherent integrated photonic neural networks under imperfections," *Journal of Lightwave Technology*, pp. 1–16, 2022.
- [44] G. Giamougiannis, A. Tsakyridis, Y. Ma, A. Totovic, M. MoralisPegios, D. Lazovsky, and N. Pleros, "A coherent photonic crossbar for scalable universal linear optics," *Journal of Lightwave Technology*, 2023.