

Experimentally realized *in situ* backpropagation for deep learning in nanophotonic neural networks

Sunil Pai,^{1,*} Zhanghao Sun,¹ Tyler W. Hughes,^{1,2} Taewon Park,¹ Ben Bartlett,³ Ian A. D. Williamson,^{1,4} Momchil Minkov,^{1,2} Maziyar Milanizadeh,⁵ Nathnael Abebe,¹ Francesco Morichetti,⁵ Andrea Melloni,⁵ Shanhui Fan,¹ Olav Solgaard,¹ and David A.B. Miller¹

¹*Department of Electrical Engineering, Stanford University, Stanford, CA 94305, USA*

²*now at Flexcompute Inc., Belmont, MA, USA*

³*Department of Applied Physics, Stanford University, Stanford, CA 94305, USA*

⁴*now at X Development LLC, Mountain View, CA USA.*

⁵*Politecnico di Milano, Milan, Italy*

Neural networks are widely deployed models across many scientific disciplines and commercial endeavors ranging from edge computing and sensing to large-scale signal processing in data centers. The most efficient and well-entrenched method to train such networks is backpropagation, or reverse-mode automatic differentiation. To counter an exponentially increasing energy budget in the artificial intelligence computing sector, there has been recent interest in analog implementations of neural networks, specifically nanophotonic optical neural networks for which no analog backpropagation demonstration exists. We design mass-manufacturable silicon photonic neural networks that alternately cascade our custom designed “photonic mesh” accelerator with digital nonlinearities to output the result of arbitrary matrix multiplication of the input signal. These photonic meshes are parametrized by reconfigurable physical voltages that tune the interference of optically encoded input data propagating through integrated Mach-Zehnder interferometer networks. Here, using our packaged photonic chip, we demonstrate *in situ* backpropagation for the first time to solve classification tasks and evaluate a new protocol to keep the entire gradient measurement and update of physical device voltages in the analog domain, improving on past theoretical proposals. This *in situ* method is made possible by introducing three changes to typical photonic meshes: (1) measurements at optical “grating tap” monitors, (2) bidirectional optical signal propagation automated by fiber switch, and (3) universal generation and readout of optical amplitude and phase. After training, our classification achieves accuracies similar to digital equivalents even in the presence of systematic error. Our findings suggest a new training paradigm for photonics-accelerated artificial intelligence based entirely on a physical analog of the popular backpropagation technique.

Neural networks (NNs) are intelligent computational graph-based models that are ubiquitous in scientific data analysis and commercial artificial intelligence (AI) applications such as self-driving cars and speech recognition software. Through deep learning, NN models are dynamically “trained” on input image, audio or language data to automatically make decisions (“inference”) for complex signal processing powering much of today’s modern technology. Due to increasing demand, these models require an ever-increasing computational energy budget, which has recently been estimated to double every 3 to 4 months, according to OpenAI [1]. An increasingly large reservoir of available data and adoption of AI in modern technology necessitates an energy-efficient solution for training of NNs.

In this paper, we experimentally demonstrate the first (to our knowledge) optical implementation of *backpropagation*, the most widely used and accepted method of training NNs [2, 3], on a scalable foundry-manufactured device. (A minimal bulk optical demonstration has been previously explored [4].) Specifically, backpropagation consists of backward propagating model errors computed for training data through the NN graph to determine

updating gradients on each element in the graph. Importantly, this can be physically implemented in linear optical devices by simply sending light-encoded errors backwards through photonic devices and performing optical measurements [5], which is a faster and more efficient calculation than a digital implementation. Our demonstration of this new physics-based backpropagation algorithm, and analysis of systematic error in gradient calculations used in backpropagation, could help ultimately offer new, possibly energy-efficient strategies to teach modern AI to make intelligent decisions on mass manufacturable silicon photonics hardware using efficient model-free training [5, 6].

As a physical platform for our backpropagation demonstration, we explore programmable nanophotonic devices called “photonic meshes” for accelerating matrix multiplication [7, 8]. Photonic meshes shown in Fig. 1 are silicon-based low-cost, commercially scalable $N \times N$ port photonic integrated circuits (PICs) consisting of Mach-Zehnder interferometers (MZIs) and programmable optical phase shifts. These PICs are capable of representing matrix-vector multiplication (MVM) through only the propagation of guided monochromatic light (1560 nm in our demonstration) through N silicon waveguide “wires” clad by silicon oxide [7–9]. Each waveguide can support a single optical mode which has two degrees of freedom:

* sunilpai@stanford.edu

amplitude and phase, yielding a complex N -dimensional vector \mathbf{x} at the input of the system. Programmable phase shift settings physically modulate the propagation speed (and relative phases) of the wave over segments of silicon wire to affect how the N propagating modes constructively or destructively interfere in each interferometer. The energy efficiency of these devices has been estimated to be up to two orders of magnitude higher than current state-of-the-art electronic application-specific integrated circuits (ASICs) in AI [10].

Assuming no light is lost in the ideal photonic circuit, the mesh can be programmed to transform optical inputs using an arbitrary programmable unitary MVM $\mathbf{y} = U\mathbf{x}$ [8, 11, 12]. The matrix U is parametrized by the programmable phase shifters on the device and transforms inputs \mathbf{x} propagating through the device to output modes \mathbf{y} . The programmed phase shifts on the device define the matrix U , and the N output mode amplitude and phase measurements \mathbf{y} represent the solution to this optical computation. This fundamental mathematical operation enables meshes to be widely employed in various analog signal processing applications such as telecommunications [9], quantum computing [14], sensing, and machine learning [7], the last of which we explore experimentally in this work via our backpropagation demonstration.

To form what we call a “hybrid” digital-photonic NN (PNN), we alternately cascade photonic meshes and digital nonlinear functions [13, 15], which ultimately forms a composite function and model capable of complex decision making. While performing inference or backpropagation, the hybrid PNN performs time- and energy-efficient MVM, converts photonic mesh output signals to the digital domain, applies nonlinearities, and then converts the data back to optical domain for MVM in the next layer. Hybrid PNNs offer more versatility over fully analog PNNs in the near term due to flexible manipulation of signals in the digital domain easing implementations for recurrent and convolutional neural networks. As a result, hybrid PNNs have been demonstrated to provide a reliable low-latency and energy-efficient analog optical solution for inference, recently in circuit sizes of up to 64×64 in commercial settings [16]. Despite this success in PNN-based inference, on-device backpropagation training of PNNs has not been demonstrated, due to significantly higher experimental complexity compared to the inference procedure.

In this paper, we address this gap by experimentally demonstrating *in situ* backpropagation in a hybrid PNN architecture. First, we propose a novel and energy-efficient implementation of the technique for measuring phase shifter updates entirely in the optoelectronic (analog) domain. Second, we experimentally validate training a backpropagation-enabled, foundry-manufactured photonic circuit using a custom optical rig setup on a multilayer neural network. Our demonstration solves machine learning tasks on this photonic hardware using this optically-accelerated backpropagation with similar accuracy compared to a conventional digital implementation,

adding new capabilities beyond existing inference or *in silico* learning demonstrations [7, 17, 18]. Our findings ultimately pave the way for a new class of approaches for energy-efficient analog training of neural networks and optical devices more broadly.

PHOTONIC NEURAL NETWORKS

At a high level, a neural network is able to transform data into useful decisions or interpretations, an example of which is shown in Fig. 1(a) and (d) where we label points in 2D based on their location within or outside of a ring. This problem (and in principle many more complex problems like audio signal processing [7]) can be solved in the optical domain using photonic neural networks (PNNs) as shown in Fig 1(a)-(d). To solve the problem, we design and evaluate a hybrid digital-optical deep PNN architecture parameterized by trainable programmable phase shifts $\boldsymbol{\eta} \in [0, 2\pi)^D$, where D represents the total number of phase shifting elements across all layers in the overall PNN.

Using a combination of photonic hardware and software implementations, our hybrid PNN can solve non-trivial tasks using alternating sequences of analog linear optical MVM operations $U^{(\ell)}(\boldsymbol{\eta}^{(\ell)})$ and digital nonlinear transformations $\mathbf{f}^{(\ell)}$ where ℓ denotes the neural network layer and we assume a total of L layers. For example, after $L = 3$ neural network layers for $N = 4$, our photonic neural network is capable of transforming the boundary function used to separate the labelled points in Fig. 1(d). To realize this implementation in a mathematical model, the following sequence of functions transforms the data, proceeding in a “feedforward” manner through the layers of the network:

$$\begin{aligned} \mathbf{y}^{(\ell)} &= U^{(\ell)}\mathbf{x}^{(\ell)} \\ \mathbf{x}^{(\ell+1)} &= \mathbf{f}^{(\ell)}(\mathbf{y}^{(\ell)}). \end{aligned} \quad (1)$$

The inputs $\mathbf{x} = \mathbf{x}^{(1)}$ to the overall system are forward-propagated to the final layer (layer L), outputting $\hat{\mathbf{z}} := \mathbf{x}^{(L+1)}$. This forward propagation and resulting output measurement of data sent through this network is called “inference” and is depicted in Fig. 1(a, b, d). The model cost or error function is represented by $\mathcal{L}(\mathbf{x}, \mathbf{z}) = c(\hat{\mathbf{z}}(\mathbf{x}), \mathbf{z})$ for a given set of ground truth labels \mathbf{z} , where c is any cost function representing the error between $\hat{\mathbf{z}}$ and \mathbf{z} . We refer to the input, label pair of training data (\mathbf{x}, \mathbf{z}) as a “training example.” Backpropagation and other gradient-based training approaches seek to update parameters $\boldsymbol{\eta}$ based on the vector gradient $\frac{\partial \mathcal{L}}{\partial \boldsymbol{\eta}} \in \mathbb{R}^D$ evaluated for a given training example (or averaged over a batch of training examples). We now explain how we implement both the inference and backpropagation training calculations directly on our core photonic neural network.

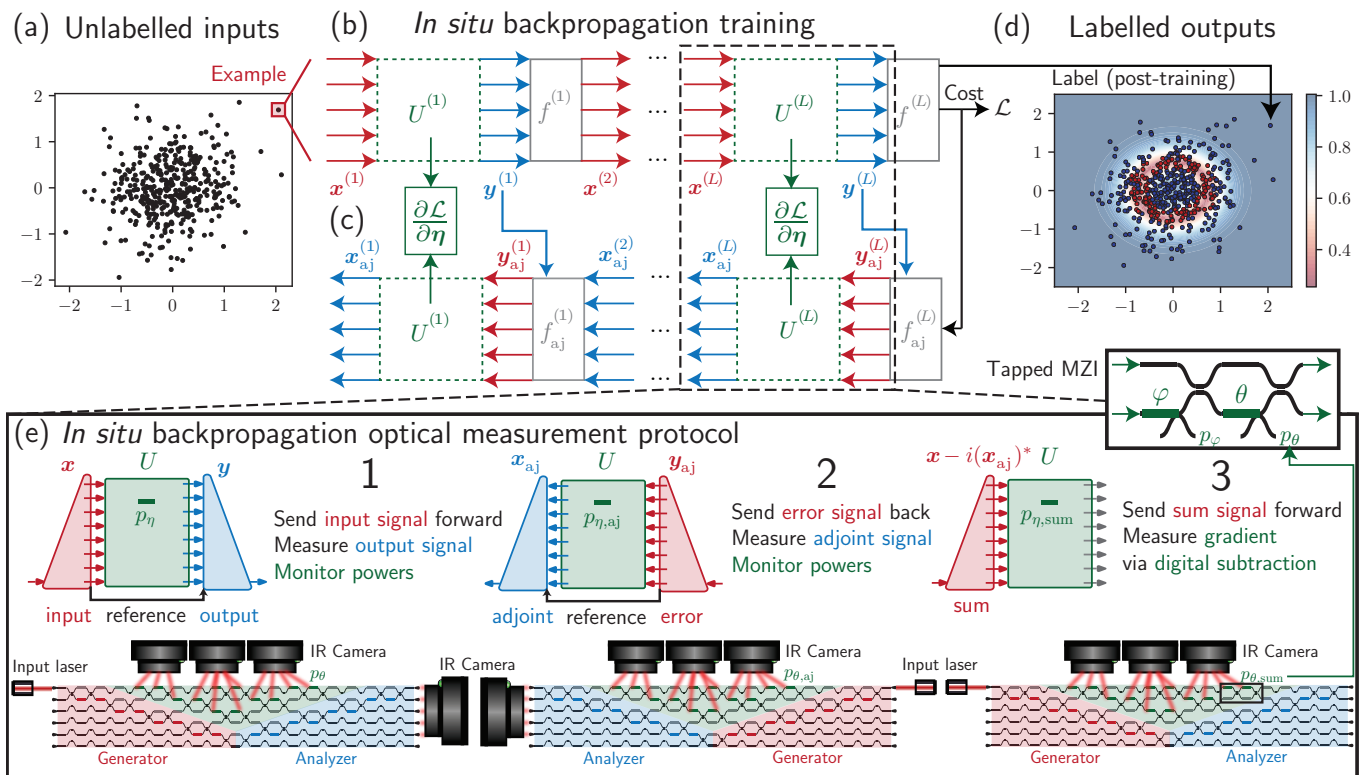


FIG. 1. (a) Example machine learning problem: an unlabelled 2D set of inputs that are formatted to be input into a photonic neural network. (b, c) *In situ* backpropagation training of an L photonic neural network for (b) the forward direction and (c) the backward direction showing the procedure for calculating gradient updates for phase shifts. (d) An inference task implemented on the actual chip results in good agreement between the chip-labelled points and the ideal implemented ring classification boundary (resulting from the ideal model) and a 90% classification accuracy. (e) We depict our proposed architecture and the three steps of *in situ* (analog) backpropagation, consisting of a 6×6 mesh implementing coherent 4×4 forward and inverse unitary matrix-vector products using a reference arm. We depict the (1) forward (2) backward (3) sum steps of *in situ* photonic backpropagation. Arbitrary input setting and complete amplitude and phase output measurement are enabled in both directions using the reciprocity and symmetries of our architecture. All powers throughout the mesh are monitored using the tapped MZI shown in the inset for each step, allowing for digital subtraction to compute the gradient [5]. These power measurements performed at phase shifts are indicated by green horizontal bars.

BACKPROPAGATION DEMONSTRATION

For practical demonstration purposes, our multilayer PNN is completely controlled by a single photonic mesh (Note that in practice, energy-efficient photonic NNs are controlled by separate photonic meshes of MZIs for each linear layer). Each MZI unit is controlled by an electronic control unit that applies voltages to set various phase shifts on the device packaged on a thermally controlled assembly as shown in Fig. 2(a, b). These phase shifts are placed at the input external arm of the MZI (ϕ , controlled by voltage v_ϕ) and in the internal arm of the MZI (θ controlled by voltage v_θ); this ultimately controls the propagation pattern of the light through the chip, enabling arbitrary unitary matrix multiplication. In our chip specifically, we embed an arbitrary 4×4 unitary matrix multiply in a 6×6 triangular network of MZIs. This configuration incorporates two 1×5 photonic meshes on either end of the 4×4 “Matrix unit”

(shown in green) capable of sending any input vector \mathbf{x} and measuring any output vector \mathbf{y} from Eq. 1. These calibrated optical I/O circuits are referred to as “Generator” and “Analyzer” circuits are shown in red and blue respectively in Figs. 1(e) and 2(b). A more complete discussion of how an MVM operation is achieved using our architecture is provided in the Methods, and similar approaches have been attempted for complex-valued photonic neural network architectures [19]. Note that for the input generation and output measurements, we need to calibrate the voltage mappings $\theta(v_\theta)$, $\phi(v_\phi)$ (equivalently for the output measurement, $v_\theta(\theta)$, $v_\phi(\phi)$), which is discussed in detail in Ref. 20. This is a standard calibration protocol [7, 21, 22] discussed at length in the Appendix and required for accurate operation of the chip.

Our core contribution in this paper, shown in Fig. 1(e), is to devise and test a photonic mesh matrix accelerator architecture that experimentally implements backpropagation as proposed in Ref. 5 within a hybrid digital-

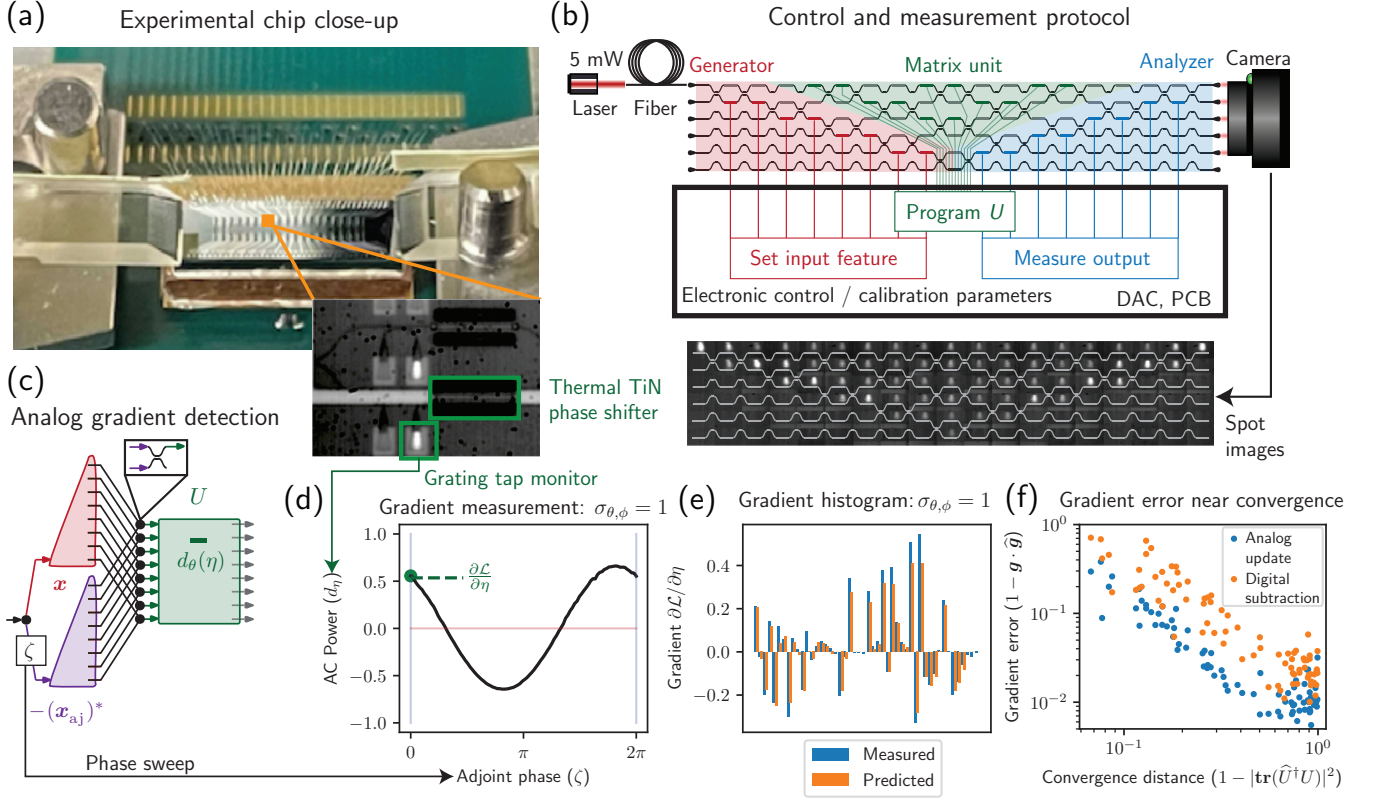


FIG. 2. (a) Image of the packaged photonic chip wirebonded to a custom PCB with fiber array for laser input and a camera overhead for imaging the chip. Zooming in reveals the core control-and-measurement unit of the chip, enabling power measurement using 3% grating tap monitors and a thermal phase shifter nearby. (b) The DAC control is used to set up inputs and perform coherent detection of the outputs, and the IR camera over the chip can be used to image all waveguide segments throughout the chip via grating tap monitors, which we use for backpropagation. (c) Analog gradient detection could be used to measure the gradient by introducing a summing interference circuit (not implemented on the chip in (b)) between the input and adjoint fields. (d) The adjoint phase ζ can be swept from 0 to 2π to perform an all-analog gradient measurement, and, ultimately in principle, update phase shifters with an optoelectronic scheme. (e) Gradients measured using our analog scheme yield approximately correct gradients when the implemented mesh is perturbed from the optimal (target) unitary $U = \text{DFT}(4)$ with phase error $\sigma_{\theta, \phi} = 1$. (f) The average normalized gradient error (averaged over 20 instances of random implemented \hat{U}) decreases with distance of the device implementation $\hat{U}(\eta)$ from the optimal $U = \text{DFT}(4)$. This “distance” is represented in terms of the fidelity error $1 - |\text{tr}(\hat{U}^\dagger U)|^2$.

analog model implementing the most expensive operations using universal linear optics. Our backpropagation-enabled architecture differs from previously-proposed photonic mesh architectures in three ways:

1. We enable “bidirectional light propagation,” the ability to send and measure light propagating left-to-right or right-to-left through the circuit (as depicted in Fig. 1(e)).
2. We implement “global monitoring,” the ability to measure optical power at any waveguide segment in the circuit using 3% grating taps (shown in the inset of Fig. 1(e) and Fig. 2(a, b)). In our proof-of-concept setup, we use an IR camera mounted on an automated stage to image these taps throughout the chip.

3. We implement both amplitude and phase detection (improving on past approaches [19]) using a self-configuring programmable Matrix unit layer [20, 23] on both the red and blue Generator and Analyzer subcircuits of Fig. 1(e) and Fig. 2(b), which by symmetry works for sending and measuring light that propagates forward or backward through the mesh.

These improvements on an already versatile photonic hardware architecture enable backpropagation-based machine learning implemented entirely using optical measurement to optimize programmable phase shifters in PNNs. As shown in Fig. 1(e), all three steps of backpropagation 5 require monitoring of optical powers at each phase shifter and measurement of complex field outputs at the left and right sides of the mesh. Furthermore,

the bidirectionality of the monitoring and optical I/O is required to switch between forward and backward propagation of signals required for *in situ* backpropagation to be experimentally realized. Equipped with these additional elements, our protocol can be implemented on any feedforward photonic circuit [6] with the requisite Analyzer and Generator circuitry, though we use a triangular mesh in this work to enable the chip to be used in other applications [24].

Here we give a quick summary of the procedure (fully described in the Appendix). For each layer ℓ and training example pair, a “forward inference” signal $\mathbf{x}^{(\ell)}$ is sent forward and a corresponding “backward adjoint” signal $\mathbf{x}_{\text{aj}}^{(\ell)}$ is sent backward through a mesh implementing $U^{(\ell)}$. The backward pass is in a sense a mirror image of the forward pass (error signal is sent from final layer to input layer) which algorithmically computes an efficient “reverse mode” chain rule calculation. The final step sends what we call a “sum” vector $\mathbf{x}^{(\ell)} - i(\mathbf{x}_{\text{aj}}^{(\ell)})^*$. Previously [5], it was shown that global monitoring in all three steps enables us to calculate the gradient by subtracting the backward and forward measurements from sum measurements in the digital domain, in what we call an “optical vector-Jacobian product (VJP)” (Appendix).

ANALOG UPDATE

Going beyond an experimental implementation of the theoretical proposal of Ref. 5, we additionally explore a more energy-efficient fully analog gradient measurement update for the final step that avoids the digital subtraction update. The key difference is in the final “sum” step where we instead sweep the adjoint phase ζ (giving $\mathbf{x}^{(\ell)} - i(\mathbf{x}_{\text{aj}}^{(\ell)})^* e^{i\zeta}$) from 0 to 2π repeatedly (e.g., using a sawtooth signal). During the sweep, we record $d_\eta(\zeta)$, the AC component of the measured power monitored through phase shifter θ , $p_{\eta, \text{sum}}(\zeta)$. It is straightforward to show that gradient is $d_\eta(0)$, the AC component evaluated when no adjoint phase is applied (Appendix). To achieve the ζ sweep physically, we can employ the summing architecture in Fig. 2(c) which sums $\mathbf{x}^{(\ell)}, i(\mathbf{x}_{\text{aj}}^{(\ell)})^*$ interferometrically with a constant loss factor of 1/2 in power ($1/\sqrt{2}$ in amplitude). Then, using a boxcar gated integrator and high pass filter, we can physically compute $d_\eta(\zeta)$ and update the phase shift voltage entirely in the analog domain (Appendix). Ultimately, this approach potentially avoids a costly analog-digital conversion and additional memory complexity required to program N^2 elements. Since the PNN has L layers, the gradient calculation step requires local feedback circuits at each phase shifter η that update the parameters using the measured gradient:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \eta} &= -\mathcal{I}(x_\eta x_{\eta, \text{aj}}) \\ &= (|x_\eta - ix_{\eta, \text{aj}}^*|^2 - |x_\eta|^2 - |x_{\eta, \text{aj}}|^2)/2 \\ &= (p_{\eta, \text{sum}} - p_\eta - p_{\eta, \text{aj}})/2 = d_\eta(0)/2, \end{aligned} \quad (2)$$

where the last equation indicates the equivalence of “digital subtraction,” shown in Fig. 1 and our proposed “analog update” scheme $d_\eta(0)/2$ in Fig. 2(c, d) (Appendix). Pseudocode and the complete enumerated backpropagation protocol are discussed in the Appendix. Note that the digital and analog gradient updates can both be implemented in parallel across all photonic layers of the network.

Now that we have defined the analog *in situ* backpropagation update, we experimentally evaluate the accuracy of the analog gradient measurement for a matrix optimization problem in Fig. 2(b, d). Since our circuit does not have an explicit backprop unit architecture, we experimentally simulate the “backprop unit” of Fig. 2(c) by programming a sequence of summing vectors in the Generator unit of our chip and recording $d_\eta(\zeta)$ to compute the gradient with respect to η . We implement backpropagation in a single photonic mesh layer optimizing a linear cost function $\mathcal{L}_m = 1 - |\hat{\mathbf{u}}_m^T \mathbf{u}_m^*|^2$, where \mathbf{u}_m is row m of U , a target matrix that we choose to be the four-point discrete Fourier transform (DFT), and $\hat{\mathbf{u}}_m$ is row m of \hat{U} , the implemented matrix on the device. Each phase shifter in the photonic network implements a phase shift $\theta + \delta\theta$, where θ is the optimal phase shift for U and $\delta\theta$ is some random phase error with standard deviation $\sigma_{\theta, \phi}$, which can serve as a measure of “distance to convergence” during training of the device. For our gradient measurement step, we send in the derivative $\mathbf{y}_{\text{aj}} = \frac{\partial \mathcal{L}_m}{\partial \mathbf{y}} = -2(\hat{\mathbf{u}}_m^T \mathbf{u}_m^*)^* \mathbf{e}_m$ to achieve an adjoint field \mathbf{x}_{aj} , where \mathbf{e}_m is the m th standard basis vector (1 at position m , 0 everywhere else). We find in Fig. 2(f) that analog gradient measurement is increasingly less accurate when calculated near convergence, likely due to uncorrected photonic circuit error (e.g. due to loss and/or thermal crosstalk) resulting in large gradient measurement errors.

PHOTONIC NEURAL NET TRAINING

To test overall training within our photonic mesh chip, we assess the accuracy of *in situ* backpropagation in Fig. 3 to train L -layer photonic neural networks to solve multiple 2D classification tasks using the digital subtraction protocol in Ref. 5. The classification problem assigns points in 2D space to a 0 or 1 label (red or blue coloring) based on whether the point is in a region of space, and the neural network implements the nonlinear boundary (for instance circle-, moon- or ring-shaped) separating points of different labels standardized using the Python package Sklearn and specified in our code [26]. The points are randomly synthetically generated and are “noisy,” meaning some training example points have a small probability of being assigned a label despite being on the wrong side of the ideal boundary.

To solve this task, we use a three layer PNN where each linear layer uses 4 optical ports (4×4 MVM), i.e. $L = 3$ with $N = 4$ inputs and outputs. The inference

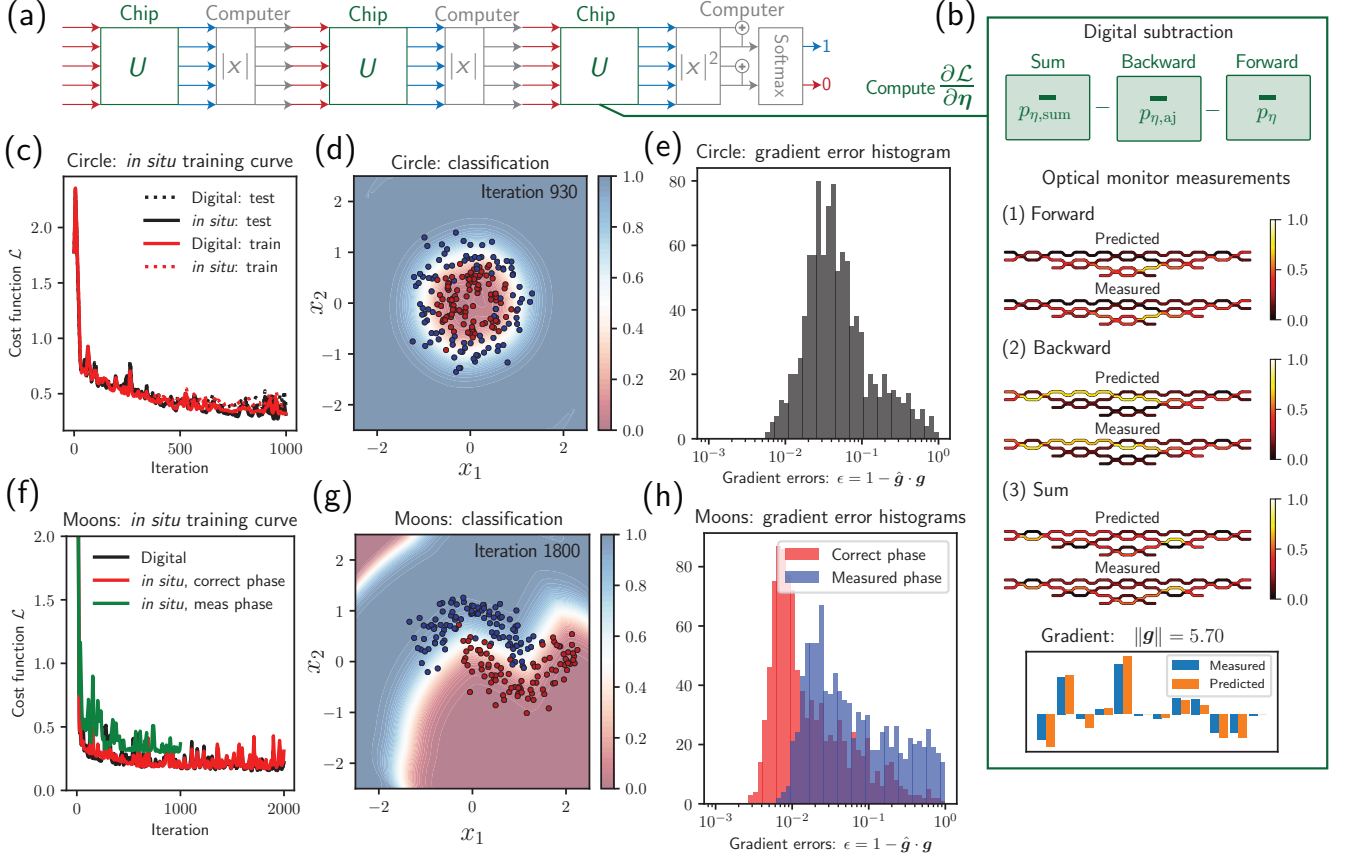


FIG. 3. We perform *in situ* backpropagation training based on stochastic gradient descent and Adam update [25] (learning rate 0.01) for two classification tasks solvable by (a) a three layer hybrid photonic neural network consisting of absolute value nonlinearities and a softmax (effectively sigmoid) decision layer. For the circle dataset, (b) the training curve for Adam update for stochastic gradient descent shows excellent agreement between test and train for both digital and *in situ* backpropagation updates resulting in (c) a classification curve showing the true labels and the classification curve based on the learned parameters resulting in 96% model test accuracy and 93% model train accuracy. (d) For the moons dataset, we compare the test cost curves for the digital and the *in situ* updates, where we find that our phase measurements are sufficiently inaccurate to impact training leading to a lower model train accuracy of 87%. If we use ground truth phase measurements (red) instead of measured phases (blue), we ultimately arrive at (e) a sufficiently high model test accuracy of 98% (model train accuracy is lower at 95%). When using ground truth phases instead of measured phases, (f) the gradient error reduces considerably by roughly an order of magnitude. (g) While training the circle classification was successful, the measured gradient error is similarly large as that measured in the moons training experiment. This suggests that the importance of accurate gradients can be problem-dependent.

operation of our photonic neural network consists of programming the inputs in the red Generator circuit and measuring outputs on the blue Analyzer circuit, reprogramming the unitary for each Matrix unit layer on the same chip, and square-rooting the output power measurement to achieve absolute value nonlinearities of the form $|\mathbf{y}|$ (see Appendix for more detailed description). This unitary layer reprogramming is only intended for a proof-of-concept; the ultimate implementation would dedicate a separate optical device to each linear layer.

The neural network inference model outputs probability of 0 or 1 (red or blue assignment) of each point based on the following model:

$$\hat{\mathbf{z}}(\mathbf{x}) = \text{softmax2}(|U^{(3)}|U^{(2)}|U^{(1)}\mathbf{x}|) \quad (3)$$

where we define the softmax2 $: \mathbb{C}^4 \rightarrow [0, 1]^2$ as two-element vector representing the probability of 0 or 1 label to be $\text{softmax2}(\mathbf{y}) = (e^{|y_1|^2 + |y_2|^2}, e^{|y_3|^2 + |y_4|^2}) / (e^{|y_1|^2 + |y_2|^2} + e^{|y_3|^2 + |y_4|^2})$. We then apply a softmax cross entropy (SCE) cost function $\mathcal{L}(\mathbf{x}) = \text{SCE}(\hat{\mathbf{z}}(\mathbf{x}), \mathbf{z}) = z_0 \log \hat{z}_0 + z_1 \log \hat{z}_1$. The ultimate goal is to apply automatic differentiation and *in situ* analog gradient measurement on our photonic device to optimize \mathcal{L} .

Input data to our device is formatted into the form (x_1, x_2, p, p) , where x_1, x_2 correspond to the location in 2D space and p is some power ensuring that all inputs are normalized to the same power P , i.e. $x_1^2 + x_2^2 + 2p^2 = P$; this convention follows the simple example in Ref. 5.

We perform a 80/20% train-test split (200 train points, 50 test points), holding out test data from training to ensure no “overfitting” takes place, though this is unlikely for our simple model. We generally find higher test than train accuracy in our results since there are fewer “noisy examples” in our randomly generated test sets.

Our single photonic chip is used to perform the data input, data output and matrix operations for all three layers of our photonic neural network shown in Fig. 3(a). After each pass through the photonic chip, we measure the output power and digitally perform a square-root operation to effectively implement absolute value nonlinearities on the computer. Throughout the process (forward, backward and sum steps), we also perform digital simulations so we can compare the experimental and simulated performance at each step as shown in Fig. 1(b). Minimizing the cost \mathcal{L} ultimately leads to maximum accuracy in classifying points to the appropriate labels.

When performing training, the most critical information is in the gradient direction, so we compute gradient direction error using $1 - \mathbf{g} \cdot \hat{\mathbf{g}}$ comparing normalized measured and predicted $\mathbf{g} = \partial\mathcal{L}/\partial\boldsymbol{\eta} \cdot \|\partial\mathcal{L}/\partial\boldsymbol{\eta}\|^{-1}$. An important distinction to make for metric reporting is the difference in “model” versus “device” cost function and accuracy. In Fig. 3, we report “model metrics” by evaluating device parameters learned on our chip on the true model. Thus, actual training of the physical parameters is performed on the device itself, and the result of training is evaluated on the computer.

Our first task is to verify that inference works on our platform, which we show for a randomly generated “ring” dataset to have 90% device test set accuracy on our physical platform shown previously in Fig. 1(c). Once we confirm that the inference performance is acceptable, we then perform training of 2D classification problems using our digital subtraction approach on our randomly generated datasets. We use standard gradient update Adam stochastic gradient descent [25] with a learning rate of 0.01, with all non-linear automatic differentiation performed off the chip via Python libraries JAX and Haiku [27, 28].

We first report our training update model metrics for the circle dataset for the photonic neural network shown in Fig. 3(a). In Fig. 3(b), we show the grating tap-to-camera measurements of normalized field magnitudes in the final layer for all three passes required for digital subtraction across all layers of our device at iteration 930 (near the optimum), which show excellent agreement between predicted and measured fields. The training curves in Fig. 3(c) indicate that stochastic gradient descent is a highly noisy training process due to the noisy synthetic dataset about the boundary; this phenomenon can be observed for both the digital and analog approaches. Due to these outliers, we only observe convergence in the time (iteration)-averaged curves because even at convergence, updates based on outliers or some incorrectly labelled points can result in large swings in the cost function. These large swings appear roughly correlated between

the simulated and measured training curves. Despite these swings and phase shift gradient errors shown in Fig. 2(e), our results of 96% model test accuracy and 93% model train accuracy indicate successful training as shown in Fig. 2(d).

We then train a moons dataset, where we apply the same procedure to achieve a model train accuracy of 87% and model test accuracy of 94%, which suggests that training occurs but there is room for improvement as shown in green in Fig. 3(f). Upon further investigation, we find that if we use the ground truth phase for the phase measurement but keep the amplitude measurements, we reduce the phase shift gradient error by roughly an order of magnitude on average as shown in Fig. 3(h). This results in the successfully trained classification of Fig. 3(g) and the red curve in in Fig. 3(f) which shows excellent correlation with the black digital training curve. When using the corrected ground truth phase measurement, we achieve a model train accuracy of 95% and model test accuracy of 97% for the full backpropagation demonstration based on measured phase (stopped early at 1000 iterations), an improvement that underscores the importance of accurate phase measurement for improved training efficiency.

DISCUSSION AND OUTLOOK

In this paper, we have laid the foundation for the analysis of our new *in situ* backpropagation proposal and tolerance to gradient errors for the design of practically useful photonic mesh accelerators. Our proof-of-principle experiments suggests that even in the presence of such gradient error, gradient measurement and training photonic neural networks using analog backpropagation updates is efficient and feasible.

Although there exist many approaches for training photonic neural networks, our demonstration and energy calculations (Appendix) suggest that *in situ* backpropagation is the most practical and efficient approach for training deep multilayer hybrid photonic neural networks. Our hybrid approach to training optically accelerates the most computationally intensive operations (both in energy and in time complexity), specifically $O(N^2)$ matrix-vector products and matrix gradient computations (backward pass). On the other hand, all other $O(N)$ computations such as nonlinearities and their derivatives are implemented on the computer directly, which is reasonable because $O(N)$ time is needed to modulate and measure optical inputs and outputs anyway. Other techniques such as population-based methods [29], direct feedback alignment [30, 31], and perturbative approaches require fewer components to implement but are ultimately less efficient for training deep neural networks compared to backpropagation.

Our main finding is that gradient accuracy plays an important role in reaching optimal results during training. As we find in Fig. 3, more accurate gradients result

in training convergence speeds and oscillations comparable to digital calculations of the gradients updated over the same training example sequence. This accuracy is vital for *in situ* backpropagation to be a viable competitor to existing purely digital training schemes; in particular, even if individual updates are faster to compute, high error would result in longer training times that mitigate that benefit. In the Appendix, we frame this error scaling in terms of a larger scale PNN simulation on the MNIST dataset originally as explored in Ref. [32], where we consider errors in gradient measurement due to optical I/O errors and photodetector noise at the global monitoring taps.

Our findings ultimately have wide ranging implications because backpropagation is the most efficient and widely used neural network training algorithm in conventional machine learning hardware used today. Our analog approach for machine learning thus opens up a vast opportunity for energy-efficient artificial intelligence applications using photonic hardware. We additionally provide seamless integration into current machine learning training protocols (e.g. autodifferentiation frameworks such as JAX [27] and TensorFlow [33]). A particularly impactful opportunity is in data center machine learning where optical signals already store data that can be fed into PNNs for inference and training tasks. In such settings, our demonstration presents a key new opportunity for both inference and training of hybrid PNNs to dramatically reduce carbon footprint and counter the exponentially increasing costs of AI computation.

ACKNOWLEDGEMENTS

We would like to acknowledge Advanced Micro-Foundries (AMF) in Singapore for their help in fabricating and characterizing the photonic circuit for our demonstration and Silitronics for their help in packaging our chip for our demonstration. We would also like to acknowledge funding from Air Force Office of Scientific Research (AFOSR) grants FA9550-17-1-0002 in collaboration with UT Austin and FA9550-18-1-0186 through which we share a close collaboration with UC Davis under Dr. Ben Yoo. Thanks also to Payton Broaddus for helping with wafer dicing, Simon Lorenzo for help in fiber splicing the fiber switch for bidirectional operation, Nagaraja Pai for advice on electrical and thermal control packaging, and finally Carsten Langrock and Karel Urbanek for their help in building our movable optical breadboard.

DATA AND SOFTWARE

All software and data for running the simulations and experiments are available through Zenodo [34] and Github through the Phox framework, including our ex-

perimental code via Phox [26], simulation code via Simphox [35], and circuit design code via Dphox [36].

CONTRIBUTIONS

SP taped out the photonic integrated circuit and ran all experiments with input from ZS, TH, TP, BB, NA, MM, OS, SF, DM. SP and ZS wrote code to control experimental device. TP designed the custom PCB with input from SP. SP wrote the manuscript with input from all coauthors. All coauthors contributed to discussions of the protocol and results.

CONFLICTS OF INTEREST

SP, ZS, TH, IW, MM, SF, OS, DM have filed a patent for the analog backpropagation update protocol discussed in this work with Prov. Appl. No.: 63/323743. The authors declare no other conflicts of interest.

METHODS

A. Circuit design and packaging

Our photonic integrated circuit is a 6×6 triangular photonic mesh consisting of a total of 15 MZIs fabricated at the AdvancedMicroFoundry (AMF) in Singapore designed using our photonic library DPhox [36] which is a custom automated photonic design library in Python. Each of the MZIs in the mesh is controlled using programmable phase shifters in the form of $80 \mu\text{m} \times 2 \mu\text{m}$ titanium nitride heaters with 10.5 ohm/sq sheet resistance surrounded by deep trenches that are $80 \mu\text{m} \times 10 \mu\text{m}$ and a total of $7 \mu\text{m}$ away from the waveguide, which use resistive heating to control the interference of light propagating in the chip. The MZIs consist of two 50/50 directional couplers, with S-bends consisting of $30 \mu\text{m}$ radius arc turns and $40 \mu\text{m}$ long interaction lengths with a 300 nm gap. Next to each of the phase shifters is a bidirectional grating tap monitor, which is a directional coupler tap that couples 3% of the light propagating either forward or backward through the waveguide attached to the tap and feeds that light to a grating to be imaged on a camera focused on the grating. Traces for one of the terminals of each of the phase shifters are routed to separate individual pads on the edge of the chip, and the ground connections across all phase shifters in a column of MZIs are shared and connected to a single ground pad. The trace widths need to be thick enough to handle high thermal currents, so we use $15 \mu\text{m}$ wide traces and $15N_{\text{wire}} \mu\text{m}$ wide traces when multiple connections are connected to a shared ground contact.

The photonic chip is attached using silver paint to a 1.5mm thick copper shim and a custom Advanced Circuits PCB designed in KiCAD consisting of ENIG

coated metal traces to interface the phase shifters with an NI PCIe-6739 controller for setting programmable phase shifts throughout the device. Our PCB is wirebonded using two-tier wirebonding to the chip by Silitronics Solutions, made possible by fanout to NI SCB-68 connectors that interface directly to our PCIe-6739 system. The input optical source is a Agilent 81606A tunable laser with a tunable range of 1460 nm to 1580 nm. The laser light is coupled into a single-mode fiber and optically interfaced to the chip using W2 Optronics 127 micron pitch fiber array interposers at the left and right sides of the mesh, with a mirror facet designed to couple optical signals at 10 degrees from the normal as we only need to couple into a single grating coupler for each fiber array coupler. Optical stray reflections from light not coupled into the chip generally interfere with grating tap signals forming extra streaks in the camera; these stray reflections are blocked using pieces of paper carefully placed above the fiber arrays that act as lightweight removable stray light blockers.

For thermal stability, this chip-PCB assembly is thermally connected to a thermoelectric cooler (TEC). This thermal connection is made possible by metal vias connecting rectangular ENIG-coated copper patches on the top of the PCB to the bottom of the PCB, with thermal paste between an aluminum heat sink mount and the bottom rectangular metal patch. For feedback control, a thermistor placed near the chip and the TEC under the chip are attached to a TEC controller unit, allowing stable chip temperature (kept at 30°C) for training.

B. Optical rig design

Our optical rig consists of an Ethernet cable-connected Xenics Bobcat 640 IR camera and microscope assembly mounted on an XY stage and six-axis stages for free space fiber alignment. The IR camera and microscope image individual grating taps throughout a photonic integrated circuit (PIC) and is responsible for all measurement on the chip (both optical I/O and optical gradient monitoring).

The microscope uses an ∞ -corrected Mitutuyo IR 10x objective and a 40cm tube lens leading to a dichroic connected to visible and IR optical paths for simultaneous visible and infrared imaging. The optical rig is also outfitted with additional paths for LEDs to illuminate the actual chip features. This allows us to find the optimal focus for the grating spots, an image shown in Fig. 2(a). In order to measure intensities directly using the IR camera, the Bobcat camera “Raw” mode is turned on and autogain features are turned off. The integration time is set to 1 millisecond, and the input laser power is set to 3 mW; note that higher integration times are required for lower input laser powers. We take an initial reference image to get a baseline and then to measure the spots intensities or powers, we sum up the pixel values that “fill” the appropriate grating taps throughout the device. The

triangular mesh circuit is constructed such that the grating taps lie along columns of devices, which means the optical rig images a 6×19 array of spots. The infrared path has roughly a $700 \times 600 \mu\text{m}$ field of view, allowing simultaneous measurements of 6×3 grating spots on the chip (MZIs are $625 \mu\text{m}$ long in total given roughly $165 \mu\text{m}$ long directional couplers), which necessitates an XY translation stage to image multiple spots simultaneously on the chip.

The speed of backpropagation is limited by the mechanics of the XY stage required to image spots throughout the chip, so our demonstration training experiments took up to 31 hours of real time to run, limited primarily by the wait time for the stage to settle on various groups of spots on the chip. Assuming T iterations, the stage needs to move a total of $15T$ times (5 for each of the three *in situ* backpropagation steps to be able to image all of the spots). For 1000 iterations, the stage needs to move a total of 15000 times which necessitates the need of automation for the stage of our proof-of-concept demonstration. In a final commercial implementation, the grating taps would be replaced by integrated photodetectors; there would in principle be no separate optical rig system in a fully packaged hybrid digital-analog photonic circuit.

C. Forward inference operation

Forward inference proceeds as follows for layer ℓ (see Fig. 1(a) in the main text) where each step is $O(N)$:

1. Compute the sets of phase shifter settings $\theta_X^{(\ell)}, \phi_X^{(\ell)}$ for the Generator to give the desired vector $\mathbf{x}^{(\ell)}$ of complex input amplitudes for the Matrix unit in layer ℓ .
2. Set these as the actual phase shifts in the Generator phase shifters using calibration curves for $v_\theta(\theta), v_\phi(\phi)$ and shine light into the Generator circuit to create the corresponding actual vector of optical input amplitudes for the Matrix unit.
3. After the propagation of light through the Matrix unit, the system has optically evaluated the vector of complex optical output amplitudes $\mathbf{y}^{(\ell)} = U^{(\ell)}\mathbf{x}^{(\ell)}$. Now self-configure [37] the output Analyzer circuit to give all the output power in the “top” output waveguide, and note the corresponding sets of voltages v_ϕ and voltages v_θ now applied to each phase shifter in the Generator circuit.
4. Deduce the phase shifts $\theta_Y^{(\ell)}, \phi_Y^{(\ell)}$ in the Analyzer circuit using calibration curves for $\theta(v_\theta), \phi(v_\phi)$, and hence compute the corresponding measured output amplitudes $\mathbf{y}^{(\ell)}$.
5. Compute $\mathbf{x}^{(\ell+1)} = f^{(\ell)}(\mathbf{y}^{(\ell)})$ on the computer.

The first four steps are also used in cases where light is sent backwards (see Fig. 1(g, h)), switching the role

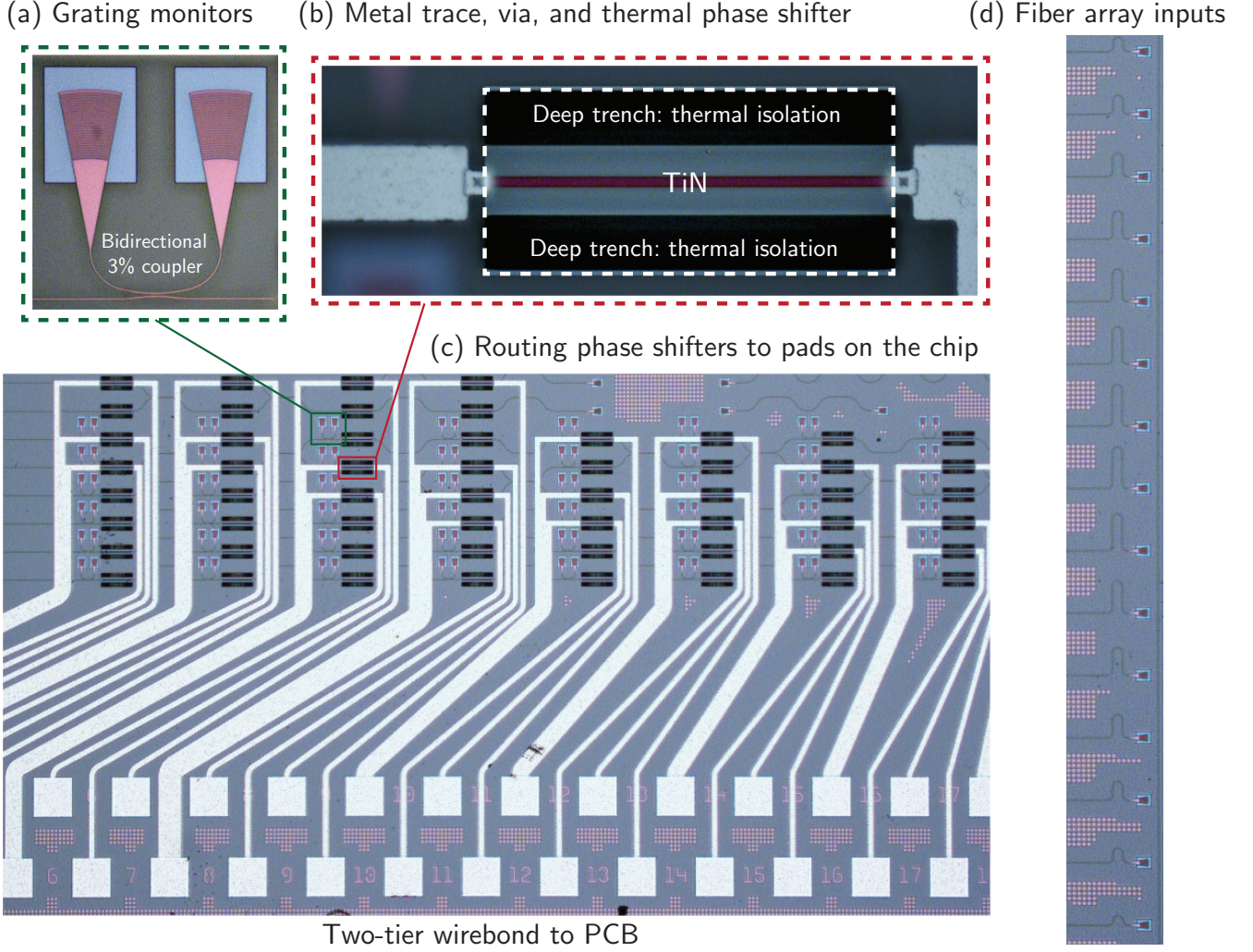


FIG. 4. Microscope images of the photonic mesh used in this paper. (a) Grating monitor closeup showing the bidirectional grating tap we use to perform the backpropagation protocol. (b) Metal trace, via, and TiN (titanium nitride) phase shifter is colocated with the grating monitor and is used to control the interference by changing optical phase in the mesh programmatically. Deep trenches are used for thermal isolation. Here, we show an overlay of phase shifter focal plane on the top metal trace and via used to connect each phase shifter to the pads. (c) A large scale view of a section of the chip (d) Fiber array inputs to the photonic mesh are spaced $127 \mu\text{m}$ apart and are used for interfacing fiber arrays.

of the input and output vector units from Generator to Analyzer and vice versa. Pseudocode for the forward operation of the PNN is provided in the Appendix, and code for the actual implementation is provided in our photonic simulation and control framework Phox [?].

D. Backpropagation protocol

For each training example (\mathbf{x}, \mathbf{z}) , we calculate gradient updates to phase shifts $\boldsymbol{\eta}$ using a “backward pass” corresponding to the inference “forward pass” for that data. More formally, we define a “vector-Jacobian product” or VJP for each function $U^{(\ell)}, f^{(\ell)}$ to algorithmically compute the gradient of our cost function \mathcal{L} . As shown in

Fig. 1(c), each transformation from the forward step is mapped to a VJP in the corresponding backward step (defined in decreasing order from layer L to 1) which depends on intermediate function evaluations in both forward and backward passes. The *in situ* backpropagation step implements the costly intermediate VJP evaluations (i.e. matrix multiplications) directly in the analog optical domain. We define the VJP for nonlinearity $f^{(\ell)}(\mathbf{y}^{(\ell)})$ as $f_{\text{vjp}}^{(\ell)}(\mathbf{y}^{(\ell)}, \mathbf{x}_{\text{aj}}^{(\ell+1)})$:

$$\begin{aligned} \mathbf{y}_{\text{aj}}^{(\ell)} &= f_{\text{vjp}}^{(\ell)}(\mathbf{y}^{(\ell)}, \mathbf{x}_{\text{aj}}^{(\ell+1)}) \\ \mathbf{x}_{\text{aj}}^{(\ell)} &= (U^{(\ell)})^T \mathbf{y}_{\text{aj}}^{(\ell)} \end{aligned} \quad (4)$$

Finally, we synthesize Eqs. 1 and 4 and the results of Ref. 5 to get the backpropagation update based on

applying the chain rule evaluating the cost function at a random training example $\mathbf{x}_t, \mathbf{z}_t$ at iteration t :

$$\begin{aligned}
\frac{\partial \mathcal{L}}{\partial \boldsymbol{\eta}^{(\ell)}} &= \frac{\partial \mathbf{y}^{(\ell)}}{\partial \boldsymbol{\eta}^{(\ell)}} \overbrace{\frac{\partial \mathbf{x}^{(\ell+1)}}{\partial \mathbf{y}^{(\ell)}} \cdots \frac{\partial \hat{\mathbf{z}}}{\partial \mathbf{y}^{(L)}} \frac{\partial \mathcal{L}}{\partial \hat{\mathbf{z}}}}^{\mathbf{x}_{\text{aj}}^{(\ell)}} \bigg|_{\mathbf{x}_t, \mathbf{z}_t} \\
&= \underbrace{\frac{\partial \mathbf{y}^{(\ell)}}{\partial \boldsymbol{\eta}^{(\ell)}}}_{D_\ell \times N \text{ Jacobian}} \cdot \underbrace{\mathbf{x}_{\text{aj}}^{(\ell)}}_{N \times 1 \text{ vector}} \quad (5) \\
&= \underbrace{(\mathbf{x}^{(\ell)})^T \frac{\partial U^{(\ell)}}{\partial \boldsymbol{\eta}^{(\ell)}} \mathbf{x}_{\text{aj}}^{(\ell)}}_{\text{"optical VJP"}} = \underbrace{-\mathcal{I}(x_{\boldsymbol{\eta}^{(\ell)}}) x_{\boldsymbol{\eta}^{(\ell)}, \text{aj}}}_{D_\ell \times 1 \text{ in situ gradient}} \\
\boldsymbol{\eta}_t &:= \boldsymbol{\eta}_{t-1} + \alpha \frac{\partial \mathcal{L}}{\partial \boldsymbol{\eta}} \bigg|_{\mathbf{x}_t, \mathbf{z}_t}
\end{aligned}$$

where $x_{\boldsymbol{\eta}}$ represents a vector of intermediate fields at the input of phase shifters in layer ℓ $\boldsymbol{\eta}^{(\ell)}$ at iteration t , D_ℓ is the number of phase shifts parametrizing the device at layer ℓ , \mathcal{I} refers to imaginary part, and α is the learning rate. The main idea is that if enough training examples are supplied (i.e., after T updates), the device will automatically discover or “learn” a function that performs the task we desire.

Based on Eq. 4, the steps of our optical VJP step, as depicted in Fig. 1(c), is as follows in order from layer $\ell = L$ to 1 of the photonic neural network:

1. Compute the “adjoint” vector $\mathbf{y}_{\text{aj}}^{(\ell)} = f_{\text{vjp}}^{(\ell)}(\mathbf{y}^{(\ell)}, \mathbf{x}_{\text{aj}}^{(\ell+1)})$. For the last layer, set $\mathbf{y}_{\text{aj}}^{(\ell)}$ to be the *error signal* $\mathbf{y}_{\text{aj}}^{(L)} = (\frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(L+1)}})^*$.
2. Perform the backward “adjoint” pass $\mathbf{x}_{\text{aj}}^{(\ell)} = U^T \mathbf{y}_{\text{aj}}^{(\ell)}$ by sending light backwards through layer ℓ of the mesh and measuring the resulting vector of amplitudes $\mathbf{x}_{\text{aj}}^{(\ell)}$ emerging backwards from the mesh.
3. Send the vector of optical amplitudes $\mathbf{x}^{(\ell)} - i(\mathbf{x}_{\text{aj}}^{(\ell)})^*$ forward into layer ℓ of the mesh.
4. Measure gradient $\partial \mathcal{L} / \partial \theta$ for any phase shifter θ :
 - (a) If using digital subtraction measurement [5], measure the sum power $p_{\theta, \text{sum}}$ and subtract p_θ and $p_{\theta, \text{aj}}$ (monitored power from forward and backward steps) to get the gradient.
 - (b) If using analog gradient measurement, sweep the adjoint global phase ζ (giving $\mathbf{x}^{(\ell)} - i(\mathbf{x}_{\text{aj}}^{(\ell)})^* e^{i\zeta}$) from 0 to 2π repeatedly (e.g., using a sawtooth signal). Measure $d_\theta(\zeta)$, the AC component of the measured power through phase shifter θ , $p_{\theta, \text{sum}}(\zeta)$. The gradient is $d_\theta(0)/2$.

5. Update $\boldsymbol{\eta}$ using measured gradients $\partial \mathcal{L} / \partial \boldsymbol{\eta}$.

Note that Step 1 can be simplified to $\mathbf{y}_{\text{aj}}^{(\ell)} = (f^{(\ell)})'(\mathbf{y}^{(\ell)}) \odot \mathbf{x}_{\text{aj}}^{(\ell+1)}$ in the case that $f^{(\ell)}$ is holomorphic, or complex-differentiable. In this paper for the neural network parametrized by Eq. 3, we specifically care about the nonlinearity $f^{(\ell)}(\mathbf{y}) = |\mathbf{y}|$, which has the associated VJP:

$$f_{\text{vjp}}^{(\ell)}(\mathbf{y}, \mathbf{x}_{\text{aj}}) = \frac{\mathbf{y}}{|\mathbf{y}|} \cdot \mathcal{R}(\mathbf{x}_{\text{aj}}). \quad (6)$$

The other VJP required to calculate $\partial \mathcal{L} / \partial \mathbf{y}^{(L)}$ from the final softmax cross entropy and power measurement at the end of the network is handled by our automatic differentiation framework JAX [27, 28].

Steps 3 and 4 can be parallelized over all layers (i.e., parameters of the network) for both the digital and analog update schemes. Pseudocode for the overall protocol (using digital subtraction), along with an energy-efficient proposal for analog gradient computation, is discussed in the Appendix. The final step can be achieved using “stochastic gradient descent” (which independently updates the loss function based on randomly chosen training examples) or adaptive learning where the update vector depends both on past updates and the new gradient. A successful and commonly used implementation of this, which we use in this paper, is called the Adam update [25].

Appendix A: Energy and latency analysis

In this section, we justify why the analog *in situ* update discussed in the main text may be chosen over the digital update proposed in Ref. [5] and used in our main backpropagation training demonstration.

In our hybrid scheme, most of the computation is concentrated in sending in N input modes using modulators (each taking energy E_{inp}) and digital-analog converters and measuring the N output mode powers and amplitudes using photodetectors and analog-digital converters (each taking energy E_{meas}). Therefore, the various approaches for a given matrix-vector product cost roughly $N \cdot (E_{\text{inp}} + E_{\text{meas}})$, equivalent to the cost for setting up the input/output behavior for the photonic mesh. A digital electronic computer, on the other hand, requires N^2 sequential operations (i.e., multiply-and accumulate operations that are *not* parallel) to compute any given matrix-vector product.

Beyond inference tasks, the additional backward and sum steps required for *in situ* backpropagation adds additional energy and latency contributions. The analog update explored in Fig. 2 requires N^2 optoelectronic units for energy-efficient operation, each of which is outfitted with a photodetector, a lock-in amplifier, and high-pass filter consuming energy E_{grad} to measure $d_\theta(0)$ for a total energy consumption of $N^2 E_{\text{grad}} + N \cdot (3E_{\text{inp}} + 2E_{\text{meas}})$ for all three steps of the full backpropagation measurement. The $2E_{\text{meas}}$ comes from the output measurements

in the first two steps, and the E_{grad} comes from an analog gradient measurement in the final step.

In comparison, the digital subtraction described in Fig. 1 can be useful in adaptive updates that require storing information about previous gradients (such as Adam [25] which we exploit for training), but there are a couple of drawbacks. First, a digital update is less memory efficient since N^2 elements need to be stored using analog memory to be able to run the “digital subtraction” computation in backpropagation. Additionally, the total energy consumption becomes $3N^2 E_{\text{grad,digital}} + N \cdot (3E_{\text{inp}} + 2E_{\text{meas}})$, with $E_{\text{grad,digital}} \gg E_{\text{grad,analog}}$ due to large numbers of analog-digital conversions required to implement the analog-digital conversions and digital subtraction calculations. Analog-digital conversions are among the most energy- and time-consuming operations in a hybrid photonic device; when operating at GHz speeds, the best individual comparators generally require up to 40 fJ [38, 39] (versus around 1 fJ/bit for input modulators [40]) and therefore should ideally be reserved for optical input/output in the photonic meshes.

A final energy consideration is the phase shift modulation. These voltage-controlled modulators may be controlled by thermal actuation [41], microelectromechanical (MEMS) actuation [42] or phase-change materials such as barium titanate (BTO) [43]. Of these options, MEMS actuation is among the most promising because unlike thermally actuated phase shifters, they cost no energy to maintain a given programmed state (“static energy”), dramatically improving the energy efficiency of operation compared to thermal phase shifters which constantly dissipate large amounts of heat. Additionally, unlike phase change materials, MEMS phase shifters use CMOS materials such as silicon or silicon nitride. Furthermore, such devices can be designed to operate in the linearly with voltage [44, 45] which ensures that the gradient update applied to the voltage is the same as that of the phase shift without a calibration curve. This helps with gradient accuracy as we discuss now.

Appendix B: Gradient accuracy

As shown in Fig. 3(f, h) and in Fig. 6(h), gradient accuracy can affect the optimization and decrease as the optimization approaches convergence. As we find in the main text, accurate phase measurement plays an important role in measuring accurate gradients. This is true even when the nonlinearity (as in our case with absolute value) removes the need to measure phases in the inference step. Since in the main text, we evaluate the model accuracy (device-trained parameters evaluated on a theoretical computer model), we also show some evidence that the device and model classifications match quite well in Fig. 6(i, j).

One popular type of update is based on “minibatch gradient descent,” a machine learning technique that calculating gradients based on multiple training examples.

This would dramatically smooth out the noisy training curves shown in Fig. 6(a-f), as the resulting averaged gradients would actually be much smaller in magnitude. However, as shown in Fig. 6(g), we find that the normalized error of a minibatch gradient is generally significantly higher than that of the gradient for a single training example which can have negative implications for training. This is because the variance of the gradient error remains the same when averaged over many examples, but the contribution of the gradient error is much larger over a batch. This phenomenon might be problem-dependent; if the average gradient for the minibatch is not closer to zero than the gradient for individual training examples, this error may not be an issue. This underscores the importance of accurate gradient measurement, which can be improved using more accurate output phase measurements; our output phase measurement alone results in an order-of-magnitude increase in gradient error.

Finally, a linear relationship between phase and voltage can help to improve gradient update accuracy without requiring nontrivial scaling complexity in the hardware. In other words, we ensure $\partial \mathcal{L} / \partial v_\theta = \partial \theta / \partial v_\theta \cdot \partial \mathcal{L} / \partial \theta$ with constant $\partial \theta / \partial v_\theta$ which simplifies the required analog circuitry. The $\partial \theta / \partial v_\theta$ term is calculated using calibration curves, and this assumption is more-or-less valid in our case as we operate the phase shifters in the linear regime as shown in Fig. 8(e).

Appendix C: Usage in machine learning software

Backpropagation is also known as automatic differentiation (AD) because any program that uses backpropagation registers a “backward” gradient function for any forward function, which is used by AD Python engines such as JAX[27], TensorFlow2 [33], and PyTorch. We demonstrate that our protocol can be easily coupled with an existing automatic differentiation framework (JAX and Haiku [27, 28]), which can register a backward step and adaptive update based on Adam [25] for all unitary matrix operations as an analog *in situ* backpropagation gradient calculation rather than an expensive digital operation. In this way, the digital side of our hybrid PNN never needs to store or have any knowledge of parameters in the photonic mesh architecture. However, in cases where adaptive gradient updates are used, such as Adam, aggregated knowledge based on past gradient updates needs to be stored; non-volatile memory may be required to energy-efficiently store these additional parameters.

Appendix D: Comparison with other training algorithms

Backpropagation is the most widely used and efficient known algorithm for training multilayer neural network models, though it is far from the only method for calculating gradient-based updates.

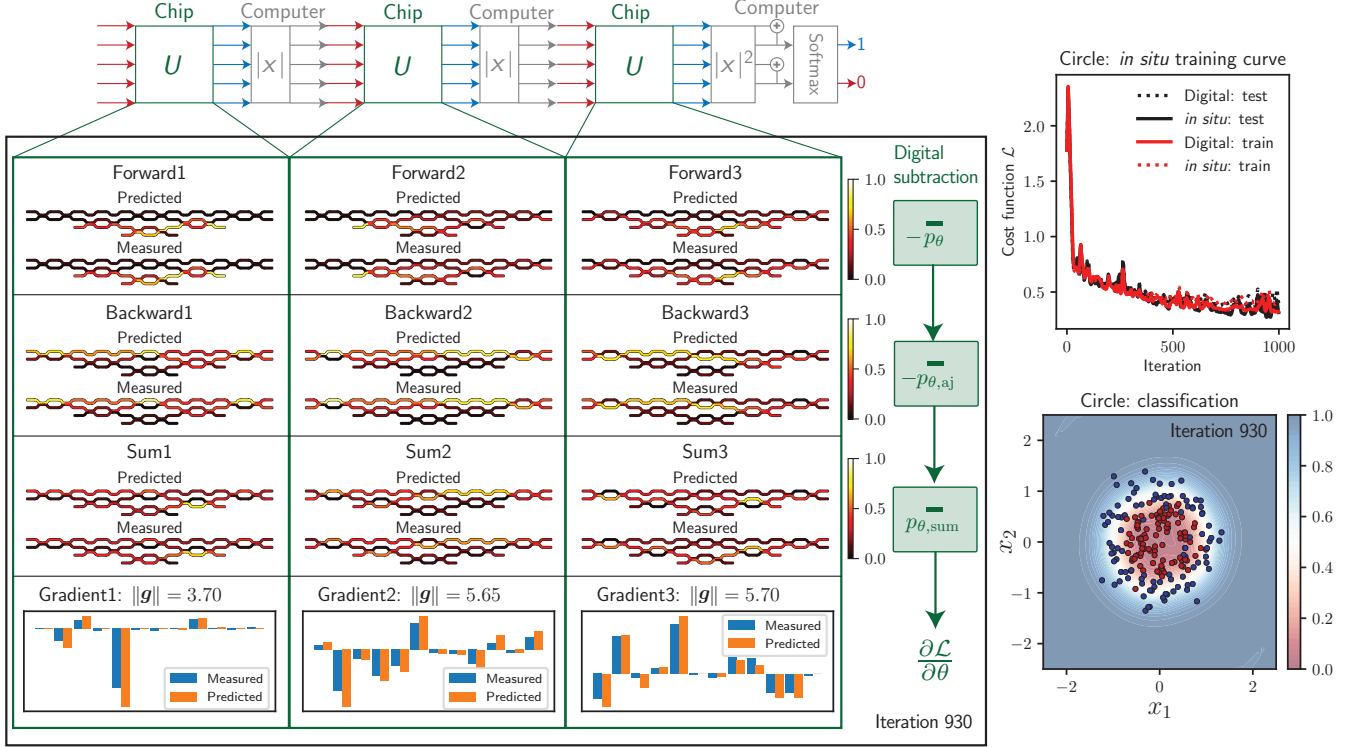


FIG. 5. Here, we simply extend the data shown for iteration 930 in Fig. 3, specifically showing the intermediate power measurements at each point in the photonic neural network. As indicated by “digital subtraction,” we directly subtract the sum measurements by the top forward and backward measurements.

Finite differences-based training has been proposed as a method of training photonic neural networks [7]. Finite differences falls under the umbrella of perturbative learning, a well known and model-free analog machine learning technique for analog neural networks that works by perturbing each element by a small amount, or perturbing many elements simultaneously, and measuring the resulting change in the overall loss function [46–48]. Perturbative learning is most useful in the context of *fully* optical neural networks that implement nonlinearities directly on the device, an example of which has previously been proposed for all-photonic neural networks [32]. The PNN architecture relevant for perturbative learning is therefore fundamentally different from the hybrid PNN we propose that can benefit from *in situ* backpropagation.

It is worth noting that hybrid PNNs can be more versatile and useful to a larger range of traditional AI applications compared to fully analog PNNs [32]. Many complex models (e.g. as transformers, convolutional networks, word embedding layers and recurrent neural networks) used in machine intelligence today are more easily implemented in hybrid rather than all-analog systems due to the sheer complexity and logic implemented in the model architectures.

Additionally, backpropagation is significantly more ef-

ficient than finite differences and other similar adaptive approaches. In backpropagation, the time complexity of the “forward-propagated” inference pass or direct evaluation of the model is roughly the same as that of the “backpropagated” gradient calculation pass. In contrast, a perturbative gradient calculation is significantly more costly since it cannot be computed on a layer-by-layer basis; the forward propagation must continue on to the end of the network, which does not favor our hybrid approach.

Other alternatives to backpropagation include direct-feedback alignment (DFA) [30, 31], derivative-free optimization and population-based learning, which include evolutionary-based (genetic algorithm or GA) [29] and swarm-based methods. The GA and DFA training approaches have been recently experimentally demonstrated to successfully train optical devices at moderately challenging machine learning tasks [29, 31]. However, these are generally regarded to be less efficient at training models compared to backpropagation and are have not proven to scale to more challenging image and word processing machine learning benchmarks like ImageNet [49]. Work is still required to test the scalability of photonic machine learning to solve problems of such complexity as ImageNet.

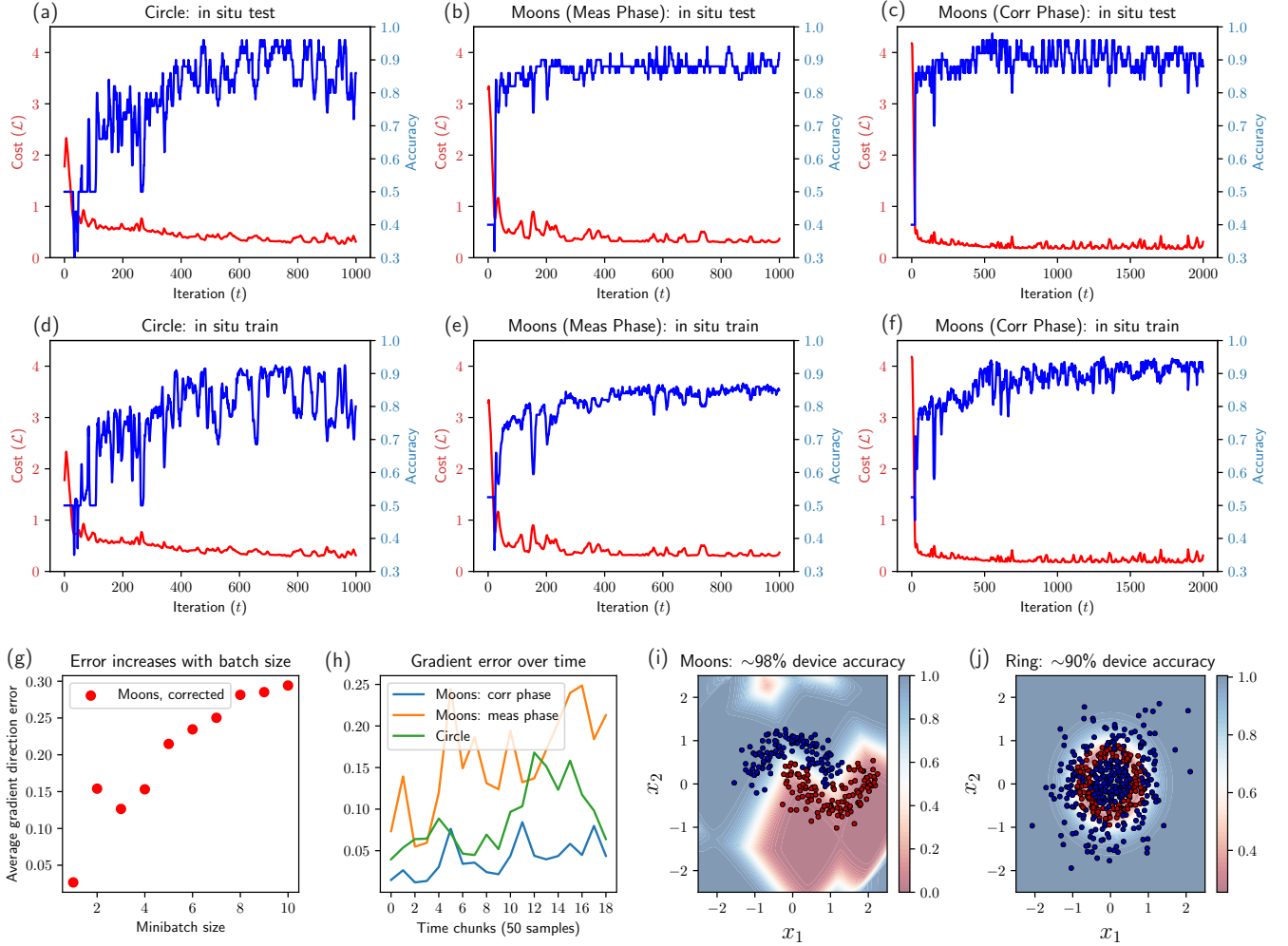


FIG. 6. (a-f) A comparison of the model cost and accuracy curves between circle, moons (measured) and moons (corrected) experiments comparing test (a-c) and train (d-f) data. (g) The error in the gradient increases with the batch size. (h) The gradient error increases over the course of the optimization here shown as a time-averaged series averaged over 50-sample time (iteration) chunks. This increase has to do with the fact that near convergence, the measured gradient is smaller, leading to slightly larger errors. (i, j) Device accuracy for moons and ring dataset inference tasks showing the model boundary (calculated on the computer) in the background and device-classified points (in red, blue). With this evidence, we make an assumption that model and device metrics are close enough, so we use model metrics throughout the paper.

Appendix E: Photonic mesh operation

1. Bidirectional matrix multiplication

In photonic neural networks, programmable photonic meshes act to perform compute-intensive linear operations that preserve the overall power in the form of unitary transmission operator U . Meshes are configured using three subunits: an input vector generator network (generating \mathbf{x}), a matrix network (multiplying by U), and an output vector analyzer network (measuring \mathbf{y}). Our mesh is “bidirectional” in the sense that it can represent matrix-vector operations regardless of whether the light is shined in the forward (left-to-right) or backward (right-to-left) direction as depicted in Fig. 1(a) of the

main text, where in the latter case the output analyzer and input generator switch places.

2. Tunable splitter

A tunable splitter, the basic building block of a photonic mesh, is a 2×2 element that consists of a tunable split ratio region and a differential phase shifter at the input or output. For straightforward calibration, we may use Mach-Zehnder interferometer building blocks that consist of a differential ϕ phase shift, 50/50 splitter, differential θ phase shift, and then a final 50/50 splitter, giving us the following mathematical representation

acting on modes x_1, x_2 and yielding outputs y_1, y_2 :

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = i \begin{bmatrix} e^{i\phi} \sin \frac{\theta}{2} & \cos \frac{\theta}{2} \\ e^{i\phi} \cos \frac{\theta}{2} & -\sin \frac{\theta}{2} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad (\text{E1})$$

$$\mathbf{y} = T_2(\theta, \phi)\mathbf{x},$$

where $\theta \in [0, \pi]$ and $\phi \in [0, 2\pi)$. In practice, due to the nonlinear relationship between the phase shifts θ, ϕ and the respective voltage drives v_θ, v_ϕ , we instead may need to represent T_2 with an additional global phase as:

$$\tilde{T}_2(\theta, \phi) = e^{-i\frac{\theta}{2}} T_2(\theta, \phi) \quad (\text{E2})$$

where we use a single phase shift θ instead of a differential phase shift in the internal phase shift of the MZI. The fundamental function of the MZI is to be able “nullify” (minimize to zero) power in either of its output powers given any input vector. In mathematical terms, given any \mathbf{x} , we should be able to generate an output of the form $\mathbf{y} = (y_1, 0)$. As defined in Refs. 6, 20, we can perform the nullification of y_2 for any MZI $T_2(\theta, \phi)$ with inputs x_1, x_2 :

$$\begin{aligned} \theta &= 2 \arctan \left| \frac{x_1}{x_2} \right| \\ \phi &:= -\arg \left(\frac{x_1}{x_2} \right), \end{aligned} \quad (\text{E3})$$

with the convention for θ, ϕ being internal and external phase shifters as defined in Fig. 1 of the main text.

3. Vector units

As proposed in Ref. 20, a vector unit is a may be used as a $1 \times N$ input vector generator or an $N \times 1$ output vector analyzer (the flipped version of the input generator). An N -vector unit is a “tree network” of $N - 1$ splitters θ_X and N output phases ϕ_X , with the fully balanced binary tree and the maximally unbalanced linear cascade (diagonal line) as extreme cases (see the red and blue structures of Fig. 7(b)). An input generator generates optical modes representing any N -dimensional complex vector given a single input to the system up to a (nonphysical) global phase and can be either balanced or unbalanced as shown in Fig. 7(c). Operated in reverse, the analyzer allows for the $N - 1$ splitters to route all input light into a single port.

The overall mathematics can be represented in either vector or bra-ket notation as follows (where X^\dagger represents analysis and X represents generation):

$$\begin{aligned} \mathbf{x} &= X\mathbf{e}_1 \\ |\mathbf{x}\rangle &= X|0\rangle \\ X^\dagger|\mathbf{x}\rangle &= |0\rangle \end{aligned} \quad (\text{E4})$$

The algorithm required to set up an *output* vector unit analyzer requires first establishing a path between the

root MZI and all other vector unit MZIs (known as a topological order) such that all the light exits the output of the vector unit Refs. 6, 20.

Experimentally, as proven in Refs. 6, 8, 20, this can be achieved using self-configuration by minimizing the power (first sweeping ϕ and then sweeping θ) for $N - 1$ open ports of the device to maximize output port power. All devices belonging to a given column can be programmed simultaneously (in parallel), so for binary tree architectures, this measurement can be done in $O(\log N)$ steps.

In this work, however, since we use a camera for all photodetection measurements, this protocol can be relatively slow. Therefore, we instead use four measurements, with $\theta = \pi/2$ and $\phi = 0, \pi/2, \pi, 3\pi/2$ to deduce the powers $p_0, p_{\pi/2}, p_\pi, p_{3\pi/2}$ and compute relative phase as $\arctan \left(\frac{p_{3\pi/2} - p_{\pi/2}}{p_\pi - p_0} \right)$. This is shown in Fig. 7(g).

Output detection can be made faster if necessary using homodyne coherent detection as shown in Fig. 7(c) where nominally 50 percent of the total input light is split into N waveguides and sent directly to the output of the matrix unit implementing U . In the analog domain this protocol requires only a single step. As with our self-configuration phase measurement protocol, it requires additional computation on the digital end to deduce the phase.

4. Matrix unit

The matrix unit, shown in green in Fig. 1 in the main text, is any suitable arrangement of interferometers needed to represent a subset of unitary matrices in $U(N)$; a *universal* (unitary) matrix unit can implement any unitary matrix in $U(N)$. Examples of universal matrix units are triangular [8, 12], rectangular [50], cascaded binary tree [37], and cosine-sine decomposition [51] (which is more useful for quantum applications of this scheme, but can be represented classically). Such devices consist of $O(N^2)$ parameters: $N(N - 1)/2$ MZIs with 2 phase shifters each θ_U, ϕ_U and N output phase shifters γ_U .

Because multiplying by γ_U is an $O(N)$ operation, all computation for γ_U (both forward and backward passes in the gradient computation) is performed on the computer. In the protocol shown in Fig. 7 and in the main text, we do not include any γ_U phase shifts due to the assumption that those computations are relatively inexpensive and can be fully accounted for off-chip.

The matrix unit is represented by an operator U that performs the following operation (in vector notation and bra-ket notation):

$$\begin{aligned} U\mathbf{x} &= \mathbf{y} \\ U|\mathbf{x}\rangle &= |\mathbf{y}\rangle = Y|0\rangle \end{aligned} \quad (\text{E5})$$

In vector notation, the relative phases given by $\arg\left(\frac{U\mathbf{x}}{\mathbf{x}}\right)$ can be measured only up to an overall phase, so an ad-

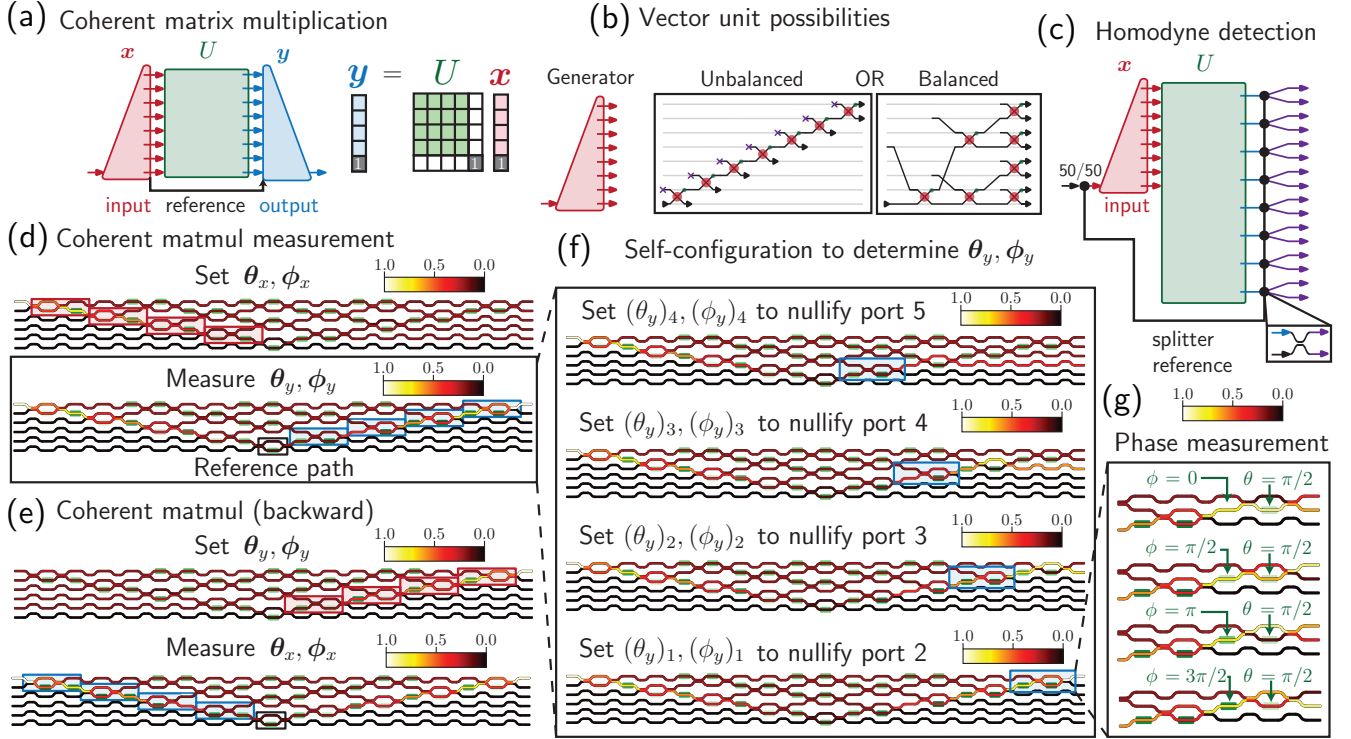


FIG. 7. (a) In our 6×6 MZI network, forward coherent matrix 4×4 multiplication is performed using a generator and analyzer (input/output vector unit) configuration and a fifth reference dimension. (b) A vector unit can be unbalanced or balanced [20]. (c) An alternative and likely faster approach is to use coherent detection or homodyne detection to measure amplitudes and phases. (d) Coherent matmul operation is performed by sending in an input based on the calibrated phase shifts and then performing self-configuration on the output fields including a reference path dimension, and (e) backward coherent matrix multiplication follows the same procedure but backwards. These plots are derived from actual measurements performed on the device and inset colorbars represent powers in the device corresponding to the colored waveguides shown. (f) Self-configuration proceeds by nullifying ports 5 through 2 in descending order. (g) Nullification is achieved using phase measurement rather than analog feedback minimization as in [9] as this is more efficient in our device configuration.

ditional measurement is required to measure this overall phase. In bra-ket notation, we typically can only ensure $\langle 0|Y^\dagger U X|0\rangle = e^{i\phi_0}$, where ϕ_0 is some phase that depends on the effective overall path length in the device, which is a function of all the phase shifts. In theory, we could figure out what this overall path length is by some $O(N^2)$ mathematical computation, but in practice, this can be measured directly in $O(1)$.

5. Reference arm

Phase shifts in physical systems typically have no meaning without a reference, and this is ultimately crucial for designing and programming a photonic mesh. Adding a reference arm waveguide to an N -waveguide photonic mesh (mathematically, embedding all N -dimensional Hilbert space operations in an $N + 1$ -dimensional Hilbert space), an example of which has previously been demonstrated in coherent detection for complex optical neural networks [19].

Independent of reference arm placement, we treat the

unitary operator (U embedded in $N + 1$ -dimensional Hilbert space as shown in Fig. 7(d) for $N = 4$) as follows:

$$\begin{bmatrix} \mathbf{y} \\ z \end{bmatrix} = \begin{bmatrix} U & \begin{matrix} 0 \\ \vdots \\ 0 \end{matrix} \\ 0 & \dots & 0 & 1 \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ z \end{bmatrix}, \quad (\text{E6})$$

which allows us to calculate all phases in the matrix-vector multiplication relative to the phase shift in the added spatial mode (reference waveguide path length). We now can program and/or measure the full input and output \mathbf{x}, \mathbf{y} no matter what settings are used for U . Assuming a total power of 1, the constant phasor z here denotes a constant amplitude, such as $1/\sqrt{N+1}$ or whatever is deemed sufficient.

To properly measure phases for an $N \times N$ operation, we set the phase for the $(N + 1)$ th output of any vector unit as the “reference phase arm” (shown throughout Fig. 7) and connect the reference arm to the waveguide where this phase is defined. If the magnitude of the N th element is zero, we choose that the reference phase of the

vector is also zero. After storing the calibration curve of this reference phase in the computer, we can always set or measure this reference phase by maximizing power output of the reference arm MZI on the appropriate side of the device (e.g. as in the first step of Fig. 7(f)). This is generally a standard technique in phase detection in photonic circuits and similar schemes have been previously explored [19].

Note that in the case of homodyne detection of Fig. 7(c), the math of the phase measurement is a bit different. A separate reference path is still provided, but instead of an analyzer with an additional reference dimension, the reference path is split and interfered at each output to determine the phase. This is a potentially faster and more “standard” method for measuring phases but the circuitry for bidirectional operation is a bit more complex.

6. Calibration

An important protocol for both vector units and matrix units is calibration of phase shifts for accurate inference and phase measurement. For our calibration protocol, we sweep phase shifts while recording an MZI split ratio measured using camera spots immediately after an assigned MZI depending on the calibrated phase shifter. An MZI split ratio can be represented in terms of a transmissivity $t = \sin^2 \theta$, where θ is twice the phase shift in the internal arm, which is used for calibration:

$$t = \frac{p_t}{p} \approx \frac{p_t}{p_r + p_t} \quad (\text{E7})$$

where t is the transmissivity, p is the total power at the input, p_t is the cross state grating power and p_r is the bar state grating power determined by summing up pixel values from the camera.

The model is:

$$\begin{aligned} \theta &= p_0 v^3 + p_1 v^2 + p_2 v + p_3 \\ t &= a \sin \theta + b. \end{aligned} \quad (\text{E8})$$

Empirically, it suffices to fit $v^2 = q_0 \theta^3 + q_1 \theta^2 + q_2 \theta + q_3$ to convert voltage to phase.

For algorithmically calibrating the phase shifts, we use interferometers within the mesh to first calibrate all θ internal phase shifts from left to right by routing light via “lightwires” to all MZIs in the device, as shown in Fig. 8(a) [52].

We then use “meta-MZI” structures within the mesh to calibrate all of the ϕ external phase shifters as shown in Fig. 8(b). For this calibration, after we calibrate each of the ϕ phase shifters, we set $\phi = 0$ so that the other ϕ phase shifter in the meta-MZI has a consistent calibrated phase. Repeating this procedure for all ϕ phase shifts is sufficient to ensure that phase calibrations are all mutually consistent [52].

Algorithm 1 VECTOR UNIT PHASE CONVERSION

```

1: function VEC2PHASE( $\mathbf{x}$ ) ▷ Fig. 7(f)
2:   require  $\mathbf{x} \in \mathbb{C}^N, \|\mathbf{x}\| = 1$ .
3:   for  $m \in [1, 2, \dots, N - 1]$  do
4:      $\phi_m \leftarrow -\arg\left(\frac{x_1}{x_2}\right)$  ▷ Fig. 7(g)
5:      $\theta_m \leftarrow 2 \arctan\left|\frac{x_1}{x_2}\right|$  ▷ nullify at  $m + 1$ 
6:      $x_m \leftarrow e^{i\phi_m} \sin\frac{\theta_m}{2} x_m + \cos\frac{\theta_m}{2} x_{m+1}$ 
7:      $x_{m+1} \leftarrow 0$ 
8:   end for
9:   return  $\boldsymbol{\theta}, \boldsymbol{\phi}$ 
10: end function

11: function PHASE2VEC( $\boldsymbol{\theta}, \boldsymbol{\phi}$ )
12:   require  $\boldsymbol{\theta} \in [0, \pi]^N$ .
13:   require  $\boldsymbol{\phi} \in [0, 2\pi]^N$ .
14:    $\mathbf{x} = [1, 0, \dots, 0] \in \mathbb{C}^N$ 
15:   for  $m \in [1, 2, \dots, N - 1]$  do
16:      $\begin{pmatrix} x_m \\ x_{m+1} \end{pmatrix} \leftarrow \tilde{T}_2(\theta_m, \phi_m)^T \begin{pmatrix} x_m \\ x_{m+1} \end{pmatrix}$ 
17:   end for
18:   return  $\mathbf{x} \exp(-i \arg(x_N))$  ▷ Zero phase for  $x_N$ 
19: end function

```

Algorithm 2 FORWARD STEP

```

1: function MESHFORWARD( $\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\phi}, \boldsymbol{\gamma}$ )
2:   require  $\mathbf{x} \in \mathbb{C}^N, \|\mathbf{x}\| = 1$ .
3:   require  $\boldsymbol{\theta} \in [0, \pi]^{N(N-1)/2}$ .
4:   require  $\boldsymbol{\phi} \in [0, 2\pi]^{N(N-1)/2}$ .
5:   require  $\boldsymbol{\gamma} \in [0, 2\pi]^N$ .
6:    $\mathbf{x} \leftarrow [\mathbf{x} \cdot \sqrt{1 - 1/N}, \sqrt{1/N}]$  ▷ add reference path
7:    $\boldsymbol{\theta}_X, \boldsymbol{\phi}_X = \text{VEC2PHASE}(\mathbf{x})$  ▷ off-chip
8:    $i \leftarrow 1$ 
9:    $\mathbf{p} \leftarrow \mathbf{0}$ 
10:   $\mathbf{w} \leftarrow \text{SENDFORWARD}(\boldsymbol{\theta}_X, \boldsymbol{\phi}_X)$  ▷ on-chip
11:  for  $n \in [1, 2, \dots, N - 1]$  do ▷ on-chip
12:    for  $m \in [1, 2, \dots, N - m]$  do
13:       $\begin{pmatrix} w_m \\ w_{m+1} \end{pmatrix} \leftarrow \tilde{T}_2(\theta_i, \phi_i) \begin{pmatrix} w_m \\ w_{m+1} \end{pmatrix}$  ▷ forward prop
14:      measure  $p_{\theta_i}, p_{\phi_i}$  ▷ detect phase shift powers
15:       $i \leftarrow i + 1$ 
16:    end for
17:  end for
18:   $\boldsymbol{\theta}_Y, \boldsymbol{\phi}_Y = \text{READFORWARD}(\mathbf{w})$  ▷ self-configuration
19:   $\mathbf{y} \leftarrow \text{PHASE2VEC}(\boldsymbol{\theta}_Y, \boldsymbol{\phi}_Y) \cdot e^{i\boldsymbol{\gamma}}$  ▷ off-chip
20:   $\mathbf{y} \leftarrow \mathbf{y} \cdot N / \sqrt{1 - 1/N}$  ▷ remove reference
21:  return  $\mathbf{y}, \mathbf{p}$ 
22: end function

```

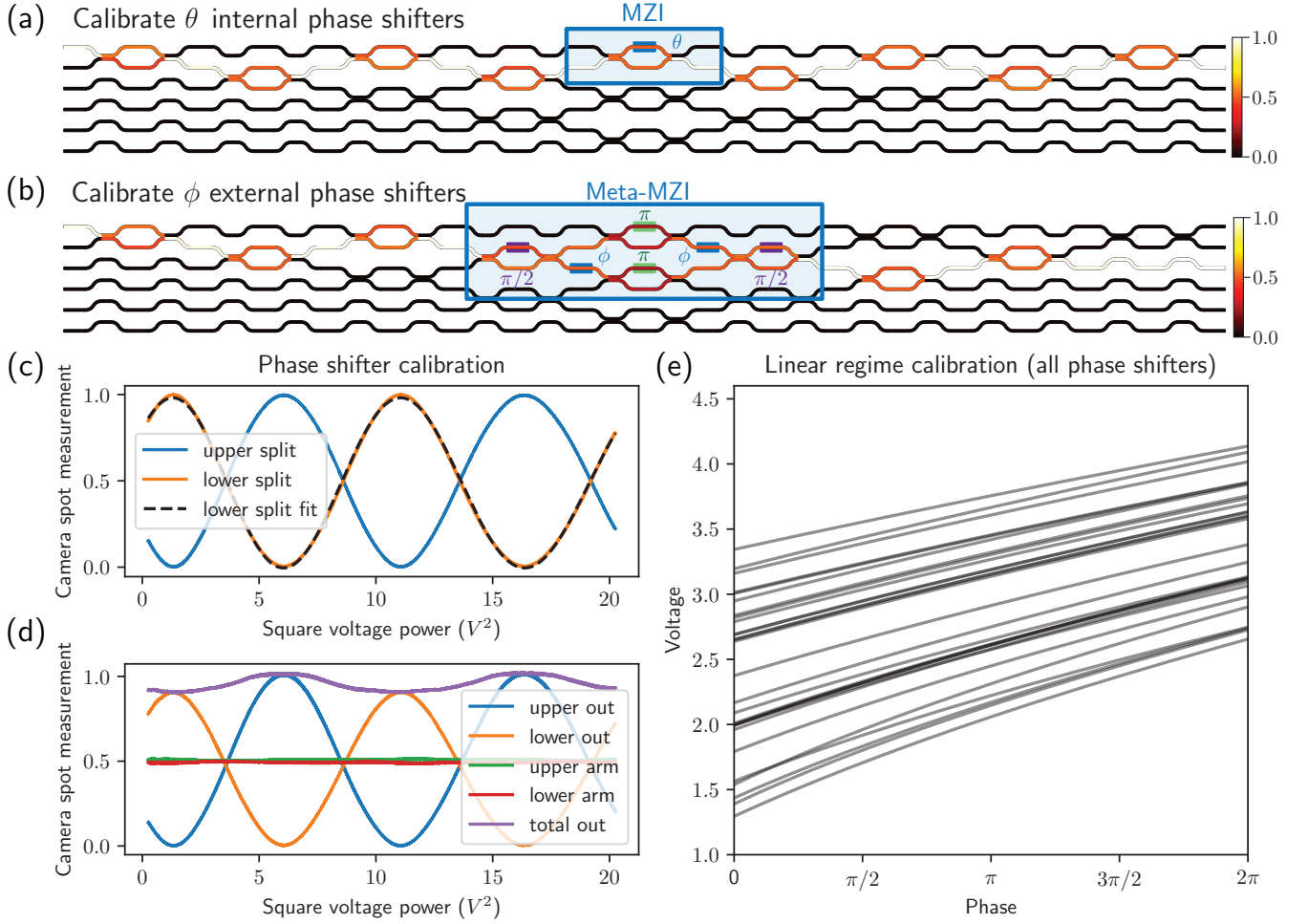


FIG. 8. (a) Calibration of θ internal phase shifts using lightwires leading to MZIs. (b) Calibration of ϕ phase shifts using lightwires leading to meta-MZI structures created out of four neighboring MZIs. (c) Phase shifter calibration protocol shows excellent fit (d) Raw camera spot measurement shows that different grating taps have different coupling efficiencies (a source of error in gradient measurements). (e) Linear regime of the calibration curve shows the range of voltages that need to be applied to our phase shifters to ensure that a full $[0, 2\pi]$ range can be achieved.

Algorithm 3 BACKWARD STEP

```

1: function MESHBACKWARD( $\mathbf{y}; \boldsymbol{\theta}, \boldsymbol{\phi}, \boldsymbol{\gamma}$ )
2:   require  $\mathbf{y} \in \mathbb{C}^N, \|\mathbf{y}\| = 1$ .
3:   require  $\boldsymbol{\theta} \in [0, \pi]^{N(N-1)/2}$ .
4:   require  $\boldsymbol{\phi} \in [0, 2\pi]^{N(N-1)/2}$ .
5:   require  $\boldsymbol{\gamma} \in [0, 2\pi]^N$ .
6:    $\mathbf{y} \leftarrow [\mathbf{y} \cdot \sqrt{1-1/N}, \sqrt{1/N}]$   $\triangleright$  add reference path
7:    $\boldsymbol{\theta}_Y, \boldsymbol{\phi}_Y = \text{VEC2PHASE}(\mathbf{y}^* \cdot e^{i\boldsymbol{\gamma}})$   $\triangleright$  off-chip
8:    $i \leftarrow N(N-1)/2$ 
9:    $\mathbf{p} \leftarrow \mathbf{0}$ 
10:   $\mathbf{w} \leftarrow \text{SENDERBACKWARD}(\boldsymbol{\theta}_Y, \boldsymbol{\phi}_Y)$   $\triangleright$  on-chip
11:  for  $n \in [1, 2, \dots, N-1]$  do  $\triangleright$  on-chip
12:    for  $m \in [1, \dots, m]$  do
13:       $\begin{pmatrix} w_m \\ w_{m+1} \end{pmatrix} \leftarrow \tilde{T}_2(\theta_i, \phi_i)^T \begin{pmatrix} w_m \\ w_{m+1} \end{pmatrix}$   $\triangleright$  back prop
14:      measure  $p_{\theta_i}, p_{\phi_i}$   $\triangleright$  detect phase shift powers
15:       $i \leftarrow i - 1$ 
16:    end for
17:  end for
18:   $\boldsymbol{\theta}_X, \boldsymbol{\phi}_X = \text{READBACKWARD}(\mathbf{w})$   $\triangleright$  self-configuration
19:   $\mathbf{x} \leftarrow \text{PHASE2VEC}(\boldsymbol{\theta}_X, \boldsymbol{\phi}_X)$   $\triangleright$  off-chip
20:   $\mathbf{x} \leftarrow \mathbf{x}_{:N} / \sqrt{1-1/N}$   $\triangleright$  remove reference
21:  return  $\mathbf{x}, \mathbf{p}$ 
22: end function

```

Algorithm 4 IN SITU BACKPROPAGATION

```

1: function INSITUGRADIENT( $\mathbf{x}, \mathbf{z}, \ell$ )
2:   require  $\mathbf{x} \in \mathbb{C}^N$ .
3:   require  $\mathbf{z} \in \mathbb{R}^N$ .
4:    $\boldsymbol{\theta}, \boldsymbol{\phi}, \boldsymbol{\gamma} \leftarrow \text{PHASES}(U^{(\ell)})$   $\triangleright$  phases of mesh  $\ell$ 
5:    $P \leftarrow \|\mathbf{x}\|^2$   $\triangleright$  input power scaling
6:    $\mathbf{y}, \mathbf{p} \leftarrow \text{MESHFORWARD}(\mathbf{x}/\sqrt{P}, \boldsymbol{\theta}, \boldsymbol{\phi}, \boldsymbol{\gamma})$ 
7:   if  $\ell = L$  then  $\triangleright$  End of the neural net
8:      $\mathbf{y}_{\text{aj}} = \partial c(\hat{\mathbf{z}}, \mathbf{z}) / \partial \mathbf{z}$   $\triangleright$  or  $\partial \mathcal{L} / \partial \mathbf{z}|_{\mathbf{x}, \mathbf{z}}$ 
9:      $\mathbf{g}_{\text{tot}} = \emptyset$   $\triangleright$  Empty gradient set
10:  else
11:     $\mathbf{x}_{\text{aj}}, \mathbf{g}_{\text{tot}} \leftarrow \text{INSITUGRADIENT}(f_\ell(\mathbf{y}), \mathbf{z}, \ell + 1)$ 
12:     $f_{\text{vjp}}^{(\ell)} \leftarrow \text{VJP}(f^{(\ell)})$   $\triangleright$  Autodiff, JAX/Haikou
13:     $\mathbf{y}_{\text{aj}} \leftarrow f_{\text{vjp}}^{(\ell)}(\mathbf{y}, \mathbf{x}_{\text{aj}})$ 
14:  end if
15:   $P_{\text{aj}} \leftarrow \|\mathbf{y}_{\text{aj}}\|^2$   $\triangleright$  adjoint power scaling
16:   $\mathbf{x}_{\text{aj}}, \mathbf{p}^{\text{aj}} \leftarrow \text{MESHBACKWARD}(\mathbf{y}_{\text{aj}}/\sqrt{P_{\text{aj}}}, \boldsymbol{\theta}, \boldsymbol{\phi}, \boldsymbol{\gamma})$ 
17:   $\dots, \mathbf{p}^{\text{sum}} \leftarrow \text{MESHFORWARD}(\mathbf{x} - i\mathbf{x}_{\text{aj}}^*, \boldsymbol{\theta}, \boldsymbol{\phi}, \boldsymbol{\gamma})$ 
18:   $\mathbf{g} \leftarrow \mathbf{p}^{\text{sum}} - \mathbf{p}^{\text{aj}} - \mathbf{p}$   $\triangleright$  analog or digital optical VJP
19:   $\mathbf{g} \leftarrow \mathbf{g} \cdot \sqrt{PP_{\text{aj}}/2}$   $\triangleright$  scaling factor
20:  return  $\mathbf{x}_{\text{aj}}, [\mathbf{g}, \mathbf{g}_{\text{tot}}]$   $\triangleright$  Append new gradients  $\mathbf{g}$  to  $\mathbf{g}_{\text{tot}}$ 
21: end function

```

Algorithm 5 IN SITU BACKPROPAGATION TRAINING

```

1: function INSITU MINIBATCH TRAIN( $X, Z, B$ )
2:   require  $X \in \mathbb{C}^{N_{\text{train}} \times N}$ .
3:   require  $Z \in \mathbb{R}^{N_{\text{train}} \times N_{\text{label}}}$ .
4:    $\mathbf{h} \leftarrow \mathbf{0}$  ▷ tracks gradient history
5:   for  $t \in [1, 2, \dots, T]$  do ▷ on-chip
6:     randomly sample  $X_t, Z_t$  from  $X, Z$ .
7:     require  $X_t \in \mathbb{C}^{B \times N}$ .
8:     require  $Z_t \in \mathbb{R}^{B \times N_{\text{label}}}$ .
9:     for  $\mathbf{x}_b, \mathbf{z}_b \in X_t, Z_t$  do
10:       $\dots, \mathbf{g}_b \leftarrow \text{INSITU GRADIENT}(\mathbf{x}_b, \mathbf{z}_b, \ell + 1)$ 
11:     end for
12:      $\mathbf{g}_t \leftarrow \sum_{b=1}^B \mathbf{g}_b / B$  ▷ minibatch average
13:      $\delta \boldsymbol{\eta}, \mathbf{h} \leftarrow \text{OPTIMIZER}(\mathbf{g}_t, \mathbf{h})$ 
14:   end for
15: end function

```

Appendix F: Pseudocode

In this section, we specifically provide some pseudocode required to implement various algorithms required for *in situ* backpropagation on our triangular mesh platform. Note that these approaches can be implemented on any matrix unit provided that the vector units can be used to generate any input fields. For output field generation, one can self-configure for backward and forward measurements on the existing vector units (Fig. 7(f)) or use a homodyne vector unit for measurement (Fig. 7(g)).

Our recursively defined algorithm for backpropagation on photonic meshes using the call $\mathbf{g} = \text{INSITU GRADIENT}(\mathbf{x}, \mathbf{z}, 1)$ where \mathbf{g} here represents gradients taken over all $\boldsymbol{\eta}$ in the network, as defined in Alg. 4, is based on Algs. 1, 2, 3 for generator/analyzer operation and the forward/backward steps for backpropagation. Note that some of the procedures such as READBACKWARD, SENDBACKWARD, READFORWARD, SENDFORWARD, PHASES do not have pseudocode, but these are explained in our Methods section and in Refs. 20, 22.

As previously discussed (Methods), the VJP (or vector Jacobian product) function is often used in neural networks and autodifferentiation frameworks (e.g., JAX) to automatically carry out chain rule steps used in measuring gradients. As defined in Alg. 4, a VJP calculation based on nonlinearity derivatives is performed in the digital domain since the nonlinearity itself is also performed in the digital domain. We have already defined VJP in the context of optical backpropagation (“optical VJP”) in the Methods section in terms of physical measurement; in general, nonlinear VJPs are more straightforward to compute digitally. Computing nonlinear VJPs does not offer much benefit in the optical domain for our purpose (energy efficient computation) since the energy to define inputs and outputs is already $O(N)$ in the digital-analog conversion which is also the complexity of an elementwise digital nonlinearity.

Finally, now that we have defined all of the gradient measurement pseudocode, we are ready to define the final

training protocol, which we use throughout this paper to achieve photonic *in situ* training. We define the full training set of N_{train} training examples as a $N_{\text{train}} \times N$ data matrix X and associated label set $N_{\text{train}} \times N_{\text{label}}$ Z :

Note that there are two nontrivial implementations in Alg. 5: the Adam optimizer [25] and minibatch training protocols. In practice, we leverage autodifferentiation packages to implement much of this needed functionality (e.g., we use JAX’s optim package for the Adam optimizer). We choose a minibatch size of 1 implementing a purely “stochastic” update which does not average over many training examples. This helps to avoid errors in the gradient which as we have found can accumulate over a large batch of training examples. This further underscores the importance of reducing gradient error to enable minibatch training.

Additionally, further investigation is warranted to explore *analog* adaptive update schemes that store previous gradients in nonvolatile memory. This would be important in cases where a purely analog update is required; otherwise a potentially more energy-consuming digital subtraction update would be needed to compute the history aggregation vector \mathbf{h} at each step of the optimization.

Appendix G: Analog update**1. Equivalence of digital and analog update**

Here, we prove the equivalence of $d_\eta(0)$ (our new analog measurement proposal) and the numerically evaluated gradient $\mathcal{I}(x_\eta x_{\eta, \text{aj}})$ which has been shown to be equivalent to the digital subtraction update [5].

The idea is to input a varying sum signal $\mathbf{x} - i\mathbf{x}_{\text{aj}}^* e^{i\zeta}$ and analyze the varying or AC component $d_\eta(\zeta)$ of the power $p_\eta(\zeta)$ measured at phase shifter η which has the fields x_η when sending \mathbf{x} alone and $i\mathbf{x}_{\text{aj}, \eta}^* e^{i\zeta}$ when sending \mathbf{x}_{aj}^* alone:

$$\begin{aligned}
 p_\eta(\zeta) &= |x_\eta|^2 + |i\mathbf{x}_{\text{aj}, \eta}^* e^{i\zeta}|^2 - 2\mathcal{R}(ix_\eta \mathbf{x}_{\text{aj}, \eta}^* e^{i\zeta}) \\
 d_\eta(\zeta) &= -2\mathcal{R}(ix_\eta \mathbf{x}_{\text{aj}, \eta}^* e^{i\zeta}) = 2\mathcal{I}(x_\eta \mathbf{x}_{\text{aj}, \eta} e^{i\zeta}),
 \end{aligned} \tag{G1}$$

which is equivalent to the gradient iff $\zeta = 0$.

2. Analog update protocol

Now that we have shown the equivalence of the digital and analog updates, we discuss the physical implementation of the analog gradient update implementation in hybrid photonic neural networks. As discussed in the main text, the analog signal processing to implement the gradient updater involves (1) a high pass filter and (2) a gated integrator implemented using a summing amplifier feedback (with gate width specified in the original signal synchronized to $\zeta(t) = 2\pi n$). The output of the

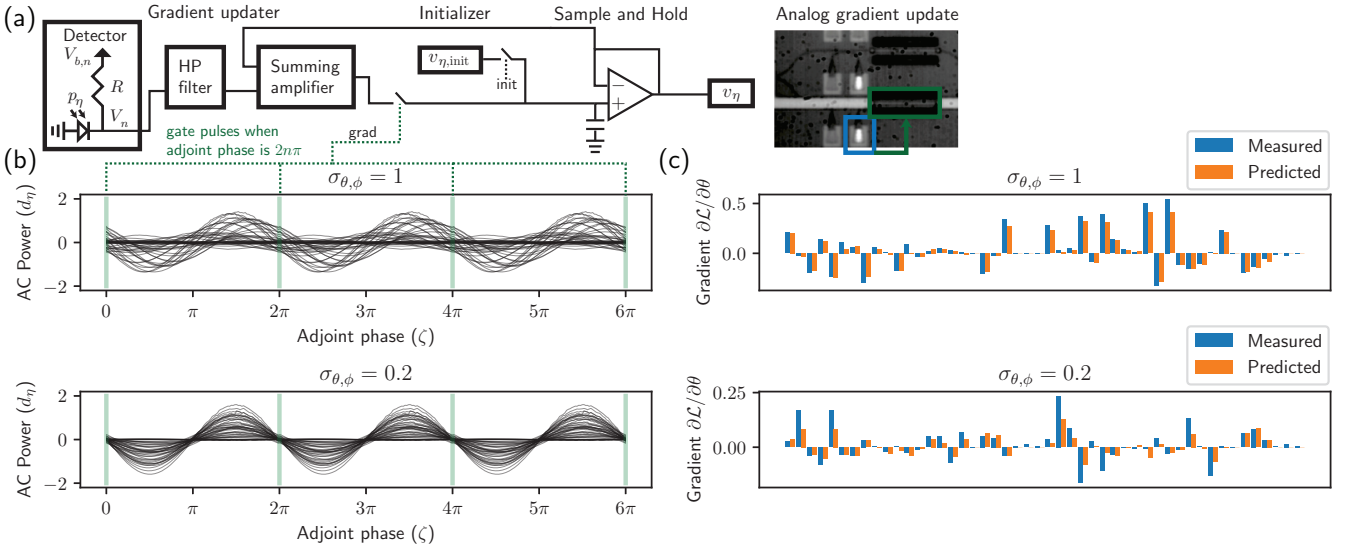


FIG. 9. (a) Our conceptual analog gradient update flow for updating phase shifter η based on power signal $p_\eta(\zeta)$, which varies according to adjoint phase. An integrated detector is connected via a “gradient updater” circuit consisting of a high-pass (HP) filter and summing amplifier to a sample-and-hold circuit scheme with an initializer for direct setting of voltage. (b-c) Standard deviation in the phase shift error is used to specify either far from or close to convergence ($\sigma_{\theta,\phi} = 1, 0.2$ respectively). (b) The AC power signal d_η versus the adjoint phase is experimentally measured on our chip across all relevant grating tap monitors, showing a decrease in gradient magnitude and more “in-phase” behavior near convergence. (c) As the distance to convergence decreases, there is more error in the gradient computation as expected, which is more explicitly shown in Fig. 2(f).

integrator is the gradient that can be directly applied as a control signal to the sample-and-hold phase shifter voltage. This is shown in more detail in Fig. 9(a).

Constant scaling factors required for gradient updates may be reflected in the analog signal processing, e.g. in the integrator step. Note that during *in situ* backpropagation, the forward- and backward- propagating optical signals in each of the photonic mesh accelerator chips are normalized to the same power. The computer stores the actual vector norms of the input and output vectors $\mathbf{x}, \mathbf{x}_{aj}$ as P, P_{aj} as defined also in Alg. 5. The sum vector $\mathbf{x} - i\mathbf{x}_{aj}^*$ is trickier to rescale. In this case, the input light is split equally into two input vectors implementing the normalized $\mathbf{x}, \mathbf{x}_{aj}^*$ and then interfered to yield the (lossy) vector sum $(\mathbf{x} - i\mathbf{x}_{aj}^*)/\sqrt{2}$ as shown in Fig. 2(b) of the main text. To recover the gradient, all that is needed is to multiply by the normalized factor $\sqrt{PP_{aj}}$. This can be applied as a uniform scaling factor to all gradient updaters used to determine the gradient in the analog domain. This is the only scaling factor that varies according to the training example sent through the device; all other scaling factors can be grouped in with the overall learning rate of the system.

Appendix H: Simulated error analysis

The processing capability of our proposed experimental prototype is limited by the size of the photonic circuit since the circuit size is just $N = 4$. For completeness, we

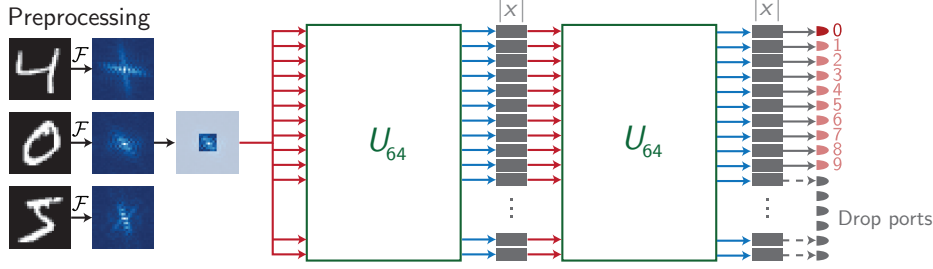
run simulations with a larger photonic circuit ($N = 64$) that uses a two-layer “triangular mesh” architecture with the same absolute value digital nonlinearity. We train this more expressive model on MNIST dataset [54] for hand-written digits recognition. As shown in Fig. 10(a), we follow the pre-processing procedure in [53] to convert the input 28×28 images into 64-dimensional complex vectors that are then input into the photonic circuit. We use the Adam optimizer [25] with learning rate $\alpha = 0.001$ to train the model following Alg. 5 for 30 epochs. We use the digit classification accuracy on the testset (with unseen data samples) as the metric for model performance.

To evaluate the robustness of the in-situ backpropagation process with respect to hardware errors, we add three types of errors in the simulations:

1. a_{error} , which represents the amplitude errors in field generation and analysis (READFORWARD, SENDFORWARD in Alg. 2 and Alg. 3).
2. p_{error} , which represents phase errors in field generation and analysis.
3. s_{noise} , which represents photon shot noise in optical power monitoring (line 14 in Alg. 2 and Alg. 3).

We calculate the shot noise signal-noise-ratio (SNR) snr_s at signal intensity $\sim 1/N$ since the input into optical neural network is normalized and each port has same average optical power. As shown in Fig. 10(b)-(d), with moderate level of noise (consistent with what is reported in current photonic circuits [55]), the model convergence

(a) Photonic neural network for the MNIST task



(b) Optical I/O phase error

(c) Optical I/O amplitude error

(d) Gradient updater noise error

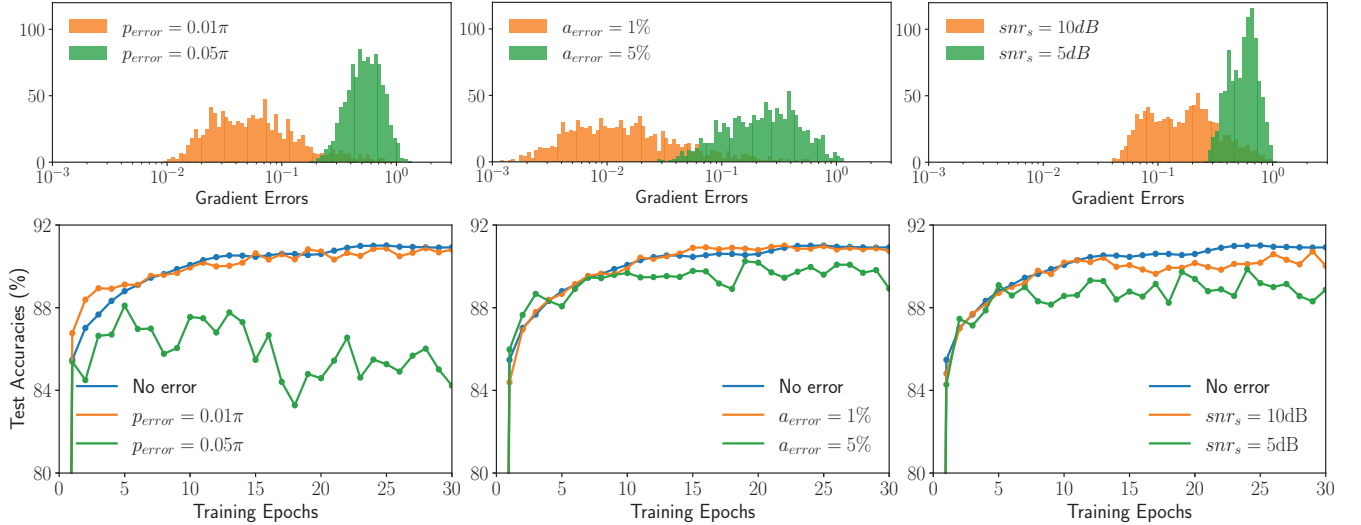


FIG. 10. (a) Two-layer triangular mesh optical neural network with $N = 64$ inputs. Images from MNIST datasets are pre-processed following the procedure in [32, 53]. (b)-(d) Normalized gradient errors (upper row) and testset accuracies (lower row) for models trained with in-situ backpropagation algorithm. Different types of hardware errors and noises are added to the training process; (b) field amplitude error a_{error} , (c) field error p_{error} , and (d) photon shot noise snr_s .

is minimally influenced, despite minor fluctuations. This demonstrates the robustness of in-situ backpropagation to noise and hardware errors, which are difficult to totally

eliminate in modern analog computing systems. All data and code to reproduce these results are provided in our data availability repository [34].

[1] Dario Amodei, Danny Hernandez, Girish Sastry, Jack Clark, Greg Brockman, and Ilya Sutskever, “AI and Compute,” (2018).
 [2] Seppo Linnainmaa, “Taylor expansion of the accumulated rounding error,” BIT **16**, 146–160 (1976).
 [3] David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams, “Learning representations by back-propagating errors,” Nature **323**, 533–536 (1986).
 [4] Alvaro A. Cruz-Cabrera, Mingtao Yang, Guoqi Cui, Elizabeth C. Behrman, James E. Steck, and Steven R. Skinner, “Reinforcement and backpropagation training for an optical neural network using self-lensing effects,” IEEE Transactions on Neural Networks **11**, 1450–1457 (2000).
 [5] Tyler W. Hughes, Momchil Minkov, Yu Shi, and Shanhui

Fan, “Training of photonic neural networks through in situ backpropagation and gradient measurement,” Optica **5**, 864 (2018).
 [6] Sunil Pai, Ian A.D. Williamson, Tyler W. Hughes, Momchil Minkov, Olav Solgaard, Shanhui Fan, and David A.B. Miller, “Parallel Programming of an Arbitrary Feedforward Photonic Network,” IEEE Journal of Selected Topics in Quantum Electronics **26** (2020), 10.1109/JSTQE.2020.2997849.
 [7] Yichen Shen, Nicholas C. Harris, Scott Skirlo, Mihika Prabhu, Tom Baehr-Jones, Michael Hochberg, Xin Sun, Shijie Zhao, Hugo Larochelle, Dirk Englund, and Marin Soljačić, “Deep learning with coherent nanophotonic circuits,” Nature Photonics **11**, 441–446 (2017).

- [8] David A. B. Miller, “Self-configuring universal linear optical component [Invited],” *Photonics Research* **1**, 1 (2013).
- [9] Andrea Annoni, Emanuele Guglielmi, Marco Carminati, Giorgio Ferrari, Marco Sampietro, David Ab Miller, Andrea Melloni, and Francesco Morichetti, “Unscrambling light - Automatically undoing strong mixing between modes,” *Light: Science and Applications* **6** (2017), 10.1038/lsa.2017.110.
- [10] Mitchell A. Nahmias, Thomas Ferreira De Lima, Alexander N. Tait, Hsuan Tung Peng, Bhavin J. Shastri, and Paul R. Prucnal, “Photonic Multiply-Accumulate Operations for Neural Networks,” *IEEE Journal of Selected Topics in Quantum Electronics* **26** (2020), 10.1109/JSTQE.2019.2941485.
- [11] Adolf Hurwitz, “über die Erzeugung der Invarianten durch Integration,” *Nachrichten von der Gesellschaft der Wissenschaften zu Göttingen, Mathematisch-Physikalische Klasse*, 71–72 (1897).
- [12] Michael Reck, Anton Zeilinger, Herbert J. Bernstein, and Philip Bertani, “Experimental realization of any discrete unitary operator,” *Physical Review Letters* **73**, 58–61 (1994).
- [13] Wim Bogaerts, Daniel Pérez, José Capmany, David A.B. Miller, Joyce Poon, Dirk Englund, Francesco Morichetti, and Andrea Melloni, “Programmable photonic circuits,” *Nature* **586**, 207–216 (2020).
- [14] Jacques Carolan, Christopher Harrold, Chris Sparrow, Enrique Martín-López, Nicholas J. Russell, Joshua W. Silverstone, Peter J. Shadbolt, Nobuyuki Matsuda, Manabu Oguma, Mikitaka Itoh, Graham D. Marshall, Mark G. Thompson, Jonathan C.F. Matthews, Toshikazu Hashimoto, Jeremy L. O’Brien, and Anthony Laing, “Universal linear optics,” *Science* (2015), 10.1126/science.aab3642.
- [15] Nicholas C. Harris, Jacques Carolan, Darius Bunandar, Mihika Prabhu, Michael Hochberg, Tom Baehr-Jones, Michael L. Fanto, A. Matthew Smith, Christopher C. Tison, Paul M. Alsing, and Dirk Englund, “Linear programmable nanophotonic processors,” *Optica* **5**, 1623 (2018).
- [16] “Lightmatter - The photonic (super)computer company.”
- [17] Logan G. Wright, Tatsuhiro Onodera, Martin M. Stein, Tianyu Wang, Darren T. Schachter, Zoey Hu, and Peter L. McMahon, “Deep physical neural networks trained with backpropagation,” *Nature* **601**, 549–555 (2022).
- [18] James Spall, Xianxin Guo, and A. I. Lvovsky, “Hybrid training of optical neural networks,” (2022).
- [19] H. Zhang, M. Gu, X. D. Jiang, J. Thompson, H. Cai, S. Paesani, R. Santagati, A. Laing, Y. Zhang, M. H. Yung, Y. Z. Shi, F. K. Muhammad, G. Q. Lo, X. S. Luo, B. Dong, D. L. Kwong, L. C. Kwek, and A. Q. Liu, “An optical neural chip for implementing complex-valued neural network,” *Nature Communications* **12**, 1–11 (2021).
- [20] David A. B. Miller, “Analyzing and generating multimode optical fields using self-configuring networks,” (2020).
- [21] Mihika Prabhu, Charles Roques-Carmes, Yichen Shen, Nicholas Harris, Li Jing, Jacques Carolan, Ryan Hamerly, Tom Baehr-Jones, Michael Hochberg, Vladimir Čeperić, John D. Joannopoulos, Dirk R. Englund, and Marin Soljačić, “Accelerating recurrent Ising machines in photonic integrated circuits,” *Optica* **7**, 551 (2020).
- [22] David A. B. Miller, “Perfect optics with imperfect components,” *Optica* **2**, 747 (2015).
- [23] David A. B. Miller, “Setting up meshes of interferometers – reversed local light interference method,” *Optics Express* **25**, 29233 (2017).
- [24] David A. B. Miller, “Establishing Optimal Wave Communication Channels Automatically,” *Journal of Lightwave Technology*, Vol. 31, Issue 24, pp. 3987-3994 **31**, 3987–3994 (2013).
- [25] Diederik P Kingma and Jimmy Lei Ba, “Adam: A Method for Stochastic Optimization,” *International Conference on Learning Representations* (2015).
- [26] Sunil Pai, Zhanghao Sun, and Taewon Park, “phox: Base repository for simulation and control of photonic devices,” (2022).
- [27] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake Vander{P}las, Skye Wanderman-{M}ilne, and Qiao Zhang, “{JAX}: composable transformations of {P}ython+{N}um{P}y programs,” (2022).
- [28] Tom Hennigan, Trevor Cai, Tamara Norman, and Igor Babuschkin, “{H}aikU: {S}onnet for {JAX},” (2020).
- [29] Hui Zhang, Jayne Thompson, Mile Gu, Xu Dong Jiang, Hong Cai, Patricia Yang Liu, Yuzhi Shi, Yi Zhang, Muhammad Faeyz Karim, Guo Qiang Lo, Xianshu Luo, Bin Dong, Leong Chuan Kwek, and Ai Qun Liu, “Efficient On-Chip Training of Optical Neural Networks Using Genetic Algorithm,” *ACS Photonics* **8**, 1662–1672 (2021).
- [30] Arild Nøkland, “Direct Feedback Alignment Provides Learning in Deep Neural Networks,” *Advances in Neural Information Processing Systems*, 1045–1053 (2016).
- [31] Matthew J. Filipovich, Zhimu Guo, Mohammed Al-Qadasi, Bicky A. Marquez, Hugh D. Morison, Volker J. Sorger, Paul R. Prucnal, Sudip Shekhar, and Bhavin J. Shastri, “Monolithic Silicon Photonic Architecture for Training Deep Neural Networks with Direct Feedback Alignment,” (2021).
- [32] Ian A. D. Williamson, Tyler W. Hughes, Momchil Minkov, Ben Bartlett, Sunil Pai, and Shanhui Fan, “Reprogrammable Electro-Optic Nonlinear Activation Functions for Optical Neural Networks,” *IEEE Journal of Selected Topics in Quantum Electronics* **26**, 1–12 (2020).
- [33] Martin Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg, Rajat Monga, Sherry Moore, Derek G. Murray, Benoit Steiner, Paul Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng, “TensorFlow: A System for Large-Scale Machine Learning,” in *Operating Systems Design and Implementation* (Savannah, GA, 2016) pp. 265–283.
- [34] Sunil Pai and Zhanghao Sun, “solgaard-lab/photonicbackprop: Data and code for the paper “Experimentally realized in situ backpropagation for deep learning in energy-efficient nanophotonic neural networks,”” (2022), 10.5281/ZENODO.6557414.
- [35] Sunil Pai, “simphox: Another inverse design library,” (2022).
- [36] Sunil Pai and Nathnael Abebe, “dphox: photonic layout and device design,” (2022).
- [37] David A. B. Miller, “Self-aligning universal beam coupler,” *Optics Express* **21**, 6360 (2013).
- [38] Leo Filippini and Baris Taskin, “A 900 MHz Charge Recovery Comparator with 40 fJ per Conversion,” in *Pro-*

- ceedings - IEEE International Symposium on Circuits and Systems*, Vol. 2018-May (Institute of Electrical and Electronics Engineers Inc., 2018).
- [39] Masaya Miyahara, Yusuke Asada, Daehwa Paik, and Akira Matsuzawa, "A low-noise self-calibrating dynamic comparator for high-speed ADCs," in *Proceedings of 2008 IEEE Asian Solid-State Circuits Conference, A-SSCC 2008* (2008) pp. 269–272.
- [40] Cheng Wang, Mian Zhang, Brian Stern, Michal Lipson, and Marko Lončar, "Nanophotonic lithium niobate electro-optic modulators," *Optics Express* **26**, 1547 (2018).
- [41] Nicholas C. Harris, Yangjin Ma, Jacob Mower, Tom Baehr-Jones, Dirk Englund, Michael Hochberg, and Christophe Galland, "Efficient, compact and low loss thermo-optic phase shifter in silicon," *Optics Express* **22**, 10487 (2014).
- [42] Carlos Errando-Herranz, Alain Yuji Takabayashi, Pierre Edinger, Hamed Sattari, Kristinn B. Gylfason, and Niels Quack, "MEMS for Photonic Integrated Circuits," *IEEE Journal of Selected Topics in Quantum Electronics* **26** (2020), 10.1109/JSTQE.2019.2943384.
- [43] M. Wuttig, H. Bhaskaran, and T. Taubner, "Phase-change materials for non-volatile photonic applications," (2017).
- [44] Pierre Edinger, Carlos Errando-Herranz, and Kristinn Gylfason, "Low-loss MEMS phase shifter for large scale reconfigurable silicon photonics," in *The 32nd IEEE International Conference on Micro Electro Mechanical Systems* (2019).
- [45] Pierre Edinger, Carlos Errando-Herranz, Alain Yuji Takabayashi, Hamed Sattari, Niels Quack, Peter Verheyen, Wim Bogaerts, and Kristinn B Gylfason, *Conference on Lasers and Electro-Optics (CLEO 2020)*, Tech. Rep. (2020).
- [46] Amir Dembo and Thomas Kailath, "Model-Free Distributed Learning," *IEEE Transactions on Neural Networks* **1**, 58–70 (1990).
- [47] Gert Cauwenberghs, "A Fast Stochastic Error-Descent Algorithm for Supervised Learning and Optimization," in *Advances in Neural Information Processing Systems 5* (1992).
- [48] J Alspector, R Meir, B Yuhas, A Jayakumar, and D Lippe, "A Parallel Gradient Descent Method for Learning in Analog VLSI Neural Networks," in *Advances in Neural Information Processing Systems* (1992) pp. 836–844.
- [49] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei, "ImageNet: A large-scale hierarchical image database," (Institute of Electrical and Electronics Engineers (IEEE), 2010) pp. 248–255.
- [50] William R. Clements, Peter C. Humphreys, Benjamin J. Metcalf, W. Steven Kolthammer, and Ian A. Walmsley, "An Optimal Design for Universal Multiport Interferometers," *Optica* , 1–8 (2016).
- [51] Mikko Möttönen, Juha J. Vartiainen, Ville Bergholm, and Martti M. Salomaa, "Quantum Circuits for General Multiqubit Gates," *Physical Review Letters* **93**, 130502 (2004).
- [52] Jacob Mower, Nicholas C. Harris, Gregory R. Steinbrecher, Yoav Lahini, and Dirk Englund, "High-fidelity quantum state evolution in imperfect photonic integrated circuits," *Physical Review A* **92**, 032322 (2015).
- [53] Sunil Pai, Ian A. D. Williamson, Tyler W. Hughes, Momchil Minkov, Olav Solgaard, Shanhui Fan, and David A. B. Miller, "Parallel fault-tolerant programming of an arbitrary feedforward photonic network," arXiv preprint (2019).
- [54] Li Deng, "The MNIST database of handwritten digit images for machine learning research," *IEEE Signal Processing Magazine* **29**, 141–142 (2012).
- [55] Saumil Bandyopadhyay, Ryan Hamerly, and Dirk Englund, "Hardware error correction for programmable photonics," *Optica* **8**, 1247 (2021).