

# All-Optical Deep Learning with Quantum Nonlinearity

Qingyi Zhou,<sup>1,\*</sup> Jungmin Kim,<sup>1,\*</sup> Yutian Tao,<sup>2,\*</sup> Guoming Huang,<sup>1</sup> Ming Zhou,<sup>3</sup> Zewei Shao,<sup>1</sup> and Zongfu Yu<sup>1</sup>

<sup>1</sup>*Department of Electrical and Computer Engineering,  
University of Wisconsin-Madison, Madison, WI 53706, USA*

<sup>2</sup>*The Computer Sciences Department, University of Wisconsin-Madison, Madison, WI 53706, USA*

<sup>3</sup>*Department of Electrical Engineering, Stanford University, Stanford, CA 94305, USA*

(Dated: January 6, 2026)

The rapid scaling of deep neural networks comes at the cost of unsustainable power consumption. While optical neural networks offer an alternative, their capabilities remain constrained by the lack of efficient optical nonlinearities. To address this, we propose an all-optical deep learning architecture by embedding quantum emitters in inverse-designed nanophotonic structures. Due to their saturability, quantum emitters exhibit exceptionally strong nonlinearity compared with conventional materials. Using physics-aware training, we demonstrate that the proposed architecture can solve complex tasks, including nonlinear classification and reinforcement learning, which have not been realized in all-optical neural networks. To enable fair comparison across different platforms, we introduce a framework that quantitatively links nonlinearity to a network's expressive power. Analysis shows that our quantum activation, operating below  $\text{nW}/\mu\text{m}^2$  intensity, reduces the power budget by seven orders of magnitude. System-level estimates show that the optical power required for large language models scales sublinearly with model size, enabling watt-level operation. Our results indicate that quantum nanophotonics provides a route toward sustainable AI inference.

## I. INTRODUCTION

In the past decade, the rapid advancement of deep learning has profoundly transformed science and technology. Deep neural networks have achieved state-of-the-art performance across diverse fields, ranging from computer vision [1] and game-playing [2] to protein design [3] and language processing [4]. Such progress has been driven by the continuous scaling of the model size. However, this scaling trend imposes an energy cost toward unsustainable levels [5]. A growing effort has been directed towards finding alternative computing paradigms. In particular, following the pioneering work of Shen et al. [6], optical neural networks (ONNs) have emerged as a promising candidate [7], inspired by the vision that a passive optical device can implement linear transformation with high speed [8] and low energy cost [9]. Specifically, recent works have demonstrated that linear optical matrix operations can be performed with sub-photon energy consumption [10]. Despite these advantages, linear operations are insufficient for deep learning. The expressive power of ONN is severely limited by the lack of efficient optical nonlinearity [9, 11]. In conventional materials, optical nonlinearities are often perturbative [12]. As a result, existing all-optical nonlinear activation units demand high optical power and large footprint [13–17], making it difficult to scale up. This nonlinearity bottleneck has remained a long-standing challenge for the optical computing community. Existing works that rely on hybrid opto-electronic architectures [6, 18, 19] suffer from additional latency and substantial system complexity due to frequent optical-electrical-optical (O/E/O) conversions. More recently, “structural nonlinearity” schemes have been proposed [20–22], in which the input is encoded not in the optical field but in tunable parameters of a linear structure. However, the connection between such systems and standard deep learning models is often unclear.

To address the nonlinearity bottleneck, we first point out that nonlinear optical phenomena are not intrinsically restricted to high intensities. A single quantum emitter (including atom, quantum dot, or color center) can be saturated by the absorption of a few photons per lifetime, leading to extremely strong optical nonlinearity [23]. There have been both theoretical proposals [24, 25] and experimental demonstrations [26–28] of using quantum emitter media as activation units, underscoring their potential for realizing strong optical nonlinearities. However, to fully utilize the nonlinear functionality of quantum emitters, it is essential to enhance the interaction between light and quantum emitters, which can be achieved using properly designed nanophotonic structures. Recent progress in quantum nanophotonics has enabled deterministic integration of individual emitters with on-chip photonic structures [29–31]. Therefore, we believe it is the right time to systematically investigate whether quantum technologies can provide the strong nonlinearity required by optical neural networks.

In this work, we propose a low-power optical neural network architecture that takes advantage of the strong nonlinearity of individual quantum emitters. Specifically, we introduce a quantum-enhanced activation unit by embedding

---

\* These authors contributed equally to this work. Correspondence: qzhou75@wisc.edu.

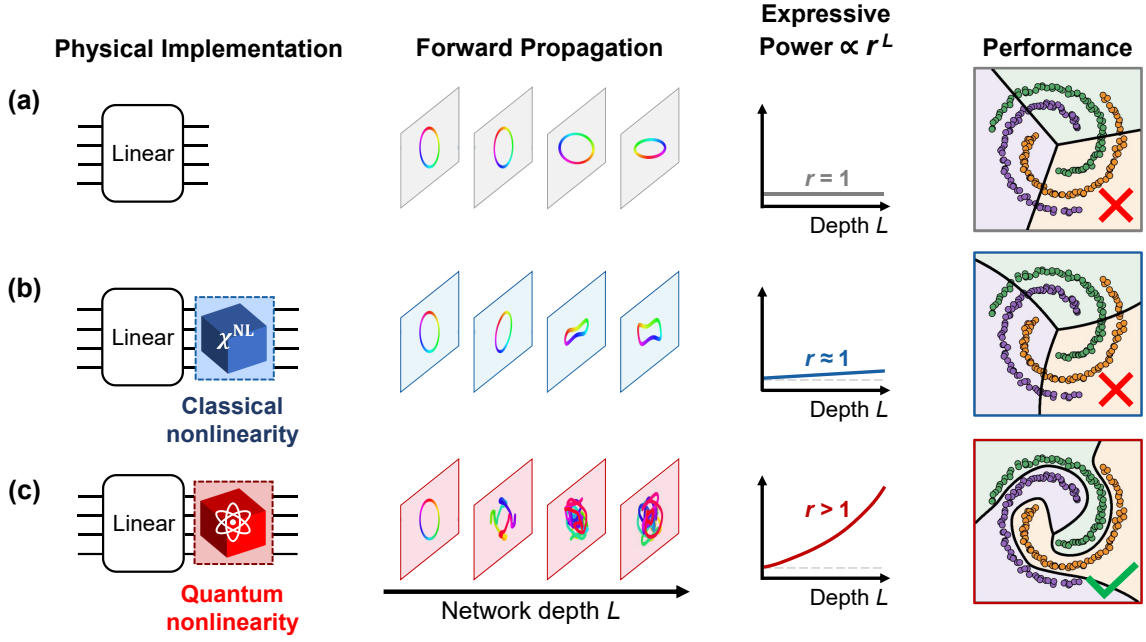


FIG. 1. **Comparison of ONN architectures with different nonlinearities.** (a) A linear ONN, whose expressive power remains constant regardless of depth  $L$ . (b) ONN with classical nonlinearity. With weak nonlinearity, the expressive power increases slightly with depth, yet is not enough for complex tasks. (c) ONN with quantum nonlinearity. The expressive power  $\propto r^L$  grows exponentially with  $r > 1$ , and is able to handle complex task.

quantum emitters into adjoint-optimized nanophotonic structures [32]. A strong nonlinear response can be achieved at intensities below  $1 \text{ nW}/\mu\text{m}^2$ . With full-wave simulations, we verify the performance on nonlinear classification task as well as a reinforcement learning tasks, demonstrating functionality beyond linear models. To enable a fair comparison across different physical nonlinearities, we develop a theoretical framework that quantifies the “expressive power” of an arbitrary activation unit. Unlike existing theoretical analyses that focus on purely linear optics [33–36] or restricted classes of unitary transformations [37], our framework directly links a nonlinear input-output response to the depth-wise growth of ONN expressivity. This allows us to translate a targeted expressive power into the required light intensity, enabling a quantitative assessment of ONN’s energy efficiency. Applying this framework, we show that conventional platforms based on Kerr effect or saturable absorption would require prohibitively high intensities to match the digital baseline. In contrast, the proposed quantum-enhanced activation reduces the required intensity by seven orders of magnitude. Finally, we look into the future by estimating the optical power consumption of running all-optical large language models (LLMs). Our analysis shows that the proposed scheme can, in principle, reduce the optical power consumption to a few watts, with a favorable sublinear scaling in model size. Taken together, our results indicate that large-scale all-optical deep learning, powered by quantum-enhanced nonlinearities, could reduce the energy footprint of AI inference, providing a path beyond the limits of electronic hardware.

## II. RESULTS

### A. Overcoming nonlinearity bottleneck with quantum activations

We consider a generic multi-layer ONN architecture. Each layer consists of a linear transformation  $W^{(i)}$  followed by a nonlinear activation  $f(\cdot)$ , mirroring the architecture of a typical multi-layer perceptron (MLP). In machine learning theory, it is well established that the expressive power of a deep neural network grows exponentially with its depth [38–40]. A purely linear network cannot enjoy this benefit since a composition of linear transformations is still linear. As a result, the overall expressive power of an ONN is limited by its nonlinear activation units. We follow the framework developed in Refs. 39, 40 and introduce a metric for quantifying expressive power. As illustrated in Fig. 1, a closed trajectory of data points serves as the input. The total curvature  $K$  of this trajectory provides a robust measure of curve complexity, and is monitored as the curve propagates through successive layers. Intuitively, the total curvature measures the degree of “folding” applied to the data manifold, a capability that is fundamentally

impossible with linear transformations. As shown in Fig. 1(a), for a purely linear network the total curvature remains constant. In contrast, in the presence of nonlinear activations, the curvature increases by a growth factor  $r > 1$  after each layer, leading to an exponential growth  $\sim r^L$  with depth  $L$  [39, 40]. The growth factor  $r$  is therefore used as a quantitative measure of expressive power (see Supplementary Note S10 for details).

Conventional optical materials possess small nonlinear susceptibilities [12]. At realistic light intensities, the response is only weakly nonlinear, which yields a growth factor  $r \approx 1$ . The expressive power, as shown in Fig. 1(b), shows little increase with depth. Existing all-optical activation units typically require  $\text{mW}/\mu\text{m}^2$  laser intensities, together with large footprints [13, 14, 16, 41] (see Supplementary Note S1 for a summary of representative designs obtained from literature). Such requirements are incompatible with large-scale ONNs. On the other hand, it has long been recognized that low-dimensional systems exhibit much stronger optical nonlinearities than bulk media [42–44], owing to enhanced oscillator strength under quantum confinement [45]. In particular, zero-dimensional quantum emitters behave as two-level systems (TLSs) that can be saturated by the absorption of only a few photons per lifetime. This leads to extremely strong nonlinear scattering responses, which have been observed in various waveguide- and cavity-QED platforms [46–50]. These observations suggest that emitter-based nonlinearities could be leveraged to overcome the nonlinearity bottleneck in ONNs (Fig. 1(c)).

## B. Device design and verification on nonlinear classification

Motivated by the above considerations, we propose an all-optical activation unit, consisting of two quantum emitters embedded inside a silicon nanophotonic structure (footprint  $5 \times 1 \mu\text{m}^2$ ), as shown in Fig. 2(a). We utilized adjoint optimization to design a nanophotonic interface that maximizes light-matter interaction and minimizes loss (see Supplementary Note S4 for design details). Each emitter is modeled as a TLS with a field-driven dipole moment  $d_z \propto \frac{\Omega/\Gamma_0}{1+2(\Omega/\Gamma_0)^2}$ , where the Rabi frequency  $\Omega$  is set by the local electric field [51, 52] and  $\Gamma_0 = 2\pi \times 94 \text{ MHz}$  is the spontaneous emission rate (parameters are obtained from  $\text{SiV}^-$  color centers, assuming lifetime-limited linewidth; see Supplementary Note S2 for details). The electric field distributions for two different input intensities are shown in Fig. 2(b), both obtained using nonlinear finite-difference frequency domain (FDFD) simulations (see Supplementary Note S3 for details). We provide 2D demonstrations, because on-chip architectures are common for constructing ONNs, and can often be approximated well using 2D effective models [53]. The device is engineered to operate in two distinct regimes: in the weak-field limit the emitters act as linear scatterers, while in the strong-field limit they become nearly transparent. In this two-port geometry, the two emitters are configured to induce a maximal change in transmission coefficient  $|\Delta t| = 2$  via interference, which results in a very strong nonlinearity. Analysis shows that such a  $|\Delta t| = 2$  change is impossible with a single emitter (see Supplementary Note S4 for the proof). From the simulated input-output curve, we extract an effective nonlinear activation function with an ultra-low intensity threshold. To evaluate the expressive power of the obtained activation function, we compute its growth factor  $r(I)$ , which reaches  $r \approx 1.2$  at intensity  $I = 1 \text{ nW}/\mu\text{m}^2$ , exceeding the digital baseline. In contrast, conventional silicon- and graphene-based nonlinearities remain near  $r = 1$  under similar operating conditions (see Supplementary Note S11 and Supplementary Figure S19 for details). We have also analyzed the effect of low quantum efficiency and demonstrate that our method remains robust even when the TLS’ quantum efficiency drops to 60% (see Supplementary Note S5 for details). These results confirm our key physical intuition: by utilizing saturable quantum emitters embedded in nanophotonic structures, strong optical nonlinearity can be realized at ultra-low optical power.

To demonstrate the practical advantage of the proposed activation, we benchmark our system on nonlinear classification task that is challenging for models that are weakly nonlinear. We adopt a physics-aware training approach to design our ONN in a physically consistent manner. We first characterize the nonlinear activation unit using nonlinear FDFD simulations to obtain its input-output transmission curve (see Supplementary Figure S8), which is then used as the activation function in a PyTorch-based ONN model. The linear weight matrices are represented as complex transmission matrices subject to energy-preserving constraints, ensuring that they can be implemented by passive structures. With this differentiable model, the ONN is trained in the digital domain via backpropagation. After training, we map the trained network to concrete photonic structures in a modular fashion. For each layer, the trained complex weight matrix is interpreted as a target transmission matrix between input and output ports. We then run a separate adjoint-based optimization to realize a compact block that implements this target matrix with high fidelity (see Supplementary Note S8 for quantitative metrics). The nonlinear layers are implemented by inserting the quantum activation units between these inverse-designed linear blocks, as illustrated in Fig. 2(a). By dissecting the network into multiple modules that can be designed individually, this approach avoids heavy full-wave simulations of the entire network during training (the training procedure is explained in Supplementary Note S7). We verify the final design using full-wave nonlinear FDFD simulations. The classification result for a three-class “spiral” dataset is shown in Fig. 2(c). The corresponding light intensities for three representative input points are also visualized, revealing how the network steers optical energy toward different output regions associated with different classes. We

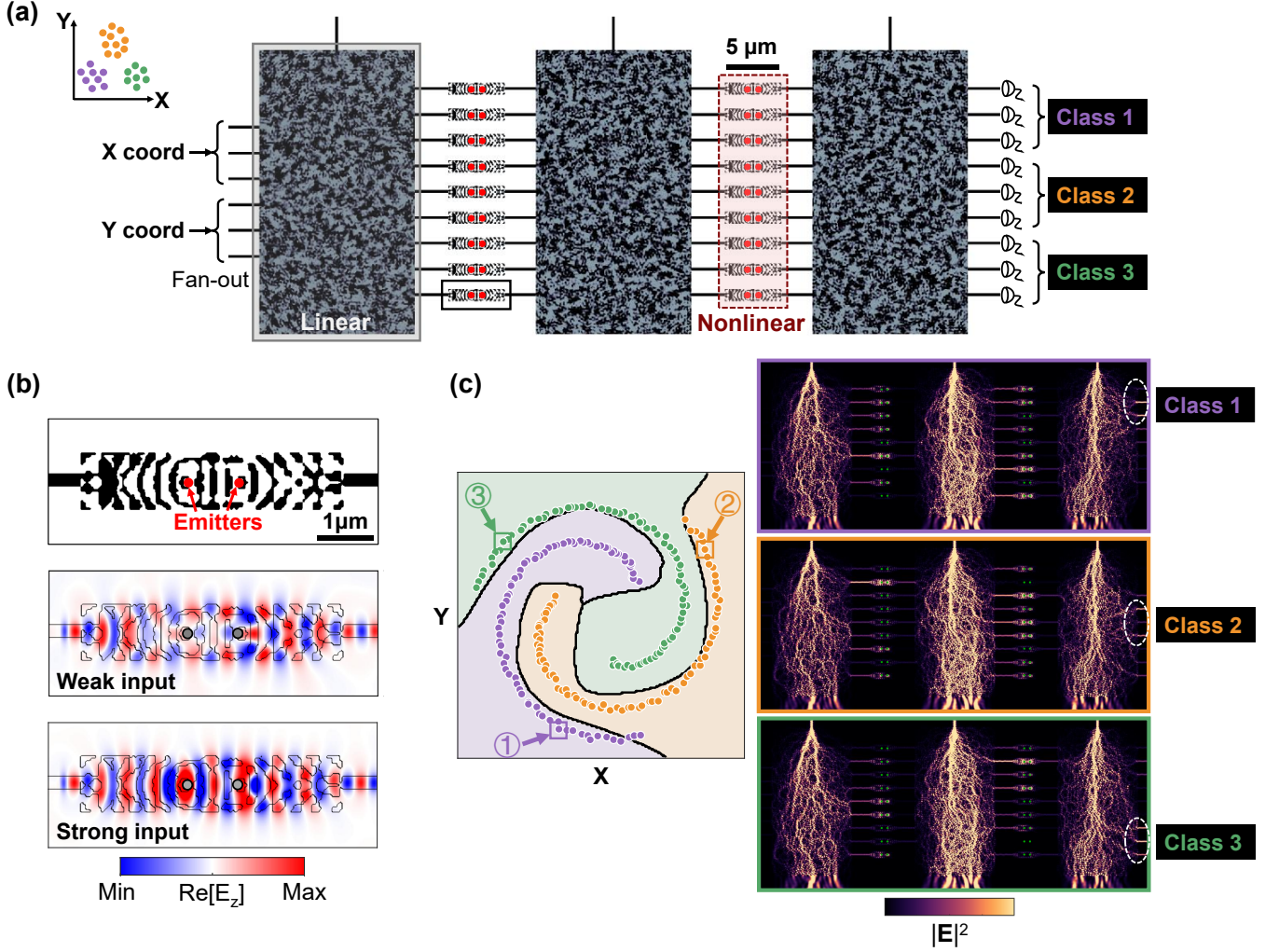


FIG. 2. **Physics-aware training and verification on nonlinear classification.** (a) All-optical neural network design, constructed by stacking nonlinear activation units between linear blocks, which are also designed through adjoint optimization. The 2D coordinates are encoded as input; the detected intensities are interpreted as classification results. (b) The proposed quantum nonlinear activation unit. Two emitters are embedded in an inverse-designed silicon structure. Simulated  $\text{Re}[E_z]$  field distributions illustrate the transition from resonant scattering to saturation, resulting in a nonlinear response at low intensity. (c) Performance verification. Classification results for the “spiral” dataset are obtained via nonlinear FDTD simulations. The  $|E|^2$  intensity distributions of 3 representative inputs are visualized. Different intensity distributions at the output port correspond to different predictions.

provide more examples in Supplementary Note S6 (performance on MNIST and FashionMNIST) and S8 (performance on nonlinear regression task). These results confirm that the physics-aware training yields physically realizable ONNs. Within realistic optical intensities well below  $1 \text{ mW}/\mu\text{m}^2$ , conventional materials cannot provide sufficient expressive power. In stark contrast, quantum-enhanced nonlinearity can function at  $1 \text{ nW}/\mu\text{m}^2$ , enabling the system to solve complex tasks.

### C. Generalizing to intelligent optical agents: Reinforcement learning

Having established that quantum-enhanced activations enable nonlinear classification, we next ask whether the same photonic building blocks can support more complex tasks. Reinforcement learning (RL), which has achieved impressive results on game-playing benchmarks [2, 54] and robotics [55], provides a natural testbed for our purpose. To demonstrate the generality of our approach, we choose two tasks, namely the “Atari Pong” game and the “HalfCheetah” control task, both provided by the Gymnasium library [56].

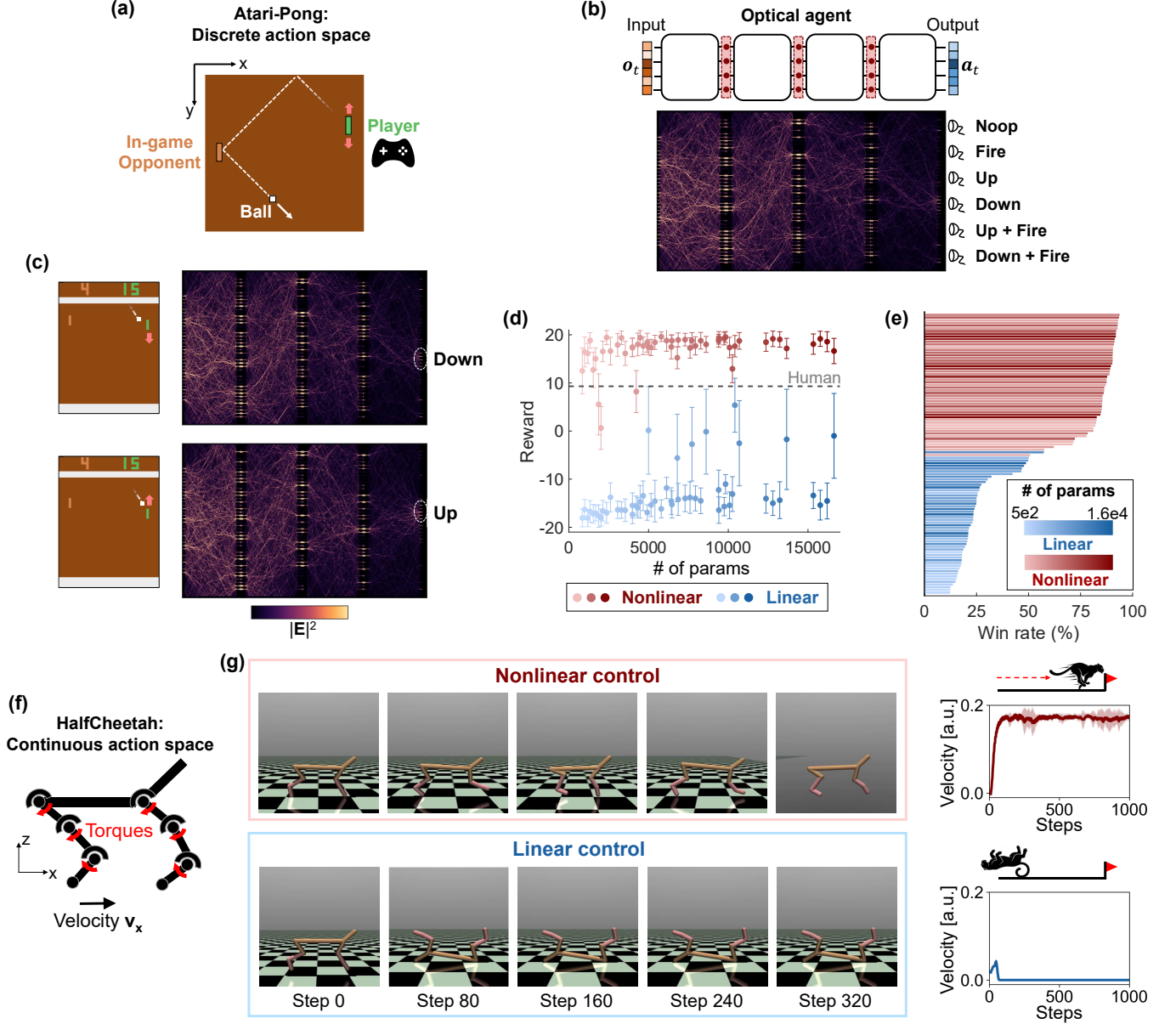


FIG. 3. **Generalizing to reinforcement learning tasks.** (a) Schematic illustration of the Pong environment. The player controls the green paddle and tries to block the ball (white) to win score. (b) The structure of our optical agent, which functions as a policy network. At each time step  $t$ , historical game frames are encoded into the input  $o_t$ . The network outputs the logits for six discrete actions. (c) Visualization of the learned policy. Light intensity distributions of two key frames are displayed. (d) Final reward versus model size. Linear models (blue) saturate at a low performance ceiling with high variance. In contrast, nonlinear models (red) converge to near-perfect play very quickly. Error bars denote standard deviations obtained over 30 episodes. The human-level performance (reward= 9.3) [2] is visualized using gray dashed line. (e) Performance ranking. Best achieved rewards for 104 trained models (52 linear, 52 nonlinear) are sorted. Nonlinear models consistently outperform linear models. (f) Schematic illustration of the HalfCheetah control task. (g) Snapshots obtained during testing. The nonlinear ONN runs stably, while the linear ONN falls down. The insets plot the corresponding velocity curves, averaged over 10 episodes. The shaded areas visualize the standard deviations.

The Atari Pong environment is illustrated schematically in Fig. 3(a). In Pong, an agent controls the right paddle against the in-game opponent. An episode ends when a player reaches 21 points, and the reward is defined as the final score difference. As shown in Fig. 3(b), the ONN acts as a policy network: at each time step  $t$ , a stack of  $F$  recent frames is encoded into an observation  $o_t$ , which serves as the input (see Supplementary Figure S14 for details). The output intensities (divided into six regions) are interpreted as logits for six discrete actions. An action  $a_t$  is



sampled from this distribution and sent back to the environment, which then advances to the next time step. The policy parameters are trained using a standard proximal policy optimization (PPO) algorithm [57], based on the same physics-aware framework described above (see Supplementary Note S9 for details). The optical intensity distributions for two representative game frame are shown in Fig. 3(c). Different spatial configurations of the ball and paddles lead to distinct activation patterns and different intensity hotspots at the output ports. We then systematically benchmark the performance of linear versus nonlinear ONNs. We train 104 models, sweeping across network width  $W$ , number of hidden layers  $L$ , and the number of input frames  $F$ . All models are trained for identical number of iterations (see Supplementary Note S9 for details). The results are summarized in Fig. 3(d), where the final reward is plotted against number of parameters. Fig. 3(e) further ranks all the trained models based on their final rewards. Linear models saturate at a low performance ceiling regardless of model size, and exhibit high variance, indicating that the learned strategies cannot win reliably. In contrast, ONNs equipped with quantum activations achieve much higher rewards as they scale up, converging to near-perfect play.

We further evaluate our ONN on the MuJoCo HalfCheetah control benchmark, as illustrated in Fig. 3(f). With a continuous action space, HalfCheetah is much more challenging than Pong. At each time step, the optical agent receives a 17-dimensional observation (joint positions and velocities, see Supplementary Figure S15) and outputs a 6-dimensional action vector that specifies torques applied at the six hinge joints. We adopt the similar ONN backbone as in Pong, with one key difference: the output uses balanced detection to support negative action values (see Supplementary Note S9 for details). Training is performed using the standard soft actor-critic (SAC) algorithm [58]. The corresponding snapshots collected during testing are shown in Fig. 3(g). With quantum activation the ONN learns to run smoothly, whereas a linear ONN fails to acquire a viable control policy and falls down early in the episode. The insets in Fig. 3(g) plot the averaged velocity  $v_x$  over 10 episodes, highlighting the stability enabled by optical nonlinearity. The above results confirm that strong nonlinearity is essential for enabling complex capabilities in deep ONNs.

#### D. Scalability and energy advantage of large-scale ONNs

Having established the importance of strong nonlinearity at the device level, we now address the central question of scalability: to what extent can such nonlinearities support truly large-scale systems, such as modern language models? As a quantitative baseline, we note that in standard digital MLPs, common activation functions typically increase the total curvature by  $r_{\text{digital}} \simeq 1.045 \sim 1.095$  per layer (see Supplementary Note S10 for details). We therefore ask: for a given physical nonlinearity, what is the minimum optical intensity  $I_{\text{min}}$  required to match this digital baseline? Using our established framework, we compute the expressive power  $r(I)$  for three representative platforms: Kerr nonlinearity in a 50  $\mu\text{m}$  long silicon waveguide [59, 60], saturable absorption in stacked graphene layers [42, 61] (15 nm total thickness), and our proposed quantum activation unit (see Supplementary Note S11 for details). We find that maintaining the target expressive power in silicon requires intensities exceeding  $72.6 \text{ W}/\mu\text{m}^2$ , while graphene requires approximately  $0.02 \text{ W}/\mu\text{m}^2$ . In contrast, the quantum activation unit achieves the same baseline at merely  $\sim 0.5 \text{ nW}/\mu\text{m}^2$ . This represents an efficiency improvement of roughly  $4.1 \times 10^7$  times relative to graphene and  $1.5 \times 10^{11}$  times relative to silicon.

Given these intensity thresholds, we next estimate the total optical power required to implement all-optical LLMs. For a standard decoder-only transformer architecture with context length  $L_{\text{seq}}$ , embedding dimension  $d_{\text{model}}$ , and  $L$  transformer layers [4], we estimate the optical input dimension per layer as  $3L_{\text{seq}}d_{\text{model}}$ , accounting for the parallel projection of query, key, and value matrices. Assuming each optical neuron occupies an effective cross-sectional area  $A \approx 0.1 \mu\text{m}^2$  [6] and is driven at  $I_{\text{min}}$ , the total optical power is estimated as

$$P \approx I_{\text{min}} A \cdot (3L_{\text{seq}}d_{\text{model}}) \cdot L. \quad (1)$$

Using the architectural parameters of representative LLMs ranging from GPT-2 to DeepSeek-V3 [62–69] (see Supplementary Note S12 for details), we evaluate Eq. (1) and summarize the results in Fig. 4(b). When using conventional nonlinearities based on silicon or graphene, the required optical power quickly reaches prohibitive levels (exceeding  $10^8 \text{ W}$  for the largest models). However, the proposed quantum architecture keeps the total optical power below 2.6 W across all investigated models. For reference, we also plot the estimated power consumption of high-end GPU (NVIDIA A100) during inference. Finally, we analyze how this power requirement scales with model size. Fig. 4(c) plots the estimated optical power against the total number of trainable parameters,  $N_{\text{param}}$ . For ONNs the data points follow a sublinear scaling law,  $P \propto N_{\text{param}}^{0.66}$ . This behavior stems from the geometric nature of the network: the parameter count grows with the “volume” of the network ( $\sim L \cdot d_{\text{model}}^2$ ), whereas the required optical power scales with the number of inputs ( $\sim L \cdot d_{\text{model}} \cdot L_{\text{seq}}$ ). Such distinction leads to a scaling exponent smaller than one, consistent with known results [70]. In contrast, the power consumption of electronic processors typically scales linearly with  $N_{\text{param}}$ . Consequently, the energy advantage of the optical approach becomes increasingly pronounced as AI models

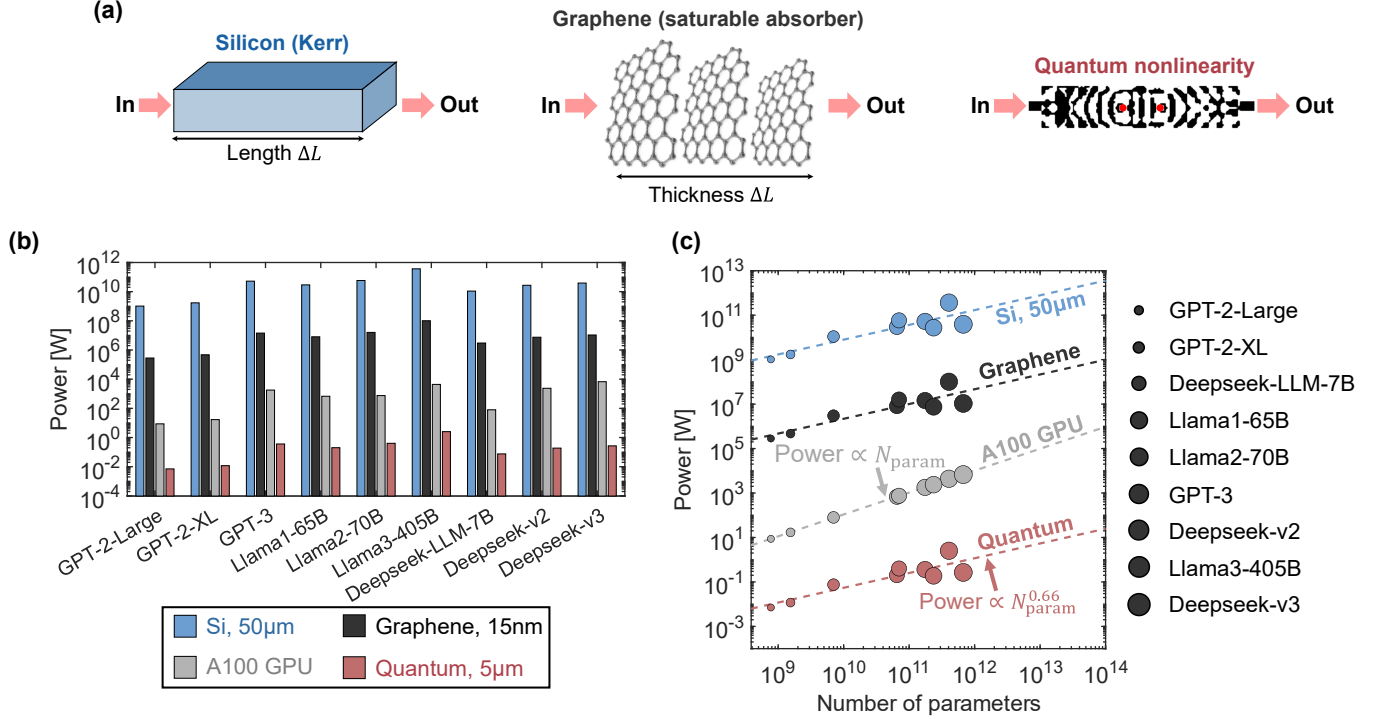


FIG. 4. **Scalability and energy advantage for large-scale optical models.** (a) Schematic illustration of nonlinear platforms: conventional bulk material (silicon, 50  $\mu$ m length), 2D saturable absorber (stacked graphene, 15 nm thickness), and the proposed quantum activation unit. (b) Histogram showing the estimated power consumption of all-optical LLMs. Conventional materials demand prohibitive power levels, while the proposed scheme enables sub-watt operation. The power consumption of running LLMs on NVIDIA A100 GPUs is shown for reference. (c) Estimated optical power versus model size  $N_{\text{param}}$ . ONNs follow a sublinear scaling  $P \propto N_{\text{param}}^{0.66}$ , indicating that optical computing has a growing advantage as models scale up.

continue to scale up. The above results identify a viable all-optical solution toward sustainable large-scale AI. By shifting to quantum nonlinearities, it should be possible to construct trillion-parameter all-optical models within a feasible power budget.

### III. DISCUSSION

In summary, we have presented a comprehensive framework to address the nonlinearity bottleneck in optical computing. At the device level, by integrating quantum emitters with inverse-designed nanophotonic structures, strong nonlinearity can be realized at intensities below 1 nW/ $\mu\text{m}^2$ . This enables complex functionalities ranging from nonlinear classification to reinforcement learning, presenting a clear performance gap over linear ONNs. Moreover, we have developed a general theoretical framework to quantify the expressive power of arbitrary nonlinear physical systems, which in turn allows us to determine the light intensity requirements. At the system level, we estimate that this architecture can, in principle, reduce the optical power consumption of LLMs to a few watts, exhibiting a favorable sublinear scaling with model size. Together, these results suggest that the lack of nonlinearity is not a fundamental limit, but rather an engineering challenge that could be overcome with quantum technologies. Despite these advances, transforming our theoretical proposal into large-scale hardware is still facing several practical challenges. A key trade-off exists between intensity threshold and operation bandwidth: the high sensitivity is inherently related to the emitter's long lifetime, which limits the response speed to GHz range. The limited bandwidth can be improved via Purcell enhancement inside optimized photonic structures [71]. Another challenge is the inhomogeneity of solid-state emitters. To tackle this issue, solutions such as DC Stark tuning have been demonstrated to tune the resonance of individual emitters [72]. Regarding integration, while deterministic placement of emitters remains difficult, recent advances in fabrication techniques offer promising solutions for large-scale integration [73, 74]. Finally, solid-state quantum emitters often require cryogenic operation to suppress dephasing and to approach lifetime-limited linewidths [75, 76], which introduces an additional power overhead for cooling.

Looking forward, this work shows that in order to unlock the full potential of optical computing, we should exploit the strong nonlinearity provided by quantum emitters rather than conventional bulk materials. By combining inverse-designed nanophotonics with modern deep learning theory, we provide a path toward all-optical deep learning that is both computationally powerful and more energy efficient than its electronic counterparts. Realizing this vision will ultimately pave the way toward sustainable, next-generation artificial intelligence.

## ACKNOWLEDGMENTS

The authors would like to thank Erfan Khoram, Zhicheng Wu, and Prof. E. Sifakis for insightful discussions.

- 
- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, Imagenet classification with deep convolutional neural networks, *Advances in neural information processing systems* **25** (2012).
  - [2] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, *et al.*, Human-level control through deep reinforcement learning, *nature* **518**, 529 (2015).
  - [3] J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K. Tunyasuvunakool, R. Bates, A. Žídek, A. Potapenko, *et al.*, Highly accurate protein structure prediction with alphafold, *nature* **596**, 583 (2021).
  - [4] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, Attention is all you need, *Advances in neural information processing systems* **30** (2017).
  - [5] D. Patterson, J. Gonzalez, Q. Le, C. Liang, L.-M. Munguia, D. Rothchild, D. So, M. Texier, and J. Dean, Carbon emissions and large neural network training, *arXiv preprint arXiv:2104.10350* (2021).
  - [6] Y. Shen, N. C. Harris, S. Skirlo, M. Prabhu, T. Baehr-Jones, M. Hochberg, X. Sun, S. Zhao, H. Larochelle, D. Englund, *et al.*, Deep learning with coherent nanophotonic circuits, *Nature photonics* **11**, 441 (2017).
  - [7] X. Lin, Y. Rivenson, N. T. Yardimci, M. Veli, Y. Luo, M. Jarrahi, and A. Ozcan, All-optical machine learning using diffractive deep neural networks, *Science* **361**, 1004 (2018).
  - [8] S. Shekhar, W. Bogaerts, L. Chrostowski, J. E. Bowers, M. Hochberg, R. Soref, and B. J. Shastri, Roadmapping the next generation of silicon photonics, *Nature Communications* **15**, 751 (2024).
  - [9] G. Wetzstein, A. Ozcan, S. Gigan, S. Fan, D. Englund, M. Soljačić, C. Denz, D. A. Miller, and D. Psaltis, Inference in artificial intelligence with deep optics and photonics, *Nature* **588**, 39 (2020).
  - [10] T. Wang, S.-Y. Ma, L. G. Wright, T. Onodera, B. C. Richard, and P. L. McMahon, An optical neural network using less than 1 photon per multiplication, *Nature Communications* **13**, 123 (2022).
  - [11] W. Shi, Z. Huang, T. Fu, and H. Chen, Review of nonlinear activation functions in optical neural networks, *Advanced Photonics* **7**, 064004 (2025).
  - [12] R. W. Boyd, A. L. Gaeta, and E. Giese, Nonlinear optics, in *Springer Handbook of Atomic, Molecular, and Optical Physics* (Springer, 2008) pp. 1097–1110.
  - [13] J. Feldmann, N. Youngblood, C. D. Wright, H. Bhaskaran, and W. H. Pernice, All-optical spiking neurosynaptic networks with self-learning capabilities, *Nature* **569**, 208 (2019).
  - [14] B. Wu, H. Li, W. Tong, J. Dong, and X. Zhang, Low-threshold all-optical nonlinear activation function based on a ge/si hybrid structure in a microring resonator, *Optical Materials Express* **12**, 970 (2022).
  - [15] T. Wu, Y. Li, L. Ge, and L. Feng, Field-programmable photonic nonlinearity, *Nature Photonics* , 1 (2025).
  - [16] A. Jha, C. Huang, and P. R. Prucnal, Reconfigurable all-optical nonlinear activation functions for neuromorphic photonics, *Optics letters* **45**, 4819 (2020).
  - [17] R. Yanagimoto, B. A. Ash, M. M. Sohoni, M. M. Stein, Y. Zhao, F. Presutti, M. Jankowski, L. G. Wright, T. Onodera, and P. L. McMahon, Programmable on-chip nonlinear photonics, *Nature* , 1 (2025).
  - [18] I. A. Williamson, T. W. Hughes, M. Minkov, B. Bartlett, S. Pai, and S. Fan, Reprogrammable electro-optic nonlinear activation functions for optical neural networks, *IEEE Journal of Selected Topics in Quantum Electronics* **26**, 1 (2019).
  - [19] M. M. Pour Fard, I. A. Williamson, M. Edwards, K. Liu, S. Pai, B. Bartlett, M. Minkov, T. W. Hughes, S. Fan, and T.-A. Nguyen, Experimental realization of arbitrary activation functions for optical neural networks, *Optics Express* **28**, 12138 (2020).
  - [20] M. Yildirim, N. U. Dinc, I. Oguz, D. Psaltis, and C. Moser, Nonlinear processing with linear optics, *Nature Photonics* **18**, 1076 (2024).
  - [21] F. Xia, K. Kim, Y. Eliezer, S. Han, L. Shaughnessy, S. Gigan, and H. Cao, Nonlinear optical encoding enabled by recurrent linear scattering, *Nature Photonics* **18**, 1067 (2024).
  - [22] C. C. Wanjura and F. Marquardt, Fully nonlinear neuromorphic computing with linear wave scattering, *Nature Physics* **20**, 1434 (2024).
  - [23] P. Lodahl, S. Mahmoodian, and S. Stobbe, Interfacing single photons and single quantum dots with photonic nanostructures, *Reviews of Modern Physics* **87**, 347 (2015).
  - [24] C. Zhu, T. Wang, P. L. McMahon, and D. Soh, Quantum optical neural networks using atom-cavity interactions to provide all-optical nonlinearity, *arXiv preprint arXiv:2511.06167* (2025).



- [25] R. Canora, X. Xu, Z. Niu, H. Alaeian, and S. Du, Engineering nonlinear activation functions for all-optical neural networks via quantum interference, arXiv preprint arXiv:2504.04009 (2025).
- [26] Y. Zuo, B. Li, Y. Zhao, Y. Jiang, Y.-C. Chen, P. Chen, G.-B. Jo, J. Liu, and S. Du, All-optical neural network with nonlinear activation functions, *Optica* **6**, 1132 (2019).
- [27] A. Ryou, J. Whitehead, M. Zhelyeznyakov, P. Anderson, C. Keskin, M. Bajcsy, and A. Majumdar, Free-space optical neural network based on thermal atomic nonlinearity, *Photonics Research* **9**, B128 (2021).
- [28] Z. Huang, W. Shi, S. Wu, Y. Wang, S. Yang, and H. Chen, Pre-sensor computing with compact multilayer optical neural network, *Science Advances* **10**, eado8516 (2024).
- [29] K. Ohno, F. J. Heremans, L. C. Bassett, B. A. Myers, D. M. Toyli, A. C. B. Jayich, C. J. Palmstrøm, and D. D. Awschalom, Engineering shallow spins in diamond with nitrogen delta-doping, *Applied Physics Letters* **101**, 082413 (2012).
- [30] Y.-C. Chen, B. Griffiths, L. Weng, S. S. Nicley, S. N. Ishmael, Y. Lekhai, S. Johnson, C. J. Stephen, B. L. Green, G. W. Morley, *et al.*, Laser writing of individual nitrogen-vacancy defects in diamond with near-unity yield, *Optica* **6**, 662 (2019).
- [31] A. M. Day, J. R. Dietz, M. Sutula, M. Yeh, and E. L. Hu, Laser writing of spin defects in nanophotonic cavities, *Nature Materials* **22**, 696 (2023).
- [32] C. M. Lalau-Keraly, S. Bhargava, O. D. Miller, and E. Yablonovitch, Adjoint shape optimization applied to electromagnetic design, *Optics express* **21**, 21693 (2013).
- [33] O. Kulce, D. Mengü, Y. Rivenson, and A. Ozcan, All-optical information-processing capacity of diffractive surfaces, *Light: Science & Applications* **10**, 25 (2021).
- [34] D. A. Miller, Why optics needs thickness, *Science* **379**, 41 (2023).
- [35] Y. Li and F. Monticone, The spatial complexity of optical computing: toward space-efficient design, *Nature Communications* **16**, 8588 (2025).
- [36] T. Onodera, M. M. Stein, B. A. Ash, M. M. Sohoni, M. Bosch, R. Yanagimoto, M. Jankowski, T. P. McKenna, T. Wang, G. Shvets, *et al.*, Arbitrary control over multimode wave propagation for machine learning, *Nature Physics* , 1 (2025).
- [37] S. Yu, X. Piao, and N. Park, Nonlinear unitary circuits for photonic neural networks, *ACS Photonics* (2025).
- [38] G. Montúfar, R. Pascanu, K. Cho, and Y. Bengio, On the number of linear regions of deep neural networks, *Advances in neural information processing systems* **27** (2014).
- [39] B. Poole, S. Lahiri, M. Raghu, J. Sohl-Dickstein, and S. Ganguli, Exponential expressivity in deep neural networks through transient chaos, *Advances in neural information processing systems* **29** (2016).
- [40] M. Raghu, B. Poole, J. Kleinberg, S. Ganguli, and J. Sohl-Dickstein, On the expressive power of deep neural networks, in *international conference on machine learning* (PMLR, 2017) pp. 2847–2854.
- [41] Y. Shi, J. Ren, G. Chen, W. Liu, C. Jin, X. Guo, Y. Yu, and X. Zhang, Nonlinear germanium-silicon photodiode for activation and monitoring in photonic neuromorphic networks, *Nature Communications* **13**, 6048 (2022).
- [42] Q. Bao, H. Zhang, Y. Wang, Z. Ni, Y. Yan, Z. X. Shen, K. P. Loh, and D. Y. Tang, Atomic-layer graphene as a saturable absorber for ultrafast pulsed lasers, *Advanced Functional Materials* **19**, 3077 (2009).
- [43] J. Shi, P. Yu, F. Liu, P. He, R. Wang, L. Qin, J. Zhou, X. Li, J. Zhou, X. Sui, *et al.*, 3r mos2 with broken inversion symmetry: a promising ultrathin nonlinear optical device, *Advanced Materials* **29**, 1701486 (2017).
- [44] B. Liu, K. Liang, Q. Zhou, A. R. Khan, Z. Lu, T. Yildirim, X. Sun, S. Rahman, Y. Liu, Z. Yu, *et al.*, Giant second harmonic generation in two-dimensional tellurene with synthesis and thickness engineering, *Applied physics reviews* **12** (2025).
- [45] E. Hanamura, Rapid radiative decay and enhanced optical nonlinearity of excitons in a quantum well, *Physical Review B* **38**, 1228 (1988).
- [46] A. Javadi, I. Söllner, M. Arcari, S. L. Hansen, L. Midolo, S. Mahmoodian, G. Kiršanskė, T. Pregnolato, E. Lee, J. Song, *et al.*, Single-photon non-linear optics with a quantum dot in a waveguide, *Nature communications* **6**, 8655 (2015).
- [47] J. Volz, M. Scheucher, C. Junge, and A. Rauschenbeutel, Nonlinear  $\pi$  phase shift for single fibre-guided photons interacting with a single resonator-enhanced atom, *Nature Photonics* **8**, 965 (2014).
- [48] I. Shomroni, S. Rosenblum, Y. Lovsky, O. Bechler, G. Guendelman, and B. Dayan, All-optical routing of single photons by a one-atom switch controlled by a single photon, *Science* **345**, 903 (2014).
- [49] B. Hacker, S. Welte, G. Rempe, and S. Ritter, A photon–photon quantum gate based on a single atom in an optical resonator, *Nature* **536**, 193 (2016).
- [50] D. M. Lukin, C. Dory, M. A. Guidry, K. Y. Yang, S. D. Mishra, R. Trivedi, M. Radulaski, S. Sun, D. Vercruysse, G. H. Ahn, *et al.*, 4h-silicon-carbide-on-insulator for integrated quantum and nonlinear photonics, *Nature Photonics* **14**, 330 (2020).
- [51] Q. Zhou, S. Gangaraj, M. Zhou, and Z. Yu, Simulating quantum emitters in arbitrary photonic environments using fdtd: beyond the semi-classical regime, arXiv preprint arXiv:2410.16118 (2024).
- [52] H. Wang and S. Fan, Lorentz–drude dipoles in the radiative limit and their modeling in finite-difference time-domain methods, *Annalen der Physik* , e00156 (2025).
- [53] V. Nikkhah, A. Pirmoradi, F. Ashtiani, B. Edwards, F. Aflatouni, and N. Engheta, Inverse-designed low-index-contrast structures on a silicon photonics platform for vector–matrix multiplication, *Nature Photonics* **18**, 501 (2024).
- [54] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, *et al.*, Mastering the game of go with deep neural networks and tree search, *nature* **529**, 484 (2016).
- [55] J. Hwangbo, J. Lee, A. Dosovitskiy, D. Bellicoso, V. Tsounis, V. Koltun, and M. Hutter, Learning agile and dynamic motor skills for legged robots, *Science Robotics* **4**, eaau5872 (2019).
- [56] M. Towers, A. Kwiatkowski, J. Terry, J. U. Balis, G. De Cola, T. Deleu, M. Goulão, A. Kallinteris, M. Krimmel, A. KG,

- et al.*, Gymnasium: A standard interface for reinforcement learning environments, arXiv preprint arXiv:2407.17032 (2024).
- [57] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, Proximal policy optimization algorithms, arXiv preprint arXiv:1707.06347 (2017).
  - [58] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor, in *International conference on machine learning* (Pmlr, 2018) pp. 1861–1870.
  - [59] M. Dinu, F. Quochi, and H. Garcia, Third-order nonlinearities in silicon at telecom wavelengths, *Applied physics letters* **82**, 2954 (2003).
  - [60] E. Dulkeith, Y. A. Vlasov, X. Chen, N. C. Panoiu, and R. M. Osgood Jr, Self-phase-modulation in submicron silicon-on-insulator photonic wires, *Optics express* **14**, 5524 (2006).
  - [61] K. Y. Lau, X. Liu, and J. Qiu, A comparison for saturable absorbers: Carbon nanotube versus graphene, *Advanced Photonics Research* **3**, 2200023 (2022).
  - [62] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, *et al.*, Language models are unsupervised multitask learners, *OpenAI blog* **1**, 9 (2019).
  - [63] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, *et al.*, Language models are few-shot learners, *Advances in neural information processing systems* **33**, 1877 (2020).
  - [64] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, *et al.*, Llama: Open and efficient foundation language models, arXiv preprint arXiv:2302.13971 (2023).
  - [65] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, *et al.*, Llama 2: Open foundation and fine-tuned chat models, arXiv preprint arXiv:2307.09288 (2023).
  - [66] A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan, *et al.*, The llama 3 herd of models, arXiv e-prints, arXiv (2024).
  - [67] X. Bi, D. Chen, G. Chen, S. Chen, D. Dai, C. Deng, H. Ding, K. Dong, Q. Du, Z. Fu, *et al.*, Deepseek llm: Scaling open-source language models with longtermism, arXiv preprint arXiv:2401.02954 (2024).
  - [68] A. Liu, B. Feng, B. Wang, B. Wang, B. Liu, C. Zhao, C. Deng, C. Ruan, D. Dai, D. Guo, *et al.*, Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model, arXiv preprint arXiv:2405.04434 (2024).
  - [69] A. Liu, B. Feng, B. Xue, B. Wang, B. Wu, C. Lu, C. Zhao, C. Deng, C. Zhang, C. Ruan, *et al.*, Deepseek-v3 technical report, arXiv preprint arXiv:2412.19437 (2024).
  - [70] M. Anderson, S.-Y. Ma, T. Wang, L. Wright, and P. McMahon, Optical transformers, *Transactions on Machine Learning Research* (2023).
  - [71] D. Englund, D. Fattal, E. Waks, G. Solomon, B. Zhang, T. Nakaoka, Y. Arakawa, Y. Yamamoto, and J. Vučković, Controlling the spontaneous emission rate of single quantum dots in a two-dimensional photonic crystal, *Physical review letters* **95**, 013904 (2005).
  - [72] A. Laucht, J. Villas-Bôas, S. Stobbe, N. Hauke, F. Hofbauer, G. Böhm, P. Lodahl, M.-C. Amann, M. Kaniber, and J. Finley, Mutual coupling of two semiconductor quantum dots via an optical nanocavity, *Physical Review B—Condensed Matter and Materials Physics* **82**, 075305 (2010).
  - [73] Y.-C. Chen, P. S. Salter, S. Knauer, L. Weng, A. C. Frangeskou, C. J. Stephen, S. N. Ishmael, P. R. Dolan, S. Johnson, B. L. Green, *et al.*, Laser writing of coherent colour centres in diamond, *Nature Photonics* **11**, 77 (2017).
  - [74] P. Laferrière, E. Yeung, I. Miron, D. B. Northeast, S. Haffouz, J. Lapointe, M. Korkusinski, P. J. Poole, R. L. Williams, and D. Dalacu, Unity yield of deterministically positioned quantum dot single photon sources, *Scientific Reports* **12**, 6376 (2022).
  - [75] A. Sipahigil, K. D. Jahnke, L. J. Rogers, T. Teraji, J. Isoya, A. S. Zibrov, F. Jelezko, and M. D. Lukin, Indistinguishable photons from separated silicon-vacancy centers in diamond, *Physical Review Letters* **113**, 113602 (2014).
  - [76] L. J. Rogers, K. D. Jahnke, T. Teraji, L. Marseglia, C. Müller, B. Naydenov, H. Schauffert, C. Kranz, J. Isoya, L. P. McGuinness, and F. Jelezko, Multiple intrinsically identical single-photon emitters in the solid state, *Nature Communications* **5**, 4739 (2014).