

Quantum-noise-limited optical neural networks operating at a few quanta per activation

Shi-Yuan Ma (✉ sm2725@cornell.edu)

Cornell University <https://orcid.org/0000-0001-8299-6742>

Tianyu Wang

Cornell University <https://orcid.org/0000-0002-6087-6376>

Jérémie Laydevant

Cornell University

Logan Wright

Cornell University <https://orcid.org/0000-0001-7696-1260>

Peter McMahon

Cornell University <https://orcid.org/0000-0002-1177-9887>

Article

Keywords:

Posted Date: October 26th, 2023

DOI: <https://doi.org/10.21203/rs.3.rs-3318262/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Additional Declarations: There is **NO** Competing Interest.

Quantum-noise-limited optical neural networks operating at a few quanta per activation

Shi-Yuan Ma^{1*}, Tianyu Wang¹, Jérémie Laydevant^{1,2}, Logan G. Wright^{1,3,†},
Peter L. McMahon^{1,4*}

¹School of Applied and Engineering Physics, Cornell University, Ithaca, 14853, NY, USA.

²USRA Research Institute for Advanced Computer Science, Mountain View, 94035, CA, USA.

³NTT Physics and Informatics Laboratories, NTT Research, Inc., Sunnyvale, 94085, CA, USA.

⁴Kavli Institute at Cornell for Nanoscale Science, Cornell University, Ithaca, 14853, NY, USA.

[†]Present address: Department of Applied Physics, Yale University, New Haven, 06511, CT, USA.

*Corresponding author(s). E-mail(s): sm2725@cornell.edu; pmcmahon@cornell.edu;

A practical limit to energy efficiency in computation is ultimately from noise, with quantum noise [1] as the fundamental floor. Analog physical neural networks [2], which hold promise for improved energy efficiency and speed compared to digital electronic neural networks, are nevertheless typically operated in a relatively high-power regime so that the signal-to-noise ratio (SNR) is large (>10). We study optical neural networks [3] operated in the limit where all layers except the last use only a single photon to cause a neuron activation. In this regime, activations are dominated by quantum noise from the fundamentally probabilistic nature of single-photon detection. We show that it is possible to perform accurate machine-learning inference in spite of the extremely high noise (signal-to-noise ratio ~ 1). We experimentally demonstrated MNIST handwritten-digit classification with a test accuracy of 98% using an optical neural network with a hidden layer operating in the single-photon regime; the optical energy used to perform the classification corresponds to 0.008 photons per multiply-accumulate (MAC) operation, which is equivalent to 0.003 attojoules of optical energy per MAC. Our experiment also used $>40\times$ fewer photons per inference than previous state-of-the-art low-optical-energy demonstrations [4, 5] to achieve the same accuracy of $>90\%$. Our training approach, which directly models the system’s stochastic behavior, might also prove useful with non-optical ultra-low-power hardware.

The development and widespread use of very large neural networks for artificial intelligence [6, 7] has motivated the exploration of alternative computing paradigms—including analog processing—in the hope of improving both energy efficiency and speed [2, 8]. Photonic implementations of neural networks using analog optical systems have experienced a resurgence of interest over the past several years [3–5, 9–16]. However, analog processors—including those constructed using optics—inevitably have noise and typically also suffer from imperfect calibration and drift. These imperfections can result in degraded accuracy for neural-network inference performed using them [9, 17–19]. To mitigate the impact of noise, noise-aware training schemes have been developed [20–27]. These schemes treat the noise as a relatively small perturbation to an otherwise deterministic computation, either by explicitly modeling the noise as the addition of random variables to the processor’s output or by modeling the processor as having finite bit precision. Recent demonstrations of ultra-low optical energy usage in optical neural networks (ONNs) [4, 5] were in this regime of noise as a small perturbation and used hundreds to thousands of photons to represent the average neuron pre-activation signal prior to photodetection. In Ref. [4], we reported

41 achieving 90% accuracy on MNIST handwritten-digit classification using slightly less than 1 photon per
42 scalar weight multiplication (i.e., per MAC)—which is already counterintuitively small—and one might
43 be tempted to think that it’s not possible to push the number of photons per MAC much lower while
44 preserving accuracy. More typically, millions of photons per activation are used [5, 13, 14, 28]. In this
45 paper we address the following question: what happens if we use such weak optical signals in a ONN that
46 each photodetector in a neural-network layer receives at most just one, or perhaps two or three, photons?

47 Physical systems are subject to various sources of noise. While some noise can be reduced through
48 improvements to the hardware, some noise is fundamentally unavoidable, especially when the system is
49 operated with very little power—which is an engineering goal for neural-network processors. Shot noise
50 is a fundamental noise that arises from the quantized, i.e., discrete, nature of information carriers: the
51 discreteness of energy in the case of photons in optics, and of discreteness of charge in the case of electrons
52 in electronics [1]. A shot-noise-limited measurement of a signal encoded with an average of N_p photons
53 (quanta) will have an SNR that scales as $\sqrt{N_p}$ [29].¹ To achieve a suitably high SNR, ONNs typically use
54 a large number of quanta for each detected signal. In situations where the optical signal is limited to just
55 a few photons, photodetectors measure and can count individual quanta. Single-photon detectors (SPDs)
56 are highly sensitive detectors that—in the typical *click detector* setting—report, with high fidelity, the
57 absence of a photon (*no click*) or presence of one or more photons (*click*) during a given measurement
58 period [31]. In the quantum-noise-dominated regime of an optical signal with an average photon number
59 of about 1 impinging on an SPD, the measurement outcome will be highly stochastic, resulting in a
60 very low SNR (of about 1).² Conventional noise-aware-training algorithms are not able to achieve high
61 accuracy with this level of noise. Is it possible to operate ONNs in this very stochastic regime and still
62 achieve high accuracy in deterministic classification tasks? The answer is *yes*, and in this work we will
63 show how.

64 The key idea in our work is that when ONNs are operated in the approximately-1-photon-per-neuron-
65 activation regime and the detectors are SPDs, it is natural to consider the neurons as binary stochastic
66 neurons: the output of an SPD is binary (*click* or *no click*) and fundamentally stochastic. Instead of
67 trying to train the ONN as a deterministic neural network that has very poor numerical precision, one
68 can instead train it as a binary stochastic neural network, adapting some of the methods from the last
69 decade of machine-learning research on stochastic neural networks [32–36] and using a physics-based
70 model of the stochastic single-photon detection (SPD) process during training. We call this *physics-aware*
71 *stochastic training*.

¹The *shot-noise limit*, which is sometimes also referred to as the *standard quantum limit* [30], can be evaded if, instead of encoding the signal in a thermal or coherent state of light, a quantum state—such as an intensity-squeezed state or a Fock state—is used. In this paper we consider only the case of *classical* states of light for which shot noise is present and the shot-noise limit applies.

²Again, this is under the assumption that the optical signal is encoded in an optical state that is subject to the shot-noise limit—which is the case for classical states of light.

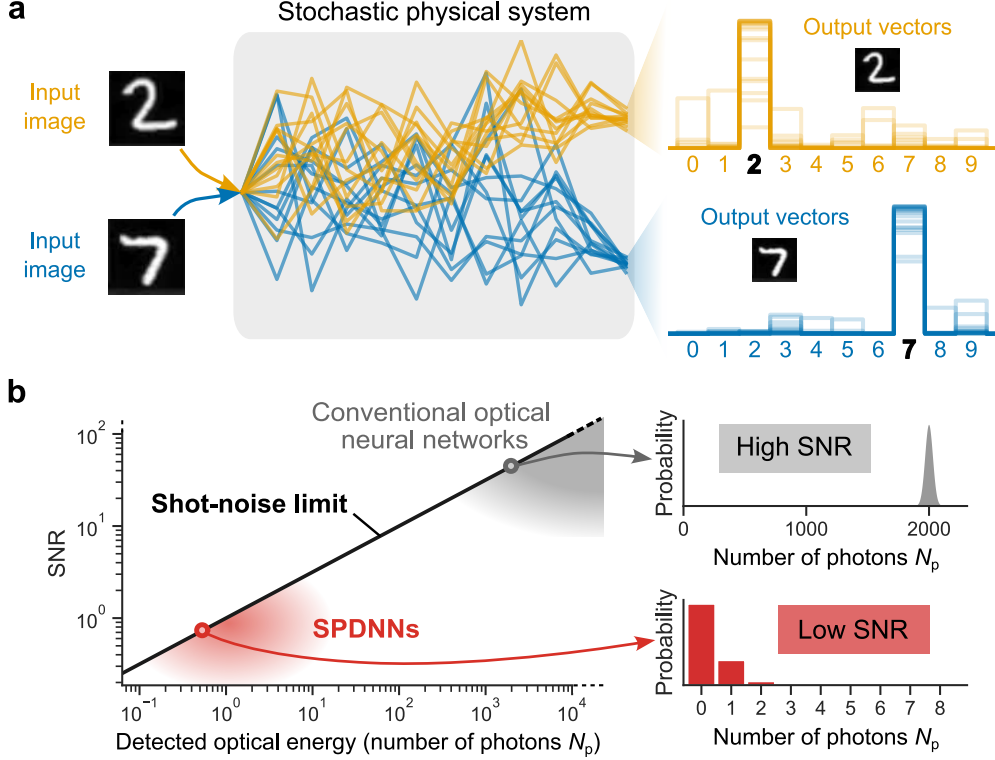


Figure 1. Deterministic inference using noisy neural-network hardware. **a**, The concept of a stochastic physical neural network performing a classification task. Given a particular input image to classify, repetitions exhibits variation (represented by different traces of the same color), but the class is predicted nearly deterministically. **b**, The single-to-noise ratio (SNR) of single-photon-detection neural networks (SPDNNs) compared to conventional optical neural networks (ONNs). Conventional ONNs operate with high photon budgets (SNR $\gg 1$) to obtain reliable results, whereas SPDNNs operate with low photon budgets—of up to just a few detected photons per shot (SNR ~ 1). The relation between the detected optical energy (in number of photons N_p) and SNR is $\text{SNR} = \sqrt{N_p}$, which is known as the shot-noise limit.

72 We experimentally implemented a stochastic ONN using as a building block an optical matrix-vector
73 multiplier [4] modified to have SPDs at its output: we call this a *single-photon-detection neural network*
74 (SPDNN). We present results showing that high classification accuracy can be achieved even when the
75 number of photons per neuron activation is approximately 1, and even without averaging over multiple
76 shots. We also studied in simulation how larger, more sophisticated stochastic ONNs could be constructed
77 and what their performance on CIFAR-10 image classification would be. While the proof-of-concept
78 experiments we report are based on a specific spatially multiplexed, free-space ONN, our approach doesn't
79 rely on details of this architecture and could be adapted for many other types of ONN, including those
80 based on diffractive optics [10, 14, 37], Mach-Zehnder interferometer (MZI) meshes [9, 38, 39], and other
81 on-chip or hybrid approaches to matrix-vector multiplication [5, 12, 13, 40].

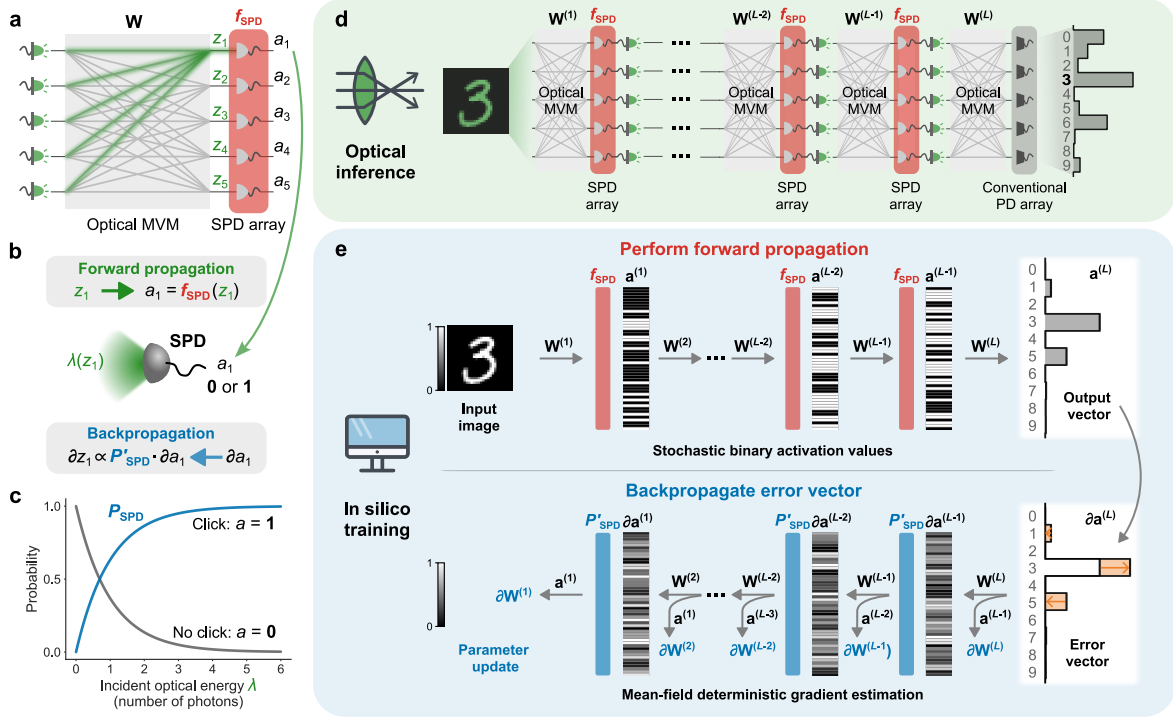


Figure 2. Single-photon-detection neural networks (SPDNNs): *physics-aware stochastic training and inference.* **a**, A single layer of an SPDNN, comprising an optical matrix-vector multiplier (optical MVM, in grey) and single-photon detectors (SPDs; in red), which perform stochastic nonlinear activations. Each output neuron’s value is computed by the physical system as $a_i = f_{\text{SPD}}(z_i)$, where z_i is the weighted sum (shown in green) of the input neurons to the i th output neuron computed as part of the optical MVM, and a_i is the stochastic binary output from a single-photon detector. **b**, Forward and backward propagation through the SPD activation function. The optical energy (λ) incident on an SPD is a function of z_i that depends on the encoding scheme used. Forward propagation uses the stochastic binary activation function f_{SPD} , while backpropagation involves the mean-field function of the probability P_{SPD} . **c**, Probability of an SPD detecting a click (output $a = 1$) or not (output $a = 0$), as a function of the incident light energy λ . **d**, Optical inference using an SPDNN with L layers. The activation values from the SPD array of each layer are passed to light emitters for the optical MVM of the next layer. The last layer uses a conventional photodetector (PD) array instead of an SPD array. **e**, *In silico* training of an SPDNN with L layers. Each forward propagation is stochastic, and during backpropagation, the error vector is passed to the hidden layers using the mean-field probability function P_{SPD} instead of the stochastic activation function f_{SPD} . In this figure, ∂x is shorthand for $\partial C/\partial x$, where C is the cost function.

82 Single-photon-detection neural networks: optical neural 83 networks with stochastic activation from single-photon detection

84 We consider ONNs in which one or more layers are each constructed from an optical matrix-vector
85 multiplier followed by an array of SPDs (Figure 2a–c), and in which the optical powers used are sufficiently
86 low that in each execution of the layer, each SPD has at most only a few photons impinging on it, leading
87 to stochastic measurement outcomes of *no click* or *click*.

88 In our setting, we aim to perform *inference* using the SPDNN—with its implementation in physical
 89 hardware—(Figure 2d) and to perform *training* of the SPDNN *in silico* (Figure 2e). That is, training is
 90 performed entirely using standard digital electronic computing.³

91 Physics-aware stochastic training

92 To train an SPDNN, we perform gradient descent using backpropagation, which involves a forward pass,
 93 to compute the current error (or loss) of the network, and a backward pass, which is used to compute the
 94 gradient of the loss with respect to the network parameters; our procedure is inspired by backpropagation-
 95 based training of stochastic and binary neural networks [32, 35]. We model the forward pass (upper part
 96 of Figure 2e) through the network as a stochastic process that captures the key physics of SPD of optical
 97 signals having Poissonian photon statistics [29]: the measurement outcome of SPD is a binary random
 98 variable (*no click* or *click*) that is drawn from the Bernoulli distribution with a probability that depends
 99 on the mean photon number of the light impinging on the detector. However, during the backward pass
 100 (lower part of Figure 2e), we employ a deterministic mean-field estimator to compute the gradients. This
 101 approach avoids the stochasticity and binarization of the SPD process, which typically pose difficulties
 102 for gradient estimation.

103 We now give a brief technical description of our forward and backward passes for training; for full
 104 details see Methods and Supplementary Notes 1A and 2A. We denote the neuron pre-activations of the
 105 l th stochastic layer of an SPDNN as $\mathbf{z}^{(l)} = W^{(l)}\mathbf{a}^{(l-1)}$, where $\mathbf{a}^{(l-1)}$ is the activation vector from the
 106 previous layer ($\mathbf{a}^{(0)}$ denotes the input vector \mathbf{x} of the data to be classified). In the physical realization of
 107 an SPDNN, $\mathbf{z}^{(l)}$ is encoded optically (for example, in optical intensity) following an optical matrix-vector
 108 multiplier (optical MVM, which computes the product between the matrix $W^{(l)}$ and the vector $\mathbf{a}^{(l-1)}$) but
 109 before the light impinges on an array of SPDs. We model the action of an SPD with a stochastic activation
 110 function, f_{SPD} (Figure 2b; Eq. 1). The stochastic output of the l th layer is then $\mathbf{a}^{(l)} = f_{\text{SPD}}(\mathbf{z}^{(l)})$.

111 For an optical signal having mean photon number λ and that obeys Poissonian photon statistics, the
 112 probability of a *click* event by an SPD is $P_{\text{SPD}}(\lambda) = 1 - e^{-\lambda}$ (Figure 2c). We define the stochastic activation
 113 function f_{SPD} as follows:

$$f_{\text{SPD}}(z) := \begin{cases} 1 & \text{with probability } p = P_{\text{SPD}}(\lambda(z)), \\ 0 & \text{with probability } 1 - p, \end{cases} \quad (1)$$

114 where $\lambda(z)$ is a function mapping a single neuron’s pre-activation value to a mean photon number. For an
 115 incoherent optical setup where the information is directly encoded in intensity, $\lambda(z) = z$; for a coherent
 116 optical setup where the information is encoded in field amplitude and the SPD directly measures the

³It is not required that the training be done *in silico* for it to succeed but is just a choice we made in this work. *Hardware-in-the-loop* training, such as used in Ref. [23], is a natural alternative to purely *in silico* training that even can make training easier by relaxing the requirements on how accurate the *in silico* model of the physical hardware process needs to be.

117 intensity, $\lambda(z) = |z|^2$. In general, the form of $\lambda(z)$ is determined by the signal encoding used in the optical
 118 MVM, and the detection scheme following the MVM. We use f_{SPD} in modeling the stochastic behavior
 119 of an SPDNN layer in the forward pass. However, during the backward pass, we make a deterministic
 120 mean-field approximation of the network: instead of evaluating the stochastic function f_{SPD} , we evaluate
 121 $P_{\text{SPD}}(\lambda(z))$ when computing the activations of a layer: $\mathbf{a}^{(l)} = P_{\text{SPD}}(\lambda(\mathbf{z}^{(l)}))$ (Figure 2b). This is an
 122 adaptation of a standard machine-learning method for computing gradients of stochastic neural networks
 123 [32].

124 Inference

125 When performing inference (Figure 2d), we can run just a single shot of a stochastic layer or we can
 126 choose to take the average of multiple shots—trading greater energy and/or time usage for reduced
 127 stochasticity. For a single shot, a neuron activation takes on the value $a^{[1]} = a \in \{0, 1\}$; for K shots,
 128 $a^{[K]} = \frac{1}{K} \sum_{k=1}^K a_k \in \{0, 1/K, 2/K, \dots, 1\}$. In the limit of infinitely many shots, $K \rightarrow \infty$, the activation
 129 $a^{[\infty]}$ would converge to the expectation value, $a^{[\infty]} = \mathbb{E}[a] = P_{\text{SPD}}(\lambda(z))$. In this work we focus on the
 130 single-shot ($K = 1$) and few-shot $K \leq 5$ regime, since the high-shot $K \gg 100$ regime is very similar
 131 to the high-photon-count-per-shot regime that has already been studied in the ONN literature (e.g., in
 132 Ref. [4]). An important practical point is that averaging for $K > 1$ shots can be achieved by counting
 133 the clicks from each SPD, which is what we did in the experiments we report. We can think of K as a
 134 discrete integration time, so averaging need not involve any data reloading or sophisticated control.

135 MNIST handwritten-digit classification with a 136 single-photon-detection multilayer perceptron

137 We evaluated the performance—both in numerical simulations and in optical experiments—of SPDNNs
 138 on the MNIST handwritten-digit-classification benchmark task with a simple, $784 \rightarrow N \rightarrow 10$ multilayer
 139 perceptron (MLP) architecture (Figure 3a). The activation values in the hidden layer were computed
 140 by SPDs. The optical power was chosen so that the SNR of the SPD measurements was ~ 1 , falling in
 141 the low-SNR regime (Figure 1b). The output layer was implemented either with full numerical precision
 142 on a digital electronic computer, or optically with an integration time set so that the measured signal
 143 comprised enough photons that a high SNR (Figure 1b) was achieved, as in conventional ONNs. Our use
 144 of a full-precision output layer is consistent with other works on binary neural networks [35, 41, 42]. In
 145 a shallow neural network, executing the output layer at high SNR substantially limits the overall energy
 146 efficiency gains from using small photon budgets in earlier layers, but in larger models, the relatively
 147 high energy cost of a high-SNR output layer is amortized. Nevertheless, as we will see, even with just a

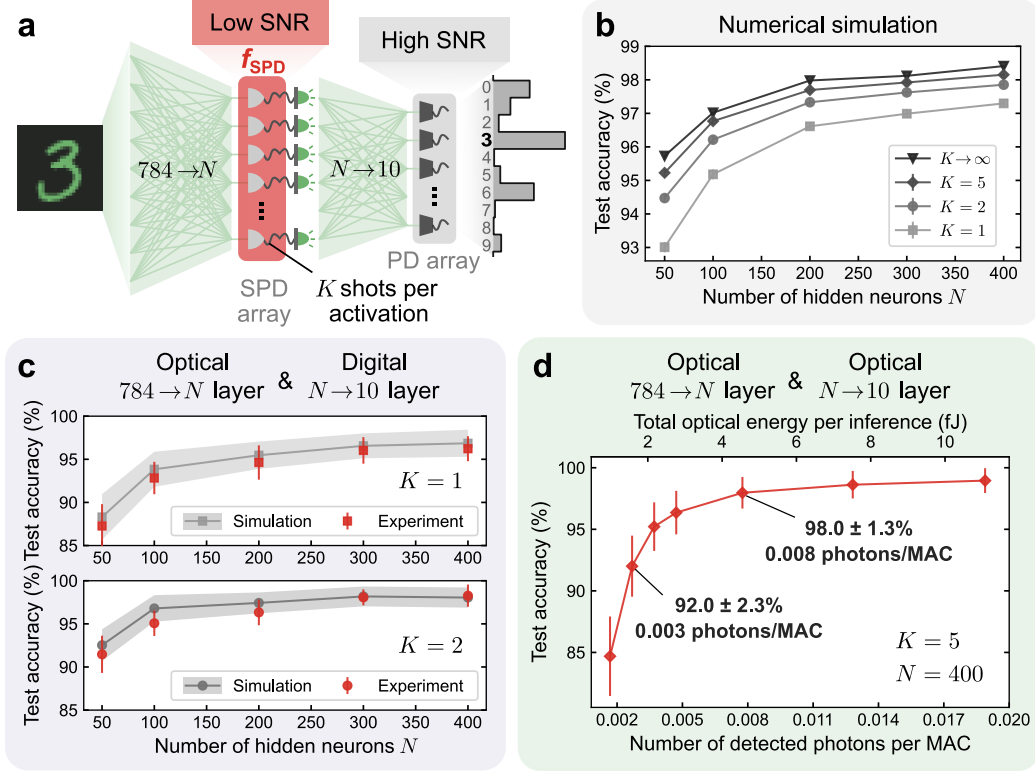


Figure 3. Performance of a single-photon-detection neural network (SPDNN) on MNIST handwritten-digit classification. **a**, An SPDNN realizing a multilayer perceptron (MLP) architecture of N neurons in the hidden layer. The hidden layer ($784 \rightarrow N$) was computed using an incoherent optical matrix-vector-multiplier (MVM) followed by a single-photon-detector (SPD) array. Each SPD realized a stochastic activation function for a single hidden-layer neuron. During a single inference, the hidden layer was executed a small number of times ($1 \leq K \leq 5$), yielding averaged activation values. The output layer ($N \rightarrow 10$) was realized either optically—using an optical MVM and high photon budget to achieve high readout SNR, as in conventional ONNs, or with a digital electronic processor, yielding a result with full numerical precision. **b**, Simulated test accuracy of MNIST handwritten-digit classification for models with different numbers of hidden neurons N and shots per activation K . Error bars have been plotted but are small enough that they are difficult to see. **c**, Experimental evaluation of the SPDNN, with the output layer performed with full numerical precision on a digital computer. Results are presented for both $K = 1$ (single-shot, i.e., no averaging; top) and $K = 2$ (bottom) shots per activation. **d**, Experimental evaluation of the SPDNN, with both the hidden and the output layer executed using the optical experimental apparatus. The average number of detected photons used per inference in the hidden layer was kept fixed and the number used per inference in the output layer was varied.

148 single-hidden-layer network, efficiency gains of $>40\times$ are possible by performing the hidden layer in the
 149 low-SNR regime.

150 The models we report on in this section used non-negative weights in the hidden layers and real-valued
 151 weights in the output layers. This allows the hidden layers to be straightforwardly realized with optical
 152 MVMs using incoherent light.⁴ In the next section and Supplementary Note 2, we report on extensions
 153 to the case of real-valued weights in coherent optical processors.

⁴A high-SNR layer with real-valued weights can be realized with an incoherent optical MVM if some digital-electronic postprocessing is allowed [4, 43]—which is the approach we take for the optical output layer executions in our experiments. However, the postprocessing strategy doesn't directly apply in the low-SNR regime because readout becomes inseparable from the application of a nonlinear activation function, so we are constrained to non-negative weights and activations in the hidden layers.

Simulation results

First, we digitally simulated the SPDNN models shown in Figure 3a. We report the simulated test accuracies in Figure 3b for the full test dataset of 10,000 images, as a function of the number of hidden neurons N and the number of shots K of binary SPD measurements integrated to compute each activation.

Due to the stochastic nature of the model, the classification output for a fixed input varies from run to run. We repeated inferences on fixed inputs from the test set 100 times; we report the mean and standard deviation of the test accuracy as data points and error bars, respectively. The standard deviations of the test accuracies are around 0.1%.

If we integrated an infinite number of SPD measurements for each activation ($K \rightarrow \infty$)—which is infeasible in experiment, but can be simulated—then the SPDNN output would become deterministic. The test accuracy achieved in this limit can be considered as an upper bound, as the classification accuracy improves monotonically with K . Notably, even with just a single SPD measurement ($K = 1$) for each activation, the mean test accuracy is around 97%. The accuracy is substantially improved with just a few more shots of averaging, and approaches the deterministic upper bound when $K \gtrsim 5$. The mean single-photon-detection probability, averaged over all neurons, is ≈ 0.5 , so the simulated number of detected photons per shot is very small: $\approx 0.5N$. As we will quantify in the next section reporting the results of optical experiments, this means high accuracy can be achieved using much less optical energy than in conventional ONNs.

Optical experimental results

In our experimental demonstrations, we based our SPDNN on a free-space optical matrix-vector multiplier (MVM) that we had previously constructed for high-SNR experiments [4], and replaced the detectors with SPDs so that we could operate it with ultra-low photon budgets (see Methods). The experiments we report were, in part, enabled by the availability of cameras comprising large arrays of pixels capable of detecting single photons with low noise [44]. We encoded neuron values in the intensity of incoherent light; as a result, the weights and input vectors were constrained to be non-negative. However, this is not a fundamental feature of SPDNNs—in the next section, we present simulations of coherent implementations that lift this restriction. A single-photon-detecting camera measured the photons transmitted through the optical MVM, producing the stochastic activations as electronic signals that were input to the following neural-network layer (see Methods, Supplementary Notes 3 and 4).

In our first set of optical experiments, the hidden layer was realized optically and the output layer was realized *in silico* (Figure 3c): the output of the SPD measurements after the optical MVM was passed through a linear classifier executed with full numerical precision on a digital electronic computer. We tested using both $K = 1$ (no averaging) and $K = 2$ shots of averaging the stochastic binary activations in the hidden layer. The results agree well with simulations, which differ from the simulation results

188 shown in Figure 3b because they additionally modeled imperfections in our experimental optical-MVM
189 setup (see Methods, Supplementary Note 7). The test accuracies were calculated using 100 test images,
190 with inference for each image repeated 30 times. The hidden layer (the one computed optically in these
191 experiments) used approximately 0.0008 detected photons per MAC, which is ≥ 6 orders of magnitude
192 lower than is typical in ONN implementations [5, 13, 14, 28] and ≥ 3 orders of magnitude lower than the
193 lowest photons-per-MAC numbers reported to date [4, 5].

194 We then performed experiments in which both the hidden layer and the output layer were computed
195 optically (Figure 3d). In these experiments, we implemented a neural network with 400 hidden neurons
196 and used 5 shots per inference ($N = 400$, $K = 5$). The total optical energy was varied by changing the
197 number of photons used in the output layer; the number of photons used in the hidden layer was kept
198 fixed (see Methods, Table S6 and Supplementary Note 9).

199 The results show that even though the output layer was operated in the high-SNR regime (Figure
200 1b), the full inference computation achieved high accuracy yet used only a few femtojoules of optical
201 energy in total (equivalent to a few thousand photons). By dividing the optical energy by the number
202 of MACs performed in a single inference, we can infer the per-MAC optical energy efficiency achieved:
203 with an average detected optical energy per MAC of approximately 0.001 attojoules (0.003 attojoules),
204 equivalent to 0.003 photons (0.008 photons), the mean and standard deviation of test accuracy achieved
205 $92.0 \pm 2.3\%$ ($98.0 \pm 1.3\%$).

206 We can also compare our results with what has been published previously. Our experiments, with
207 $N = 50$ hidden neurons and $K = 5$ shots of SPD measurements per activation (see Supplementary
208 Figure 20) achieved a test accuracy of 90.6% on MNIST handwritten-digit recognition while using only
209 an average of 1390 detected photons per inference (corresponding to ~ 0.5 fJ of detected optical energy
210 per inference). This represents a $>40\times$ reduction in the number of photons per inference to achieve $>90\%$
211 accuracy on this task versus the previous state-of-the-art [4, 5].

212 **Simulation study of possible future deeper, coherent** 213 **single-photon-detection neural networks**

214 We have successfully experimentally demonstrated a two-layer SPDNN, but can SPDNNs be used to
215 implement deeper and more sophisticated models? One of the limitations of our experimental apparatus
216 was that it used an intensity encoding with incoherent light and as a result could natively only perform
217 operations with non-negative numbers. In this section we will show that SPDNNs capable of implementing
218 signed numbers can be used to realize multilayer models (with up to 6 layers), including models with
219 more sophisticated architectures than multilayer perceptrons—such as models with convolutional layers.

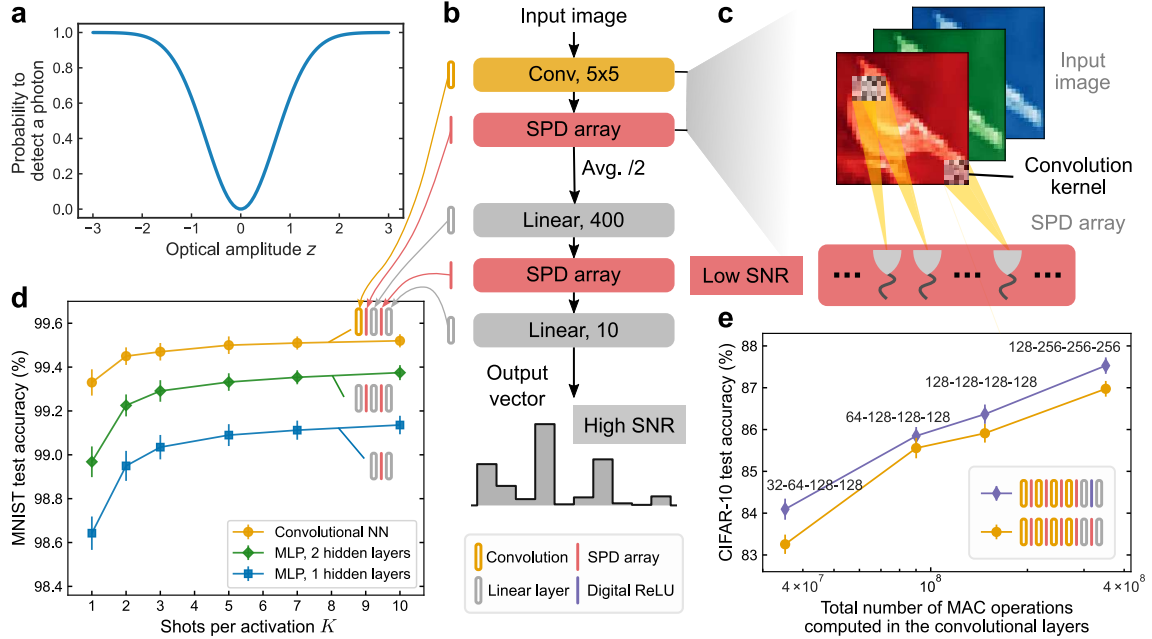


Figure 4. Simulation study predicting the performance of proposed *coherent* single-photon-detection neural networks (SPDNNs). **a**, The probability of detecting a photon as a function of the input light amplitude in a coherent SPDNN. Real-valued numbers are encoded in coherent light with either 0 phase (positive numbers) or π phase (negative numbers). Measurement by a single-photon detector (SPD) results in the probabilistic detection of a photon that is proportional to the square of the encoded value z , in comparison to intensity encodings with incoherent light. **b**, Structure of a convolutional SPDNN with a kernel size of 5×5 . Single-shot SPD measurements ($K = 1$) are performed after each layer (by an SPD array), except for the output layer. Average 2×2 pooling is applied after each convolutional operation. A digital rectified linear unit (ReLU) [45] activation function can also be used in the linear layer as an alternative. **c**, Schematic of a convolutional layer with SPD activations. **d**, Simulated test accuracy of coherent SPDNNs with varying architecture performing MNIST handwritten-digit classification. The multilayer perceptron (MLP) models had 400 neurons in each hidden layer. The convolutional model consisted of a convolutional layer with 16 output channels, followed by two linear layers with an SPD activation inbetween. **e**, Simulated test accuracy of coherent SPDNNs with varying architecture performing CIFAR-10 image classification. The models have four convolutional layers, each followed by SPD activation functions. The two linear layers can either be implemented in full-precision with a ReLU activation function (in purple) or using the SPD activation function. The number of output channels for each convolutional layer is indicated above the corresponding data point.

220 ONNs based on coherent light can naturally encode sign information in the phase of the light and
 221 have been realized in many different physical platforms [9, 10, 12, 13, 37, 46, 47]. We propose—and
 222 study in simulation—SPDNNs using coherent light. Neuron values are encoded in optical amplitudes
 223 that are constrained to have phases that are either 0 (positive values) or π (negative values). With this
 224 encoding, detection by an SPD—which measures intensity and is hence insensitive to phase—results
 225 in a stochastic nonlinear activation function that is symmetric about zero (Figure 4a; see Methods).
 226 Alternative detection schemes could be employed that would modify the activation function, but we have
 227 focused on demonstrating the capabilities of this straightforward case, avoiding introducing additional
 228 experimental complexity.

229 We performed two sets of simulation experiments: one on coherent SPDNNs trained to perform MNIST
230 handwritten-digit classification, and one on coherent SPDNNs trained to performed CIFAR-10 image
231 classification. Figure 4d shows the architectures tested and simulation results for the MNIST benchmark
232 (see Methods, Supplementary Note 2B). The accuracy achieved by MLPs with either one or two hidden
233 layers was higher than that of the single-hidden-layer MLP simulated for the incoherent case (Figure
234 3b), and an architecture with a single convolutional layer followed by two linear layers achieved >99%
235 accuracy even in the single-shot ($K = 1$) regime.

236 Figure 4e shows the results of simulating variants of a 6-layer convolutional SPDNN (comprising 4
237 convolutional layers and 2 fully connected, linear layers) on CIFAR-10 image classification. All these
238 simulation results were obtained in the single-shot ($K = 1$) regime. The number of channels in each
239 convolution layer was varied, which affects the total number of MACs used to perform an inference. We
240 observed that the test accuracy increased with the size of the SPDNN, with accuracies approaching those
241 of conventional convolutional neural networks of comparable size [48], as well as of binarized convolutional
242 neural networks [35, 49, 50]. In the models we simulated that only used SPD as the activation function
243 (i.e., the ones in which there are no ‘Digital ReLU’ blocks), the high-SNR linear output layer had only
244 4000 MAC operations, so the number of MACs in the high-SNR layer comprises less than 0.01% of the
245 total MACs performed during an inference. The models we simulated are thus sufficiently large that the
246 total optical energy cost would be dominated by the (low-SNR) layers prior to the (high-SNR) output
247 layer. Equivalently, the optical energy cost per MAC would be predominantly determined by the cost
248 of the low-SNR layers. These simulation results illustrate the ability of SPDNNs to scale to larger and
249 deeper models, enabling them to perform more challenging tasks.

250 Discussion

251 Our research is an example of realizing a neural network using a stochastic physical system. Beyond optics,
252 our work is related and complementary to recent investigations in electronic, spintronic, and quantum
253 neuromorphic computing [2, 51–58], including in training physical systems to perform neural-network
254 inference [23, 59–65]. Noise is a fundamental feature and the ultimate limit to energy efficiency in com-
255 puting with all analog physical systems. It has long been realized that noise is not always detrimental: not
256 only does it not necessarily prevent accurate computation, but can in some cases even enable fundamen-
257 tally new and more efficient algorithms or types of computation. Our work shows that using a quantum
258 physical model of a particular hardware’s noise at the software level can enable surprisingly large gains
259 in energy efficiency.

260 While there are many reasons computer science has traditionally favored the abstraction of hardware
261 from software, our work is part of a broad trend, spanning many different physical platforms [8, 66, 67], in
262 which researchers engineer computations in a physics-aware manner. By short-circuiting the abstraction

263 hierarchy—in our case, going from a physics-aware software description of a stochastic neural network
264 directly to a physical optical realization of the constituent operations—it is possible to achieve orders-of-
265 magnitude improvements in energy efficiency [3, 26] versus conventional CMOS computing. *Physics-aware*
266 *software*, in which software directly incorporates knowledge of the physics of the underlying computing
267 hardware—such as in the *physics-aware stochastic training* we used in this work—is understudied com-
268 pared to purely software-level or hardware-level innovations (i.e., “at the top” or “at the bottom” of the
269 hierarchy [68]). It is thus ripe for exploration: within the domain of neural networks, there are a multi-
270 tude of emerging physical platforms that could be more fully harnessed if the physical devices were not
271 forced to conform to the standard abstractions in modern computer architecture [23]. Beyond neural-
272 network accelerators, communities such as computational imaging [69] have embraced the opportunity
273 to improve system performance through co-optimizing hardware and software in a physics-aware man-
274 ner. We believe there is an opportunity to make gains in even more areas and applications of computing
275 technology by collapsing abstractions and implementing physics-aware software with physical hardware
276 that could be orders of magnitude faster or more energy efficient than current digital CMOS approaches
277 but that doesn’t admit a clean, digital, deterministic abstraction.

Methods

Stochastic optical neural networks using single-photon detection as the activation function

In the single-photon-detection neural networks (SPDNNs), the activation function is directly determined by the stochastic physical process of single-photon detection (SPD). The specific form of the activation function is dictated by the detection process on a single-photon detector (SPD). Each SPD measurement produces a binary output of either 0 or 1, with probabilities determined by the incident light intensity. Consequently, each SPD neuron activation, which corresponds to an SPD measurement in experiments, is considered as a binary stochastic process [70–72].

Following the Poisson distribution, the probability of an SPD detecting a photon click is given by $P_{\text{SPD}}(\lambda) = 1 - e^{-\lambda}$ when exposed to an incident intensity of λ photons per detection. Note that this photon statistics may vary based on the state of light (e.g. squeezed light), but here we only consider the Poissonian light. Therefore, the SPD process can be viewed as a Bernoulli sampling of that probability, expressed as $f_{\text{SPD}}(z) = \mathbf{1}_{t < P_{\text{SPD}}(\lambda(z))}$, where t is a uniform random variable $t \sim U[0, 1]$ and $\mathbf{1}_x$ is the indicator function that evaluates to 1 if x is true. This derivation leads to Equation 1 in the main text. In our approach, the pre-activation value z is considered as the direct output from an optical matrix-vector multiplier (MVM) that encodes the information of a dot product result. For the i th pre-activation value in layer l , denoted as $z_i^{(l)}$, the expression is given by:

$$z_i^{(l)} = \sum_{j=1}^{N_{l-1}} w_{ij}^{(l)} \cdot a_j^{(l-1)}, \quad (1)$$

where N_{l-1} is the number of neurons in layer $l - 1$, $w_{ij}^{(l)}$ is the weight between the i th neuron in layer l and the j th neuron in layer $l - 1$, $a_j^{(l-1)}$ is the activation of the j th neuron in layer $l - 1$. The intensity $\lambda(z)$ is a function of z that depends on the detection scheme employed in the optical MVM. In optical setups using incoherent light, the information is directly encoded in the intensity, resulting in $\lambda = z$. If coherent light were used in a setup, where 0 and π phases represent the sign of the amplitude, the intensity is determined by squaring the real-number amplitude if directly measured, resulting in $\lambda = z^2$. While more sophisticated detection schemes can be designed to modify the function of $\lambda(z)$, we focused on the simplest cases to illustrate the versatility of SPDNNs.

During the inference of a trained model, in order to regulate the level of uncertainty inherent in stochastic neural networks, we can opt to conduct multiple shots of SPD measurements during a single forward propagation. In the case of a K -shot inference, each SPD measurement is repeated K times, with the neuron’s final activation value $a^{[K]}$ being derived from the average of these K independent stochastic

299 binary values. Consequently, for a single shot, $a^{[1]} = a \in \{0, 1\}$; for K shots, $a^{[K]} = \frac{1}{K} \sum_{k=1}^K a_k \in$
 300 $\{0, 1/K, 2/K, \dots, 1\}$. By utilizing this method, we can mitigate the model’s stochasticity, enhancing the
 301 precision of output values. Ideally, with an infinite number of shots ($K \rightarrow \infty$), the activation $a^{[\infty]}$ would
 302 equate to the expected value without any stochasticity, that is, $a^{[\infty]} = \mathbb{E}[a] = P_{\text{SPD}}(\lambda(z))$. The detailed
 303 process of an inference of SPDNNs is described in Algorithm 2 in Supplementary Note 1A.

304 The training of our stochastic neuron models takes inspiration from recent developments in train-
 305 ing stochastic neural networks. We have created an effective estimator that trains our SPDNNs while
 306 accounting for the stochastic activation determined by the physical SPD process. To train our SPDNNs,
 307 we initially adopted the idea of the “straight-through estimator” (STE) [32, 73], which enables us to
 308 bypass the stochasticity and discretization during neural network training. However, directly applying
 309 STE to bypass the entire SPD process led to subpar training performance. To address this, we adopted a
 310 more nuanced approach by breaking down the activation function and treating different parts differently.
 311 The SPD process can be conceptually divided into two parts: the deterministic probability function P_{SPD}
 312 and the stochasticity introduced by the Bernoulli sampling. For a Bernoulli distribution, the expectation
 313 value is equal to the probability, making P_{SPD} the expectation of the activation. Instead of applying the
 314 “straight-through” method to the entire process, we chose to bypass only the Bernoulli sampling process.
 315 At the same time, we incorporate the gradients induced by the probability function, aligning them with
 316 the expectation values of the random variable. In this way, we obtained an unbiased estimator [74] for
 317 gradient estimation, thereby enhancing the training of our SPDNNs.

In the backward propagation of the l th layer, the gradients of the pre-activation $z^{(l)}$ can be computed
 as (the gradient with respect to any parameter x is defined as $g_x = \partial C / \partial x$ where C is the cost function):

$$g_{z^{(l)}} = \frac{\partial a^{(l)}}{\partial \lambda^{(l)}} \circ \frac{\partial \lambda^{(l)}}{\partial z^{(l)}} \circ g_{a^{(l)}} = P'_{\text{SPD}}(\lambda^{(l)}) \circ \frac{\partial \lambda^{(l)}}{\partial z^{(l)}} \circ g_{a^{(l)}}, \quad (2)$$

318 where $a^{(l)} = f_{\text{SPD}}(z^{(l)}) = \mathbf{1}_{t < P_{\text{SPD}}(\lambda(z^{(l)}))}$ and the gradients $g_{a^{(l)}}$ is calculated from the next layer (previous
 319 layer in the backward propagation). Using this equation, we can evaluate the gradients of the weights
 320 $W^{(l)}$ as $g_{W^{(l)}} = g_{z^{(l)}}^\top a^{(l-1)}$, where $a^{(l-1)}$ is the activation values from the previous layer. By employing
 321 this approach, SPDNNs can be effectively trained using gradient-based algorithms (such as SGD [75] or
 322 AdamW [76]), regardless of the stochastic nature of the neuron activations.

323 For detailed training procedures, please refer to Algorithm 1 and 3 in Supplementary Notes 1A and
 324 2A, respectively.

325 Simulation of incoherent SPDNNs for deterministic classification tasks

326 The benchmark MNIST (Modified National Institute of Standards and Technology database) [77] hand-
 327 written digit dataset consists of 60,000 training images and 10,000 testing images. Each image is a

328 grayscale image with $28 \times 28 = 784$ pixels. To adhere to the non-negative encoding required by incoherent
329 light, the input images are normalized so that pixel values range from 0 to 1.

330 To assess the performance of the SPD activation function, we investigated the training of the MLP-
331 SPDNN models with the structure of $784 \xrightarrow{W^{(1)}} N \xrightarrow{W^{(2)}} 10$, where N represents the number of neurons
332 in the hidden layer, $W^{(1)}$ ($W^{(2)}$) represents the weight matrices of the hidden (output) layer. The SPD
333 activation function is applied to the N hidden neurons, and the resulting activations are passed to the
334 output layer to generate output vectors (Figure 3a). To simplify the experimental implementation, biases
335 within the linear operations were disabled, as the precise control of adding or subtracting a few photons
336 poses significant experimental challenges. We have observed that this omission has minimal impact on
337 the model’s performance.

338 In addition, after each weight update, we clamped the elements of $W^{(1)}$ in the positive range in order
339 to comply with the constraint of non-negative weights of an incoherent optical setup. Because SPD is
340 not required at the output layer, the constraints on the last layer operation are less stringent. Although
341 our simulations indicate that the final performance is only marginally affected by whether the elements
342 in the last layer are also restricted to be non-negative, we found that utilizing real-valued weights in the
343 output layer provided increased robustness against noise and errors during optical implementation. As a
344 result, we chose to use real-valued weights in $W^{(2)}$.

345 During the training process, we employed the LogSoftmax function on the output vectors and used
346 cross-entropy loss to formulate the loss function. Gradients were estimated using the unbiased estimator
347 described in the previous section and Algorithm 1.

348 For model optimization, we found that utilizing the SGD optimizer with small learning rates yields
349 better accuracy compared to other optimizers such as AdamW, albeit at the cost of slower optimization
350 speed. Despite the longer total training time, the SGD optimizer leads to a better optimized model. The
351 models were trained with a batch size of 128, a learning rate of 0.001 for the hidden layer and 0.01 for
352 the output layer, over 10,000 epochs to achieve optimized parameters. To prevent gradient vanishing in
353 the plateau of the probability function P_{SPD} , pre-activations were clamped at $\lambda_{\text{max}} = 3$ photons.

354 It should be noted that due to the inherent stochasticity of the neural networks, each forward prop-
355 agation generates varying output values even with identical weights and inputs. However, we only used
356 one forward propagation in each step. This approach effectively utilized the inherent stochasticity in
357 each forward propagation as an additional source of random search for the optimizer. Given the small
358 learning rate and the significant noise in the model, the number of epochs exceeded what is typically
359 required for conventional neural network training processes. The training is performed on a GPU (Tesla
360 V100-PCIE-32GB) and takes approximately eight hours for each model.

361 We trained incoherent SPDNNs with a varying number of hidden neurons N ranging from 10 to 400.
362 The test accuracy of the models improved as the number of hidden neurons increased (see Supplementary

363 Note 1B for more details). During inference, we adjusted the number of shots per SPD activation K to
364 tune the SNR of the activations within the models.

365 For each model configuration with N hidden neurons and K shots of SPD readouts per activation, we
366 repeated the inference process 100 times to observe the distribution of stochastic output accuracies. Each
367 repetition of inference on the test set, which comprises 10,000 images, yielded a different test accuracy.
368 The mean values and standard deviations of these 100 repetitions of test accuracy are plotted in Figure
369 3b (see Supplementary Table 1 for more details). It was observed that increasing either N or K led to
370 higher mean values of test accuracy and reduced standard deviations.

371 **Experimental implementation of SPDNNs**

372 **Incoherent optical matrix-vector multiplier**

373 The optical matrix-vector multiplier setup utilized in this work is based on the design presented in [4].
374 The setup comprises an array of light sources, a zoom lens imaging system, an light intensity modulator,
375 and a photon-counting camera. For encoding input vectors, we employed an organic light-emitting diode
376 (OLED) display from a commercial smartphone (Google Pixel 2016 version). The OLED display features
377 a 1920×1080 pixel array, with individually controllable intensity for each pixel. In our experiment,
378 only the green pixels of the display were used, arranged in a square lattice with a pixel pitch of 57.5
379 μm . To perform intensity modulation as weight multiplication, we combined a reflective liquid-crystal
380 spatial light modulator (SLM, P1920-500-1100-HDMI, Meadowlark Optics) with a half-wave plate (HWP,
381 WPH10ME-532, Thorlabs) and a polarizing beamsplitter (PBS, CCM1-PBS251, Thorlabs). The SLM
382 has a pixel array of dimensions 1920×1152 , with individually controllable transmission for each pixel
383 measuring $9.2 \times 9.2 \mu\text{m}$. The OLED display was imaged onto the SLM panel using a zoom lens system
384 (Resolv4K, Navitar). The intensity-modulated light field reflected from the SLM underwent further de-
385 magnification and was focused onto the detector using a telescope formed by the rear adapter of the zoom
386 lens (1-81102, Navitar) and an objective lens (XLFLUOR4x/340, Olympus).

387 We decompose a matrix-vector multiplication in a batch of vector-vector dot products that are com-
388 puted optically, either by spatial multiplexing (parallel processing) or temporal multiplexing (sequential
389 processing). To ensure a more accurate experimental implementation, we chose to perform the vector-
390 vector dot products in sequence in most of the data collection. For the computation of an optical
391 vector-vector dot product, the value of each element in either vector is encoded in the intensity of the
392 light emitted by a pixel on the OLED and the transmission of an SLM pixel. The imaging system aligned
393 each pixel on the OLED display with its corresponding pixel on the SLM, where element-wise multiplica-
394 tion occurred via intensity modulation. The modulated light intensity from pixels in the same vector was
395 then focused on the detector to sum up the element-wise multiplication values, yielding the vector-vector
396 dot product result. Since the light is incoherent, only non-negative values can be allowed in both of the

397 vectors. For more details for the incoherent optical MVM, please refer to Supplementary Note 3. The
398 calibration of the vector-vector dot products on the optical MVM is detailed in Supplementary Note 5.

399 **Single-photon-detector array**

400 In this experiment, we used a scientific CMOS camera (Hamamatsu ORCA-Quest qCMOS Camera
401 C15550-20UP) [44] to measure both conventional light intensity measurement and SPD. This camera,
402 with 4096×2304 effective pixels of $4.6 \times 4.6 \mu\text{m}$ each, can perform SPD with ultra-low readout noise in
403 its photon counting mode. This scientific CMOS camera is capable of carrying out the SPD process with
404 ultra-low readout noise. When utilized as an SPD in the photon-counting mode, the camera exhibits an
405 effective photon detection efficiency of 68% and a dark count rate of approximately 0.01 photoelectrons
406 per second per pixel (Supplementary Note 4). We typically operate with an exposure time in the millisec-
407 ond range for a single shot of SPD readout. For conventional intensity measurement that integrates higher
408 optical energy for the output layer implementation, we chose another operation mode that used it as a
409 common CMOS camera. Further details on validating the stochastic SPD activation function measured
410 on this camera are available in Supplementary Note 6.

411 **Experimental implementation of the SPD activations**

412 We adapted our SPDNNs training methods to conform to the real-world constraints of our setup, ensuring
413 successful experimental implementation (see Supplementary Note 7). First, we conducted the implemen-
414 tation of the hidden layers and collect the SPD activations experimentally by the photon-counting camera
415 as an SPD array. Each SPD realized a stochastic activation function for a single hidden-layer neuron.
416 During a single inference, the hidden layer was executed a small number of times ($1 \leq K \leq 5$), yielding
417 averaged activation values. Then we performed the output layer operations digitally on a computer. This
418 aims to verify the fidelity of collecting SPD activations from experimental setups. Supplementary Figure
419 16 provides a visual representation of the distribution of some of the output vectors. For the experiments
420 with 1 shot per activation ($K = 1$), we collected 30 camera frames from the setup for each fixed input
421 images and weight matrix, which are regarded as 30 independent repetitions of inference. They were then
422 used to compute 30 different test accuracies by performing the output linear layer on a digital computer.
423 For the experiments with 2 shots per activation ($K = 2$), we divided the 30 camera frames into 15 groups,
424 with each group containing 2 frames. The average value of the 2 frames within each group serves as the
425 activations, which are used to compute 15 test accuracies. For additional results and details, please refer
426 to Supplementary Note 8.

427 **Optical implementation of the output linear layer**

428 Second, to achieve the complete optical implementation of the entire neural networks, we utilized our
429 optical matrix-vector multiplier again to carry out the last layer operations. For example, we first focused
430 on the data from the model with 400 hidden neurons and performed 5 shots per inference. In this case,
431 for the 30 binary SPD readouts obtained from 30 frames, we performed an averaging operation on every 5
432 frames, resulting in 6 independent repetitions of the inference. These activation values were then displayed
433 on the SLM as the input for the last layer implementation. For the 5-shot activations, the possible values
434 included 0, 0.2, 0.4, 0.6, 0.8, and 1. When the linear operation were performed on a computer with full
435 precision, the mean test accuracy was approximately 99.17%. To realize the linear operation with real-
436 valued weight elements on our incoherent optical setup, we divided the weight elements into positive
437 and negative parts. Subsequently, we projected these two parts of the weights onto the OLED display
438 separately and performed them as two different operations. The final output value was obtained by
439 subtracting the results of the negative weights from those of the positive weights. This approach requires
440 at least double the photon requirement for the output layer and offers room for optimization to achieve
441 higher energy efficiency. Nevertheless, even with these non-optimized settings, we demonstrated a photon
442 budget that is lower than any other ONN implementations known to us for the same task and accuracy.
443 For additional data and details, please refer to Supplementary Note 9.

444 **Deeper SPDNNs operating with coherent light**

445 Optical processors with coherent light have the ability to preserve the phase information of light and
446 have the potential to encode complex numbers using arbitrary phase values. In this work, we focused
447 on coherent optical computing utilizing real-number operations. In this approach, positive and negative
448 values are encoded in the light amplitudes corresponding to phase 0 and π , respectively.

449 As the intensity of light is the square of the amplitude, direct detection of the light amplitude, where
450 the information is encoded, would involve an additional square operation, i.e., $\lambda(z) = |z|^2$. This leads
451 to a “V-shape” SPD probability function with respect to the pre-activation z , as depicted in Figure 4a.
452 We chose to focus on the most straightforward detection case to avoid any additional changes to the
453 experimental setup. Our objective is to demonstrate the adaptability and scalability of SPDNN models
454 in practical optical implementations without the need for complex modifications to the existing setup.

455 **Coherent SPDNNs for MNIST classification**

456 **MLP-SPDNNs** Classifying MNIST using coherent MLP-SPDNNs was simulated utilizing similar con-
457 figurations as with incoherent SPDNNs. The only difference was the inclusion of the coherent SPD
458 activation function and the use of real-valued weights. Contrary to the prior scenario with incoherent
459 light, the input values and weights do not need to be non-negative. The models were trained using the

460 SGD optimizer [75] with a learning rate of 0.01 for the hidden layers and 0.001 for the last linear layer,
 461 over a period of 10,000 epochs.

462 **Convolutional SPDNNs** The convolutional SPDNN model used for MNIST digit classification, illus-
 463 trated in Figure 4b, consists of a convolutional layer with 16 output channels, a kernel size of 5×5 , a
 464 stride size of 1, and padding of 2. The SPD activation function was applied immediately after the con-
 465 volutional layer, followed by average pooling of 2×2 . The feature map of $14 \times 14 \times 16 = 3136$ was then
 466 flattened into a vector of size 3136. After that, the convolutional layers were followed by a linear model
 467 of $3136 \rightarrow 400 \rightarrow 10$, with the SPD activation function applied at each of the 400 neurons in the first
 468 linear layer.

469 The detailed simulation results of the MNIST test accuracies of the coherent SPDNNs can be found
 470 in Supplementary Table 2 with varying model structures and shots per activation K . For additional
 471 information, see Supplementary Note 2B.

472 **Coherent convolutional SPDNNs for Cifar-10 classification**

473 The CIFAR-10 dataset [78] has 60,000 images, each having $3 \times 32 \times 32$ pixels with 3 color channels, that
 474 belong to 10 different categories, representing airplanes, automobiles, birds, cats, deers, dogs, frogs, horses,
 475 ships and trucks. The dataset is partitioned into a training set with 50,000 images and a test set with
 476 10,000 images. The pixel values have been normalized using the mean value of (0.4914, 0.4822, 0.4465)
 477 and standard deviation of (0.2471, 0.2435, 0.2616) for each of the color channels. To boost performance,
 478 data augmentation techniques including random horizontal flips (50% probability) and random 32×32
 479 crops (with 4-pixel padding) were implemented during training.

480 The convolutional SPDNN models for Cifar-10 classification have deeper structures. Same as the
 481 convolutional models trained for MNIST, the convolutional layers use a kernel size of 5×5 , a stride size of
 482 1 and padding of 2. Each convolutional layer is followed by the SPD activation function, average pooling
 483 of 2×2 , as well as batch normalization. After N_{conv} convolutional layers ($N_{\text{conv}} = 4$ in Figure 4e) with
 484 the number of output channels of the last one to be $N_{\text{chan}}^{\text{last}}$, the feature map of $(32/2^{N_{\text{conv}}})^2 \times N_{\text{chan}}^{\text{last}}$ is
 485 flattened to a vector, followed by two linear layers of $(32/2^{N_{\text{conv}}})^2 N_{\text{chan}}^{\text{last}} \rightarrow 400 \rightarrow 10$. In the first linear
 486 layer, either SPD or ReLU [45] activation function were used for each of the 400 neurons, as depicted
 487 in Figure 4e. We vary the number of convolutional layers and number of output channels of them to
 488 change the different model size (Figure 4e and Supplementary Figure 5). In these results, we only used a
 489 single shot of SPD measurement ($K = 1$) to compute the SPD activations in the models, including the
 490 convolutional and linear layers. For additional information, please refer to Supplementary Note 2C.

491 **Data and code availability**

492 The data and code needed to reproduce the results presented in this paper are available for download
493 at <https://doi.org/10.5281/zenodo.8188270>. We have included the raw data resulting from our numerical
494 (simulation) and optical experiments, the code used to process this data, the training datasets and
495 trained-model parameters, as well as examples to demonstrate the operation of our data-collection and
496 data-processing code. We have also made available a pedagogical code repository, available at <https://github.com/mcmahon-lab/Single-Photon-Detection-Neural-Networks>, which may be adapted to train
497 models for different stochastic physical hardware setups.
498

499 **Acknowledgements**

500 We wish to thank NTT Research for their financial and technical support (S.-Y.M., P.L.M., T.W. and
501 L.G.W.). Portions of this work were supported by the National Science Foundation (award no. CCF-
502 1918549; J.L., P.L.M. and T.W.), a Kavli Institute at Cornell instrumentation grant (P.L.M. and T.W.),
503 and a David and Lucile Packard Foundation Fellowship (P.L.M.). P.L.M. acknowledges membership of
504 the CIFAR Quantum Information Science Program as an Azrieli Global Scholar. T.W. acknowledges
505 partial support from an Eric and Wendy Schmidt AI in Science Postdoctoral Fellowship. We acknowl-
506 edge discussions with M. Anderson, F. J. Chen, R. Hamerly, T. Onodera, S. Prabhu, M. M. Sohoni and
507 R. Yanagimoto. We also acknowledge Z. Eslami, V. Kremenetski, F. Presutti, C. Wan and F. Wu for
508 suggestions regarding the manuscript.

509 **Author Contributions**

510 S.-Y.M., L.G.W., T.W., and P.L.M. conceived the project. S.-Y.M. and T.W. designed the experiments
511 and built the experimental setup. S.-Y.M. and J.L. performed the neural-network training. S.-Y.M. per-
512 formed the experiments, the data analysis, and the numerical simulations. All authors contributed to
513 preparing the manuscript. T.W., L.G.W. and P.L.M. supervised the project.

514 **Additional information**

515 Supplementary information is available for this paper. Correspondence and requests for materials should
516 be addressed to Shi-Yuan Ma or Peter L. McMahon.

517 **References**

- 518 [1] Beenakker, C. & Schönberger, C. Quantum shot noise. *Physics Today* **56**, 37–42 (2003).

- 519 [2] Marković, D., Mizrahi, A., Querlioz, D. & Grollier, J. Physics for neuromorphic computing. *Nature*
520 *Reviews Physics* **2**, 499–510 (2020).
- 521 [3] Wetzstein, G. *et al.* Inference in artificial intelligence with deep optics and photonics. *Nature* **588**,
522 39–47 (2020).
- 523 [4] Wang, T. *et al.* An optical neural network using less than 1 photon per multiplication. *Nature*
524 *Communications* **13**, 1–8 (2022).
- 525 [5] Sludds, A. *et al.* Delocalized photonic deep learning on the internet’s edge. *Science* **378**, 270–276
526 (2022).
- 527 [6] LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* **521**, 436–444 (2015).
- 528 [7] Canziani, A., Paszke, A. & Culurciello, E. An analysis of deep neural network models for practical
529 applications. *arXiv:1605.07678* (2016).
- 530 [8] Christensen, D. V. *et al.* 2022 roadmap on neuromorphic computing and engineering. *Neuromorphic*
531 *Computing and Engineering* **2**, 022501 (2022).
- 532 [9] Shen, Y. *et al.* Deep learning with coherent nanophotonic circuits. *Nature Photonics* **11**, 441 (2017).
- 533 [10] Lin, X. *et al.* All-optical machine learning using diffractive deep neural networks. *Science* **361**,
534 1004–1008 (2018).
- 535 [11] Ríos, C. *et al.* In-memory computing on a photonic platform. *Science Advances* **5**, eaau5759 (2019).
- 536 [12] Xu, X. *et al.* 11 TOPS photonic convolutional accelerator for optical neural networks. *Nature* **589**,
537 44–51 (2021).
- 538 [13] Feldmann, J. *et al.* Parallel convolutional processing using an integrated photonic tensor core. *Nature*
539 **589**, 52–58 (2021).
- 540 [14] Zhou, T. *et al.* Large-scale neuromorphic optoelectronic computing with a reconfigurable diffractive
541 processing unit. *Nature Photonics* **15**, 367–373 (2021).
- 542 [15] Davis III, R., Chen, Z., Hamerly, R. & Englund, D. Frequency-encoded deep learning with speed-
543 of-light dominated latency. *arXiv:2207.06883* (2022).
- 544 [16] Ashtiani, F., Geers, A. J. & Aflatouni, F. An on-chip photonic deep neural network for image
545 classification. *Nature* **606**, 501–506 (2022).

- 546 [17] Moon, S., Shin, K. & Jeon, D. Enhancing reliability of analog neural network processors. *IEEE*
547 *Transactions on Very Large Scale Integration (VLSI) Systems* **27**, 1455–1459 (2019).
- 548 [18] Joshi, V. *et al.* Accurate deep neural network inference using computational phase-change memory.
549 *Nature Communications* **11**, 1–13 (2020).
- 550 [19] Semenova, N., Larger, L. & Brunner, D. Understanding and mitigating noise in trained deep neural
551 networks. *Neural Networks* **146**, 151–160 (2022).
- 552 [20] Klachko, M., Mahmoodi, M. R. & Strukov, D. *Improving noise tolerance of mixed-signal neural*
553 *networks*, 1–8 (IEEE, 2019).
- 554 [21] Zhou, C., Kadambi, P., Mattina, M. & Whatmough, P. N. Noisy machines: Understanding noisy neu-
555 ral networks and enhancing robustness to analog hardware errors using distillation. *arXiv:2001.04974*
556 (2020).
- 557 [22] Yang, X., Wu, C., Li, M. & Chen, Y. Tolerating noise effects in processing-in-memory systems
558 for neural networks: A hardware–software codesign perspective. *Advanced Intelligent Systems* **4**,
559 2200029 (2022).
- 560 [23] Wright, L. G. *et al.* Deep physical neural networks trained with backpropagation. *Nature* **601**,
561 549–555 (2022).
- 562 [24] Semenova, N. & Brunner, D. Noise-mitigation strategies in physical feedforward neural networks.
563 *Chaos: An Interdisciplinary Journal of Nonlinear Science* **32**, 061106 (2022).
- 564 [25] Wu, C. *et al.* Harnessing optoelectronic noises in a photonic generative network. *Science Advances*
565 **8**, eabm2956 (2022).
- 566 [26] Anderson, M. G., Ma, S.-Y., Wang, T., Wright, L. G. & McMahon, P. L. Optical transformers.
567 *arXiv:2302.10360* (2023).
- 568 [27] Jiang, Y. *et al.* Physical layer-aware digital-analog co-design for photonic convolution neural network.
569 *IEEE Journal of Selected Topics in Quantum Electronics* (2023).
- 570 [28] Bernstein, L. *et al.* Single-shot optical neural network. *Science Advances* **9**, eadg7904 (2023).
- 571 [29] Gerry, C. & Knight, P. *Introductory Quantum Optics* (Cambridge University Press, Cambridge,
572 2004).

- 573 [30] Machida, S., Yamamoto, Y. & Itaya, Y. Observation of amplitude squeezing in a constant-current-
574 driven semiconductor laser. *Physical Review Letters* **58**, 1000 (1987).
- 575 [31] Hadfield, R. H. Single-photon detectors for optical quantum information applications. *Nature*
576 *Photonics* **3**, 696–705 (2009).
- 577 [32] Bengio, Y., Léonard, N. & Courville, A. Estimating or propagating gradients through stochastic
578 neurons for conditional computation. *arXiv:1308.3432* (2013).
- 579 [33] Tang, C. & Salakhutdinov, R. R. Learning stochastic feedforward neural networks. *Advances in*
580 *Neural Information Processing Systems* **26** (2013).
- 581 [34] Gu, S., Levine, S., Sutskever, I. & Mnih, A. MuProp: Unbiased backpropagation for stochastic neural
582 networks. *arXiv:1511.05176* (2015).
- 583 [35] Hubara, I., Courbariaux, M., Soudry, D., El-Yaniv, R. & Bengio, Y. Binarized neural networks.
584 *Advances in Neural Information Processing Systems* **29** (2016).
- 585 [36] Liu, Y., Liu, S., Wang, Y., Lombardi, F. & Han, J. A survey of stochastic computing neural networks
586 for machine learning applications. *IEEE Transactions on Neural Networks and Learning Systems*
587 **32**, 2809–2824 (2020).
- 588 [37] Chang, J., Sitzmann, V., Dun, X., Heidrich, W. & Wetzstein, G. Hybrid optical-electronic convolu-
589 tional neural networks with optimized diffractive optics for image classification. *Scientific Reports*
590 **8**, 12324 (2018).
- 591 [38] Carolan, J. *et al.* Universal linear optics. *Science* **349**, 711–716 (2015).
- 592 [39] Bogaerts, W. *et al.* Programmable photonic circuits. *Nature* **586**, 207–216 (2020).
- 593 [40] Tait, A. N., Chang, J., Shastri, B. J., Nahmias, M. A. & Prucnal, P. R. Demonstration of WDM
594 weighted addition for principal component analysis. *Optics Express* **23**, 12758–12765 (2015).
- 595 [41] Rastegari, M., Ordonez, V., Redmon, J. & Farhadi, A. *XNOR-Net: Imagenet classification using*
596 *binary convolutional neural networks*, 525–542 (Springer, 2016).
- 597 [42] Zhou, S. *et al.* DoReFa-Net: Training low bitwidth convolutional neural networks with low bitwidth
598 gradients. *arXiv:1606.06160* (2016).
- 599 [43] Hayasaki, Y., Tohyama, I., Yatagai, T., Mori, M. & Ishihara, S. Optical learning neural network
600 using Selfoc microlens array. *Japanese Journal of Applied Physics* **31**, 1689 (1992).

- 601 [44] Dhimitri, K. *et al.* *Scientific CMOS (sCMOS) camera capabilities with a focus on quantum*
602 *applications*, PC122430L (International Society for Optics and Photonics, 2022).
- 603 [45] Agarap, A. F. Deep learning using rectified linear units (ReLU). *arXiv:1803.08375* (2018).
- 604 [46] Spall, J., Guo, X., Barrett, T. D. & Lvovsky, A. Fully reconfigurable coherent optical vector–matrix
605 multiplication. *Optics Letters* **45**, 5752–5755 (2020).
- 606 [47] Miscuglio, M. *et al.* Massively parallel amplitude-only Fourier neural network. *Optica* **7**, 1812–1819
607 (2020).
- 608 [48] Lee, C.-Y., Gallagher, P. W. & Tu, Z. *Generalizing pooling functions in convolutional neural networks:*
609 *Mixed, gated, and tree*, 464–472 (PMLR, 2016).
- 610 [49] Esser, S. K. *et al.* Convolutional networks for fast, energy-efficient neuromorphic computing.
611 *Proceedings of the National Academy of Sciences* **113**, 11441–11446 (2016).
- 612 [50] Qin, H. *et al.* Binary neural networks: A survey. *Pattern Recognition* **105**, 107281 (2020).
- 613 [51] Torrejon, J. *et al.* Neuromorphic computing with nanoscale spintronic oscillators. *Nature* **547**,
614 428–431 (2017).
- 615 [52] Grollier, J. *et al.* Neuromorphic spintronics. *Nature Electronics* **3**, 360–370 (2020).
- 616 [53] Cai, F. *et al.* Power-efficient combinatorial optimization using intrinsic noise in memristor Hopfield
617 neural networks. *Nature Electronics* **3**, 409–418 (2020).
- 618 [54] Harabi, K.-E. *et al.* A memristor-based bayesian machine. *Nature Electronics* **6**, 52–63 (2023).
- 619 [55] Islam, A. N. M. N. *et al.* Hardware in loop learning with spin stochastic neurons. *arXiv:2305.03235*
620 (2023).
- 621 [56] Marković, D. & Grollier, J. Quantum neuromorphic computing. *Applied Physics Letters* **117** (2020).
- 622 [57] Cerezo, M., Verdon, G., Huang, H.-Y., Cincio, L. & Coles, P. J. Challenges and opportunities in
623 quantum machine learning. *Nature Computational Science* **2**, 567–576 (2022).
- 624 [58] Roques-Carmes, C. *et al.* Biasing the quantum vacuum to control macroscopic probability
625 distributions. *Science* **381**, 205–209 (2023).
- 626 [59] Prezioso, M. *et al.* Training and operation of an integrated neuromorphic network based on metal-
627 oxide memristors. *Nature* **521**, 61–64 (2015).

- 628 [60] Romera, M. *et al.* Vowel recognition with four coupled spin-torque nano-oscillators. *Nature* **563**,
629 230–234 (2018).
- 630 [61] Mitarai, K., Negoro, M., Kitagawa, M. & Fujii, K. Quantum circuit learning. *Physical Review A* **98**,
631 032309 (2018).
- 632 [62] Hughes, T. W., Williamson, I. A., Minkov, M. & Fan, S. Wave physics as an analog recurrent neural
633 network. *Science Advances* **5**, eaay6946 (2019).
- 634 [63] Chen, T. *et al.* Classification with a disordered dopant-atom network in silicon. *Nature* **577**, 341–345
635 (2020).
- 636 [64] Cramer, B. *et al.* Surrogate gradients for analog neuromorphic computing. *Proceedings of the*
637 *National Academy of Sciences* **119**, e2109194119 (2022).
- 638 [65] Ross, A. *et al.* Multilayer spintronic neural networks with radiofrequency connections. *Nature*
639 *Nanotechnology* 1–8 (2023).
- 640 [66] Berggren, K. *et al.* Roadmap on emerging hardware and technology for machine learning,
641 *Nanotechnology* **32**, 012002 (2020).
- 642 [67] Finocchio, G. *et al.* Roadmap for unconventional computing with nanotechnology. *arXiv:2301.06727*
643 (2023).
- 644 [68] Leiserson, C. E. *et al.* There’s plenty of room at the top: What will drive computer performance
645 after Moore’s law? *Science* **368**, eaam9744 (2020).
- 646 [69] Kellman, M., Lustig, M. & Waller, L. How to do physics-based learning. *arXiv:2005.13531* (2020).
- 647 [70] Neal, R. M. Learning stochastic feedforward networks. *Department of Computer Science, University*
648 *of Toronto* **64**, 1577 (1990).
- 649 [71] Lee, V. T., Alaghi, A., Hayes, J. P., Sathe, V. & Ceze, L. *Energy-efficient hybrid stochastic-binary*
650 *neural networks for near-sensor computing*, 13–18 (IEEE, 2017).
- 651 [72] Liu, Y., Liu, S., Wang, Y., Lombardi, F. & Han, J. A stochastic computational multi-layer perceptron
652 with backward propagation. *IEEE Transactions on Computers* **67**, 1273–1286 (2018).
- 653 [73] Yin, P. *et al.* Understanding straight-through estimator in training activation quantized neural nets.
654 *arXiv:1903.05662* (2019).

- 655 [74] Williams, R. J. Simple statistical gradient-following algorithms for connectionist reinforcement
656 learning. *Machine learning* **8**, 229–256 (1992).
- 657 [75] Bottou, L. Stochastic gradient descent tricks. *Neural Networks: Tricks of the Trade: Second Edition*
658 421–436 (2012).
- 659 [76] Loshchilov, I. & Hutter, F. Decoupled weight decay regularization. *arXiv:1711.05101* (2017).
- 660 [77] De Chazal, P., Tapson, J. & Van Schaik, A. *A comparison of extreme learning machines and back-*
661 *propagation trained feed-forward networks processing the MNIST database*, 2165–2168 (IEEE, 2015).
- 662 [78] Krizhevsky, A. Learning multiple layers of features from tiny images. *Technical Report* 32–33 (2009).

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [supplementary.pdf](#)