

# Nanophotonic Media for Artificial Neural Inference

ERFAN KHORAM,<sup>1</sup> ANG CHEN,<sup>1</sup>, DIANJING LIU<sup>1</sup>, LEI YING<sup>1</sup>, QIQI WANG<sup>2</sup>, MING YUAN<sup>3</sup> ZONGFU YU<sup>1,\*</sup>

<sup>1</sup>*Department Electrical And Computer Engineering, University of Wisconsin Madison-Madison, WI53706, USA*

<sup>2</sup>*Department of Aeronautics and Astronautics, Massachusetts Institute of Technology, Cambridge, MA02139, USA*

<sup>3</sup>*Department of Statistics, Columbia University, New York, NY10027, USA*

\**zyu54@wisc.edu*

**Abstract:** We show optical waves passing through a nanophotonic medium can perform artificial neural computing. Complex information, is encoded in the wave front of an input light. The medium transforms the wave front to realize sophisticated computing tasks such as image recognition. At the output, the optical energy is concentrated to well-defined locations, which for example can be interpreted as the identity of the object in the image. These computing media can be as small as tens of wavelengths and offer ultra-high computing density. They exploit sub-wavelength scatterers to realize complex input output mapping beyond the capabilities of traditional nanophotonic devices.

© 2019 Optical Society of America under the terms of the [OSA Open Access Publishing Agreement](#)

## 1. Introduction

Artificial neural networks (ANN) have shown exciting potential in a wide range of applications, but they also require ever-increasing computing power. This has prompted an effort to search for alternative computing methods that are faster and more energy efficient. One interesting approach is optical neural computing [1–7]. This analog computing method can be passive, with minimal energy consumption, and more importantly, its intrinsic parallelism can greatly accelerate computing speed.

Most optical neural computing follow the architecture of digital ANNs, using a layered feed-forward network as shown in Fig. 1a. Free-space diffraction [4,8] or integrated waveguides [1,3,9] are used as the connections between layered activation units. Similar digital signals in ANN, optical signals pass through optical networks in the forward direction once (light reflection propagating in the backward direction is avoided or neglected). However, it is the reflection that provides the feedback mechanism which gives rise to rich wave physics. It holds the key to the miniaturization of optical devices such as laser cavities [10], photonic crystals [11], meta-materials [12], and ultra-compact beam splitters [13–15]. Here we show that by leveraging optical reflection, it is possible to go beyond the paradigm of layered feed-forward networks to realize artificial neural computing in a continuous and layer-free fashion. Fig.1b shows the proposed Nanophotonic Neural Medium(NNM). An optical signal enters from the left and the output is the energy distribution on the right side of the medium. Computation is performed by a host material, such as  $SiO_2$ , with numerous inclusions. The inclusions can be air holes, or any other material with an index different from that of the host medium. These inclusions strongly scatter light in both the forward and backward directions. The scattering spatially mixes the input light, rendering it a counterpart to linear matrix multiplication (Fig.1c) in a digital ANN. The locations and shapes of inclusions are the equivalent of weight parameters in digital ANNs, and their sizes are typically sub-wavelength. The nonlinear operation can be realized via inclusions made of dye semiconductor or graphene saturable absorbers, where they perform distributed

nonlinear activation. These nonlinearities are designed with the Rectified Linear Units (ReLU) in mind [16](Fig.1d).

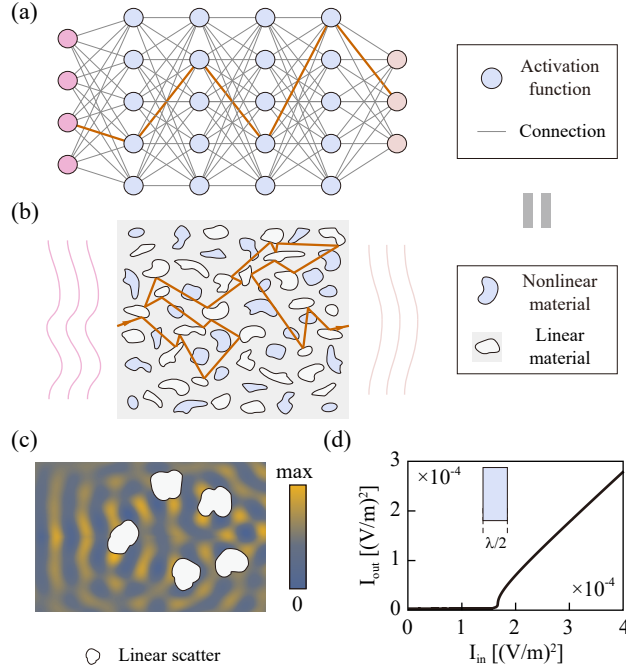


Fig. 1. (a) A conventional ANN architecture where the information propagates only in the forward direction (green arrow). (b) Proposed Nanophotonic Neural Medium (NNM). Passive neural computing is performed by light passing through nanostructured medium with both linear and nonlinear scatterers. (c) Full-wave simulation of light scattered by nanostructures, which spatially redistribute the optical energy to different directions. (d) The output intensity of light with wavelength  $\lambda$ , passing through the designed nonlinear material with a thickness of  $\lambda/2$ . It is a nonlinear function of the incident wave intensity. This material is used as nonlinear activation as indicated by light blue color.

Fig.2 shows a NNM in action, where a two-dimensional (2D) medium is trained to recognize gray-scale handwritten digits. The dataset contains 5,000 different images, representative ones of which are shown in Fig.2a. Each time, one image, represented by  $20 \times 20$  pixels, is converted to a vector, and then encoded as the spatial intensity of input light incident on the left. Inside the NNM, nanostructures create strong interferences and light is guided toward one of ten output locations depending on the digit that the image represents, where the output with the highest share of energy intensity is categorized as the inferred class. Fig.2b shows the fields created by two different hand-written 2 digits. Because of different shapes, the field patterns created by these two images are quite different but both lead to the same hot spot at the output, which correctly identifies the identity information as the number 2. As another example, Fig.2c shows the case of two hand-written 8 digits that result in another hot spot. Here, the field is simulated by solving a nonlinear wave equation using Finite-Difference Frequency-Domain (FDFD [17]) method. The size of the NNM is  $80\lambda$  by  $20\lambda$ , where  $\lambda$  is the wavelength of light used to carry and process the information. The average recognition accuracy reaches over 79% for a test set made up of 1,000 images. The limited reported accuracy is due to the heavy constraints we set during the optimization for fabrication concerns. These constraints keep the medium dense, where it would have been otherwise made up of sparse sections of air and  $SiO_2$ . By relaxing these requirements or using larger medium sizes, the accuracy can be further improved.

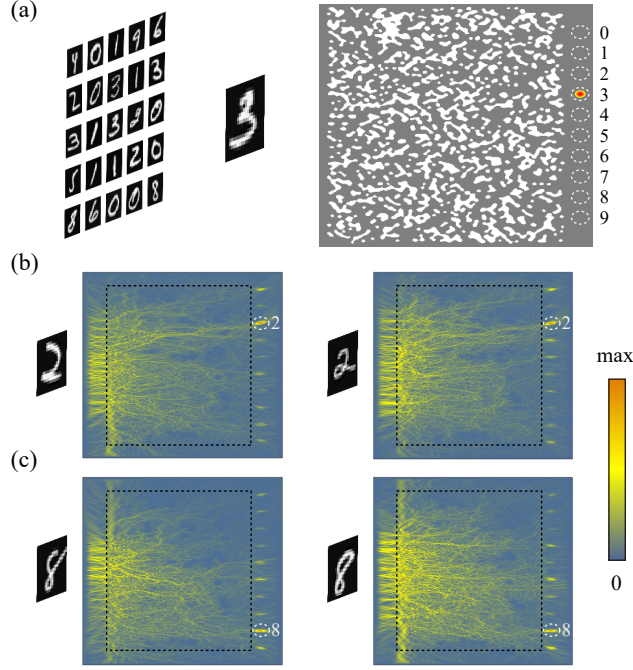


Fig. 2. (a) NNM trained to recognize handwritten digits. The input wave encodes the image as the intensity distribution. On the right side of NNM, the optical energy concentrate to different locations depending on the image's classification labels. (b) Two samples of the digit 2 and their optical fields inside NNM. As it can be seen, although the field distributions differ for the images of the same digit, they are classified as the same digit. (c) the same as (b) but for two samples of the digit 8. Also, in both (b) and (c), the boundaries of the trained medium have been shown with black borderlines.(see Visualization 1)

NNM can provide ultra-high computing density by tapping into sub-wavelength features. In theory, the number of weight parameters is infinite: every atom in this medium can be varied to influence the wave propagation. In practice, a change below 10 nm would be considered too challenging for fabrication. Even at this scale, the potential number of weights exceeds 10 billion parameters per square millimeter for a 2D implementation. This is much greater computing density than both free-space [8, 18] and on-chip optical neural networks [1, 3]. In addition, NNM has a few other attractive features. It has stronger expressive power than layered optical networks. In fact, layered networks are a subset of NNM, as a medium can be shaped into connected waveguides as a layered network. Furthermore, it does not have the issue of diminishing gradients in deep neural networks. Maxwell's equations, as the governing principle, guarantee that the underlying linear operation is always unitary, which does not have diminishing or exploding gradients [19]. Lastly, NNM does not have to follow any specific geometry, and thus it can be easily shaped and integrated into existing vision or communication devices as the first step of optical preprocessing.

We now discuss the training of NNM. Although, one could envision in-situ training of NNM using tunable optical materials [3], here we focus on training in the digital domain and use NNM only for inference. The underlying dynamics of the NNM are governed by the nonlinear Maxwell's equations, which, in the frequency domain, can be written as

$$L_{(r,E(r))}E(r) = -i\omega J(r) \quad (1)$$

where  $L_{(r,E(r))} = (\nabla \times \nabla \times) / \mu - \omega^2 \varepsilon_{(r,E(r))}$ , and  $\mu$  and  $\varepsilon$  are the permeability and permittivity.  $J$  is the current source density which represents the spatial profile of the input light and is only non-zero on the left side of the medium. Waveguide modes or plane waves can also be used as the input, which are also implemented as current sources in numerical simulation. For a classification problem, the probability of the  $i^{th}$  class label is given by  $h_i = (\int_{dr} |E(r)|^2 R_i(r)) / (\sum_{i=1}^{10} \int_{dr} |E(r)|^2 R_i(r))$ , which represents the percentage of energy at the  $i^{th}$  receiver relative to the total optical energy that reaches all receivers. Here the profile function  $R_i(r)$  defines the location of receivers, and is only non-zero at the position of the  $i^{th}$  receiver. The training is performed by optimizing the dielectric constant  $\varepsilon(r, E)$  similar to how weight parameters are trained in traditional neural networks. The cost function  $C$  is defined by the cross entropy between the output vector  $\mathbf{h}$  and the ground truth  $\mathbf{y}$ .

$$C = - \sum_{i=1}^{10} y_i \log(h_i) + (1 - y_i) \log(1 - h_i) \quad (2)$$

The ground truth  $\mathbf{y}$  is a one-hot vector. Digit 8 is represented as  $\mathbf{y} = (0, 0, 0, 0, 0, 0, 0, 0, 1, 0)$ , for instance. The gradient of the cost function with respect to the dielectric constant  $\varepsilon$  can be calculated point by point. For example, one could assess the effect of changing  $\varepsilon$  at one spatial point; the change is only kept if the loss function decreases. This method has achieved remarkable success in simple photonic devices [13]. However, each gradient calculation requires solving full-wave nonlinear Maxwell's equations. It is prohibitively costly for NNM, which could easily have millions of gradients. Here, we use Adjoint State Method(ASM) to compute all gradients in one step:

$$\frac{dC}{d\varepsilon(r)} = -2\omega^2 \text{Real}\{\lambda(r)E(r)\} \quad (3)$$

Here  $\lambda(r)$  is a Lagrangian multiplier, which is the solution to the adjoint equation (Eq.4), in which the electric field  $E(r)$  is obtained by solving Eq.1. The adjoint equation here is slightly more involved than what is generally used in inverse design, and this is due to the fact that nonlinear behavior is included in our dynamics. A similar derivation for nonlinear adjoint equation is done in [20].

$$\frac{\partial C}{\partial E(r)} + \lambda(r)(L_{(r,E(r))} + \frac{\partial L_{(r,E(r))}}{\partial E(r)} E(r)) + \overline{\lambda(r)}(\frac{\partial \overline{L_{(r,E(r))}}}{\partial E(r)} \overline{E(r)}) = 0 \quad (4)$$

The training process, as illustrated in Fig.3a, minimizes the summation of the cost functions  $C$  for all training instances through Stochastic Gradient Descent(SGD). The process starts with one input image as the light source, for which we solve the nonlinear Maxwell's equations in an iterative process (pink block in Fig.3a). The initial field is set to be random  $E_0(r)$ , which allows us to calculate the dielectric constant  $\varepsilon(r, E_0(r))$ . Then FDFD simulation is used to solve Eq.1, and the resulting electric field  $E_1(r)$  is then used to update the dielectric constant. This iteration continues until the field converges. The next step is to compute the gradient based on Eq.3. Once the structural change is updated, the training of this instance is finished.

The above process is repeated again, but for the next image in the training queue, instead of the same image. This gradient descent process is stochastic, which is quite different from typical use of ASM in nanophotonics [14, 15] where gradient descent is performed repeatedly for very few inputs until the loss function converges. In these traditional optimizations, the device needs to only function for those few specific inputs. If such processes were used here, the medium would do extremely well for particular images but fail to generalize and recognize other images.

The gradient descent process treats the dielectric constant as a continuous variable, but in practice, its value is discrete, depending on the material used at the location. For example, in the

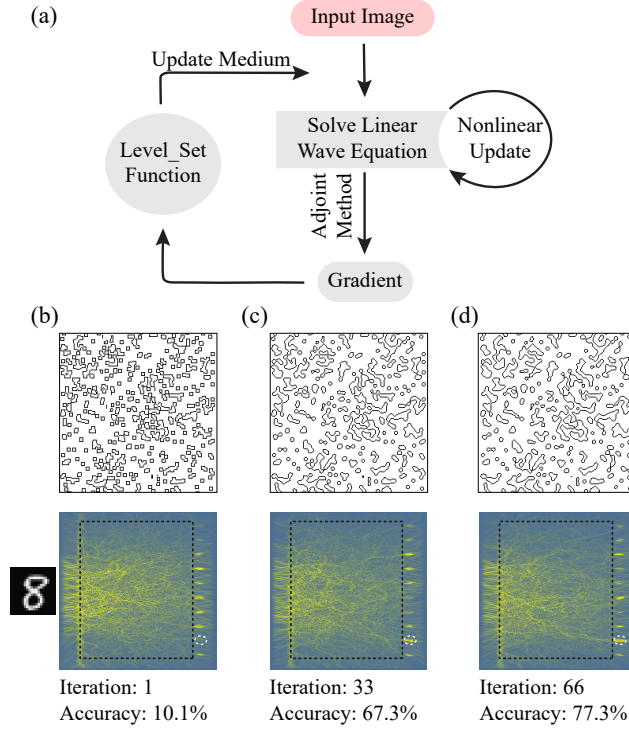


Fig. 3. (a) The training starts by encoding an image as a vector of current source densities in the FFD simulation. This step is followed by an iterative process to solve for the electric field in a nonlinear medium. Next we use the adjoint state method to calculate the gradient, which is then used to update the level-set function and consequently, the medium itself. In batch training, we sum the cost functions calculated for different images in the same batch and compute the gradients. (b)-(d) show an NNM in training after 1, 33, and 66 training iterations respectively. At each step, the boundary between the host material and the inclusions is shown, along with the field distribution for the same randomly selected digit 8. Also the accuracy of the medium on the test set can be seen for that particular stage in training.

case of a medium with  $SiO_2$  host material and linear *air* inclusions, the dielectric constants can either be 2.16 or 1. Discrete variables remain effective for neural computing [21]. Here, we need to take special care to further constrain the optimization process. This is done by using a level set function [22], where each of the two materials (host material and the linear inclusion material), is assigned to each of the two levels in the level-set function  $\phi(r)$  similar to ref [14, 15].

$$\varepsilon(r) = \begin{cases} \varepsilon_{SiO_2} & \phi(r) < 0 \\ \varepsilon_{Air} & \phi(r) > 0 \end{cases} \quad (5)$$

The training starts with randomly distributed inclusions, both linear and nonlinear, throughout the host medium. The boundaries between two materials evolve in the training. Specifically, the level set function is updated by  $-v(r)|\nabla\phi|$ , where  $v(r)$  is the gradient calculated by ASM and  $|\nabla\phi|$  indicates the boundary between the two constituent materials. Therefore, at each step, this method essentially decides whether any point on the boundary should be switched from one material to the other. Nonlinear sections perform the activation function and their location and shape are fixed in this optimization. They could also be optimized, which would be equivalent to

optimizing structural hyperparameters in layered neural networks [23–25].

As a specific example, we now discuss the training of the 2D medium shown in Fig.2. The structural evolution is shown in Fig.3b-d during the training. We start by randomly seeding the domain with dense but small inclusions. As the training progresses, the inclusions move and merge, eventually converging. The recognition accuracy for both training and test group improve during this process.

Next, we show another example based on a three-dimensional(3D) medium, whose size is  $4\lambda \times 4\lambda \times 6\lambda$ . The inputs can be an image projected on the top surface of the medium. For example, we use a plane wave to illuminate a mask with its opening shaped into a hand-written digit as shown in Fig.4(Movie S2 shows how the energy distribution on the output evolves as a handwritten digit gradually emerges as the input). Fabricating 3D inclusions is generally difficult, but it is much easier to tune the permittivity of materials using direct laser writing [26]. Thus, here we allow the dielectric constant to vary continuously. To save on computational resources, we allow 5% variation. In experimental realization, a smaller variation range can always be compensated for by using larger media. The 3D trained NNM had an accuracy of about 84% for the test set; the confusion matrix is shown in Fig.4b. The better performance in comparison with the 2D implementation is due to a higher degree of freedom we allow the dielectric constant to have.

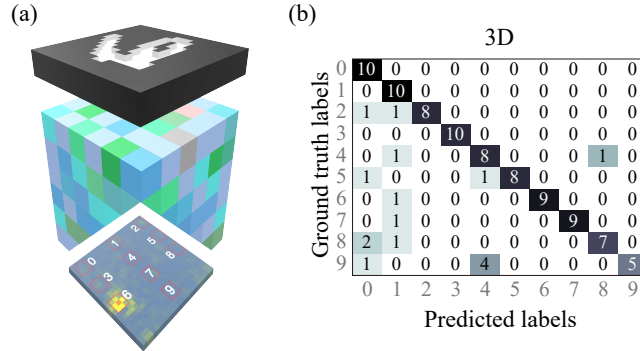


Fig. 4. (a) 3D nanophotonic neural medium case. Different colors illustrate varying values of permittivity. The input image is projected onto the top surface. Computing is performed while the wave propagates through the 3D medium. The field distribution on the bottom surface is used to recognize the image. Full-wave simulation shows the optical energy is concentrated on the location with correct class label, in this case 6.(b) The confusion matrix. The rows on the matrix show true labels of the images that have been presented as input, and the columns depict the labels that the medium has classified each input. Therefore the diagonal elements show the number of correct classifications out of every ten samples.(see Visualization 2)

## 2. Supplemental Materials

### 2.1. Gradient Derivation

In this section we will go through the derivation of the gradient of the cost function for the case of a classification problem with mini-batch gradient descent. We start with the the electromagnetic wave equation in the frequency domain:

$$\begin{aligned}
 L_{(r,\omega)}E_{(r)} &= -i\omega J_{(r)} \\
 L_{(r,\omega)} &= (\nabla \times \left(\frac{1}{\mu}\right) \times \nabla \times) - \omega^2 \epsilon_{(r)}
 \end{aligned} \tag{S1}$$

Here  $J$  is the current source density which corresponds to the values of the image pixels in our work. This equation is written for linear materials; For nonlinear materials,  $\varepsilon$  becomes a function of the electric field at that point. We will elaborate upon this point as we derive the gradient.

The receivers on the output side of the structure measure the electromagnetic energy in their position. To model the light absorption in receivers, we use a Gaussian function to define the profile of the location of the  $i^{th}$  output.

$$R_{i(r)} = e^{-\frac{\|r-r_i^{rec}\|^2}{2\sigma^2}} \quad (S2)$$

where  $r_i^{rec}$  is the location of the  $i^{th}$  receiver while  $\sigma$  represents the spatial span of the receivers. The energy inside the receiver is

$$o_i = \frac{\varepsilon_{air}}{2} \int \|E(r)\|^2 R_{i(r)} dr \quad (S3)$$

In this equation,  $\varepsilon_{air}$  is set as the permittivity of air, as the receivers are located outside the medium. Then the cross entropy cost function for a batch of size  $m$  can be written as

$$C = -\frac{1}{m} \sum_{j=1}^m \sum_{i=1}^{10} y_i^j \log(h_i^j) + (1 - y_i^j) \log(1 - h_i^j) \quad (S4)$$

$$h_i^j = \frac{o_i^j}{\sum_{n=1}^{10} o_n^j}$$

where  $j$  represents the  $j^{th}$  image in the batch. From here, the derivation of the gradient can commence.

$$\frac{dC}{d\varepsilon(r)} = \sum_{j=1}^m \left( \int \frac{\partial C}{\partial E_{(r')}^j} \frac{\partial E_{(r')}^j}{\partial \varepsilon(r)} dr' + \int \frac{\partial C}{\partial \overline{E_{(r')}^j}} \frac{\partial \overline{E_{(r')}^j}}{\partial \varepsilon(r)} dr' \right) \quad (S5)$$

where  $\overline{E_{(r')}^j}$  is the conjugate of the electric field. Here directly computing the above gradient involves calculation of  $\partial E_{(r')}^j / \partial \varepsilon(r)$ , which is computationally expensive. In order to circumvent this problem, we define a Lagrangian as

$$Lg = C + \sum_{j=1}^m \int \lambda_{(r')}^j (L_{(r')} E_{(r')}^j + i\omega J^j) dr' + \sum_{j=1}^m \int \overline{\lambda_{(r')}^j (L_{(r')} E_{(r')}^j + i\omega J^j)} dr' \quad (S6)$$

Here  $\lambda_{(r')}^j$  is the Lagrange multiplier which is equivalent to *Adjoint Field* in similar works. This Lagrangian has the same gradient as the cost function because the second and third terms are zero. We now compute its gradient with respect to the real part of the permittivity at each spatial point:

$$\frac{dLg}{d\varepsilon(r)} = \int \sum_{j=1}^m \frac{\partial C}{\partial E_{(r')}^j} \frac{\partial E_{(r')}^j}{\partial \varepsilon(r)} dr' + \int \sum_{j=1}^m \frac{\partial C}{\partial \overline{E_{(r')}^j}} \frac{\partial \overline{E_{(r')}^j}}{\partial \varepsilon(r)} dr' +$$

$$\int \sum_{j=1}^m \lambda_{(r')}^j \left( \frac{dL_{(r')}}{d\varepsilon(r)} E_{(r')}^j + L_{(r')} \frac{\partial E_{(r')}^j}{\partial \varepsilon(r)} \right) dr' + \int \sum_{j=1}^m \overline{\lambda_{(r')}^j} \left( \frac{d\overline{L_{(r')}}}{d\varepsilon(r)} \overline{E_{(r')}^j} + \overline{L_{(r')}} \frac{\partial \overline{E_{(r')}^j}}{\partial \varepsilon(r)} \right) dr' \quad (S7)$$

In Eq.S7, since the dielectric constant of nonlinear materials depends on the field,  $dL_{(r')}/d\varepsilon_{(r)}$  can be calculated as

$$\frac{dL_{(r')}}{d\varepsilon_{(r)}} = \frac{\partial L_{(r')}}{\partial \varepsilon_{(r)}} + \frac{\partial L_{(r')}}{\partial E_{(r')}^j} \frac{\partial E_{(r')}^j}{\partial \varepsilon_{(r)}} + \overline{\frac{\partial L_{(r')}}{\partial E_{(r')}^j} \frac{\partial \overline{E_{(r')}^j}}{\partial \varepsilon_{(r)}}} \quad (\text{S8})$$

Similarly,  $d\overline{L_{(r')}}/d\varepsilon_{(r)}$  can be calculated too. In order to apply the adjoint state method, we first group all the terms multiplied by  $\partial E_{(r')}^j/\partial \varepsilon_{(r)}$  together (likewise for  $\partial \overline{E_{(r')}^j}/\partial \varepsilon_{(r)}$ ). This provides us with Eq.S9.

$$\begin{aligned} \frac{dLg}{d\varepsilon_{(r)}} = 2\text{Real}\{ & \int \sum_{j=1}^m \left( \frac{\partial C}{\partial E_{(r')}^j} + \lambda_{(r')}^j (L_{(r')} + E_{(r')}^j) \frac{\partial L_{(r')}}{\partial E_{(r')}^j} + \overline{\lambda_{(r')}^j (E_{(r')}^j) \frac{\partial \overline{L_{(r')}}}{\partial E_{(r')}^j}} \right) \frac{\partial E_{(r')}^j}{\partial \varepsilon_{(r)}} dr' \} + \\ & 2\text{Real}\{ \int \sum_{j=1}^m \lambda_{(r')}^j \frac{dL_{(r')}}{d\varepsilon_{(r)}} E_{(r')}^j dr' \} \end{aligned} \quad (\text{S9})$$

we set the Lagrange multiplier  $\lambda^j$  in a manner that forces the summation of all the terms multiplied by  $\partial E_{(r')}^j/\partial \varepsilon_{(r)}$  to be zero.  $\lambda^j$  can be calculated with the *Adjoint Equation* that is produced in the equation above by setting the first term in it to zero.

$$\frac{\partial C}{\partial E_{(r')}^j} + \lambda_{(r')}^j (L_{(r')} + E_{(r')}^j) \frac{\partial L_{(r')}}{\partial E_{(r')}^j} + \overline{\lambda_{(r')}^j (E_{(r')}^j) \frac{\partial \overline{L_{(r')}}}{\partial E_{(r')}^j}} = 0 \quad (\text{S10})$$

When the adjoint equation is satisfied, the gradient can be further simplified as :

$$\frac{dLg}{d\varepsilon_{(r)}} = 2\text{Real}\{ \int \sum_{j=1}^m \lambda_{(r')}^j \frac{dL_{(r')}}{d\varepsilon_{(r)}} E_{(r')}^j dr' \} \quad (\text{S11})$$

Using the definition of the operator  $L_{(r')}$  in Eq.S1, we get  $\frac{dL_{(r')}}{d\varepsilon_{(r)}} = -\omega^2 \frac{d\varepsilon_{(r')}}{d\varepsilon_{(r)}} = -\omega^2 \delta_{(r,r')}$ , which brings us to the following equation.

$$\frac{dLg}{d\varepsilon_{(r)}} = -2\omega^2 \text{Real}\{ \sum_{j=1}^m \lambda_{(r')}^j E_{(r')}^j \} \quad (\text{S12})$$

The next step in deriving the gradient is calculating the derivative of cost with respect to the electric field.

$$\frac{\partial C}{\partial E_{(r)}^j} = -\frac{1}{m} \sum_{i=1}^{10} y_i^j \frac{1}{h_i^j} \frac{\partial h_i^j}{\partial E_{(r)}^j} - (1 - y_i^j) \frac{1}{1 - h_i^j} \frac{\partial h_i^j}{\partial E_{(r)}^j} \quad (\text{S13})$$

The term  $\partial h_i^j/\partial E_{(r)}^j$  can be calculated based on the definition of  $h_i^j$  in Eq.S4.

$$\frac{\partial h_i^j}{\partial E_{(r)}^j} = \frac{\frac{\partial \sigma_i^j}{\partial E_{(r)}^j} (\sum_{k=1}^{10} \sigma_k^j) - \sigma_i^j (\sum_{k=1}^{10} \frac{\partial \sigma_k^j}{\partial E_{(r)}^j})}{\sum_{k=1}^n \sigma_k^j} \quad (\text{S14})$$

In addition, we can easily obtain



$$\begin{aligned}
\frac{\partial \sigma_k^j}{E_{(r)}^j} &= \frac{\varepsilon_{air}}{2} \frac{\partial}{E_{(r)}^j} \int E_{(r')}^j \overline{E_{(r')}^j} R_{k(r')} dr' = \frac{\varepsilon_{air}}{2} \int \delta_{(r,r')} \overline{E_{(r')}^j} R_{k(r')} dr' = \frac{\varepsilon_{air}}{2} \overline{E_{(r)}^j} R_{k(r)} \\
&\Rightarrow \frac{\partial C}{\partial E_{(r)}^j} = -\frac{\varepsilon_{air}}{2m} \sum_{i=1}^{10} \frac{y_i^j - h_i^j}{\sigma_i^j (1 - h_i^j)} (R_{i(r)} - h_i^j \sum_{k=1}^{10} R_{k(r)}) \overline{E_{(r)}^j}
\end{aligned} \tag{S15}$$

## 2.2. Nonlinear Materials

The nonlinearity that we have implemented in this work was designed to loosely follow the behavior of a ReLU activation function [16]. Of course, applying this activation function directly to the field would not be meaningful, so it had to be modified to act upon the intensity of the field. The idea here was to design the nonlinearity so that fields with intensity below a certain threshold are blocked. This objective can be achieved by adding an imaginary part of the form  $-\alpha/|E|^2$  to the permittivity.

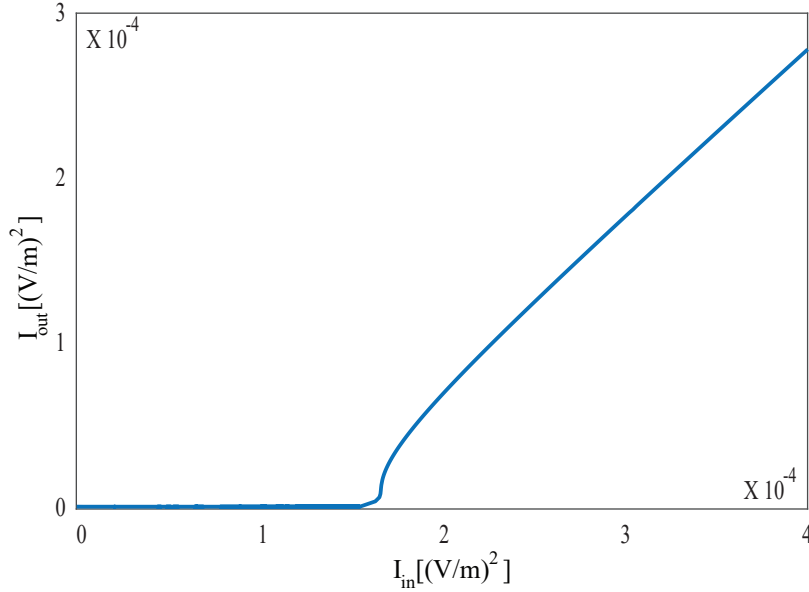


Fig. S5. The intensity response of the nonlinear layer in 1D as a function of the input intensity for a length of  $\lambda/2$  and an  $\alpha$  equal to  $5 \times 10^{-5} (\frac{V}{m})^2$ .

Although this type of nonlinearity was chosen rather arbitrarily, it alters the intensity in a similar way as optical saturable absorbers [1]. With the details of the proposed nonlinearity explained, its contribution to the gradient can now be assessed. To that end, we first define the permittivity in nonlinear sections as  $\varepsilon_{(r)} = \varepsilon_{real(r)} + i\varepsilon_{img(r)}$  where  $\varepsilon_{real(r)}$  and  $\varepsilon_{img(r)}$  are the real and imaginary parts of the permittivity respectively. As was shown earlier, the gradient has an extra term in nonlinear sections of the structure. The extra terms can now be calculated for this particular nonlinearity. The first extra term is calculated as follows

$$\begin{aligned} \frac{\partial L(r)}{\partial E(r)^j} &= \frac{\partial L(r)}{\partial \varepsilon(r)} \frac{\partial \varepsilon(r)}{\partial E(r)^j} = (\omega^2)(-i\alpha \frac{\overline{E(r)^j}}{|E(r)^j|^4}) = i\omega^2 \varepsilon_{img(r)} \frac{1}{E(r)^j} \\ L(r) + E(r)^j \frac{\partial L(r)}{\partial E(r)^j} &= (\nabla \times (\frac{1}{\mu}) \times \nabla \times) - \omega^2 \varepsilon(r) + i\omega^2 \varepsilon_{img(r)} = (\nabla \times (\frac{1}{\mu}) \times \nabla \times) - \omega^2 \varepsilon_{real(r)} \end{aligned} \quad (S16)$$

While the final term takes the following form.

$$\frac{\overline{E(r)^j} \partial \overline{L(r)}}{\partial E(r)^j} = \frac{\overline{E(r)^j} \partial \overline{L(r)}}{\partial \overline{\varepsilon(r)}} \frac{\partial \overline{\varepsilon(r)}}{\partial E(r)^j} = \omega^2 i\alpha \frac{\overline{E(r)^j} \overline{E(r)^j}}{|E(r)^j|^4} \quad (S17)$$

The final tool we need before putting all the concepts together is a way to solve the nonlinear wave equation. We use Finite-Differences Frequency-Domain(FDFD) method. In computational domain the wave equation can be written as a matrix multiplication shown in Eq.S18.

$$AE = -i\omega J \quad (S18)$$

In this equation  $A = (C_h D_\mu^{-1} C_e - \omega^2 D_\varepsilon)$ , where  $C_h$  and  $C_e$  are curl matrices,  $D_\mu$  and  $D_\varepsilon$  represent permeability and permittivity of different points in the domain respectively; while  $J$  represents the current source density matrix in the computational domain. Here we use iterative method to solve the nonlinear response. Finally, for the implementation of this method, we utilized the MATLAB-based package MaxwellFDFD [17]. At this stage, with all the necessary pieces available, we can finally formulate the problem in the forward and backward directions. Solving for the electric field in the forward direction is of this form

$$\begin{aligned} E^j &= (A^j)^{-1} b^j \\ \frac{\partial C}{\partial \varepsilon} &= -2\omega^2 \text{real}\left\{ \sum_{j=1}^m \lambda^j \odot E^j \right\} \end{aligned} \quad (S19)$$

While the adjoint equation (the problem of the adjoint field propagating backwards) takes the following form.

$$\frac{-\varepsilon}{m} \left( \left( \frac{y_i^j - h_i^j}{\sigma_i^j (1 - \sigma_i^j)} \right) (\gamma_i - h_i^j \sum_{n=1}^{10} \gamma_n) \odot \overline{E^j} \right) + \lambda^j (A^j - \omega^2 i\alpha \frac{1}{E^j \odot E^j}) + \overline{\lambda^j} (\omega^2 i\alpha \frac{1}{E^j \odot E^j}) = 0 \quad (S20)$$

In the above equations  $\odot$  represents the Hadamard product(the element wise multiplication between two matrices). Eq.S20 enables us to calculate the gradient with respect to each of the elements in the computational domain all at once. Once this gradient is computed, it is possible to choose which elements to update. Of course, since updating the permittivity of the nonlinear sections would lead to a medium that is far from any real physical interpretation, we decided to restrict the updating process to the linear sections of the medium.

### 2.3. Level-Set Evolution

As mentioned in the main paper, updating the structure at each step is done by evolving the boundary between the two constituent materials with the following equation.

$$\partial_t \phi + v(x, y) |\nabla \phi| = 0 \quad (S21)$$

Where  $\phi$  is the level-set function and  $v(x, y)$  is the velocity with which each point on the zero crossing curve of the level-set function moves normal to the the curve(which is set equal to the gradient with respect to permittivity of different points in our work). However, simply implementing this equation in the digital domain does not maintain a stable curve evolution process. To overcome this issue, we implemented the method introduced in ref [22] where an extra term is added to the evolution equation to ensure the level-set function remains a signed distance function through its evolution.

$$\partial_t \phi = \mu \nabla \cdot (p(|\nabla \phi|) \nabla \phi) - v(x, y) |\nabla \phi|$$

$$p(s) = \begin{cases} \frac{1}{(2\pi)s} \sin(2\pi s) & s \leq 1 \\ 1 - \frac{1}{s} & s \geq 1 \end{cases} \quad (S22)$$

Where  $\mu$  is a constant(set to 0.2 in this work) and  $p(s)$  is a double-well potential for distance regularization. With this approach, we do not need to concern ourselves with matters of stability for the level-set function.

#### 2.4. Training with FDFD Simulation

As discussed in the paper, both a 2D and a 3D model with the proposed method were trained. This training process required solving the wave equation for different inputs and in this section the details the FDFD simulations and training process are explained.

The 2D version of the medium was trained in a domain with PML boundary condition on all sides, while the 3D simulations were done with the PEC boundary condition on planes parallel to the direction of wave propagation and PML in the direction perpendicular to it. The different boundary conditions for the 3D version were implemented to reduce the computational load of the simulation. The wavelength for both simulations was set to  $1\mu m$ , while the spatial resolution of the simulation was set to  $\lambda/20$  and  $\lambda/5$  for 2D and 3D versions respectively.

In the main manuscript we went on to explain how the iterative process of solving the nonlinear wave equation worked. We explained that we start with a random initial wave distribution that effectively gave us the initial permittivities, and then use these initial values to solve for the new field distribution. This process was then repeated until the electric field converged to a fixed distribution. However, in practice, we start with a field solution corresponding to one of our input images. This helps the field to converge more quickly .

As was mentioned in the paper, we used SGD for the training process. More precisely we used a variation of SGD called mini-batch gradient descent. While SGD uses all the input data to calculate the gradient for each training iteration, mini-batch gradient descent only uses a small batch of the input data. This significantly increases the convergence speed of the optimization process. Of course, the use of batches was already hinted at in the main text, so we continue by explaining some of the implementation details regarding batches and learning rate. We have used batches of 100 for the 2D training and batches containing 50 samples for the 3d version. The larger batch size for the 2D version was chosen so, because the selected method of training showed that it was more sensitive to batch size and therefore it was better to train with a larger batch. Finally, the learning rate for the 2D version was set to 500 which was reduced after two epochs of training, while its 3D counterpart was set to 30.

Fig.S3 shows the cost values for the training set and the accuracy values for the test set on the 2D binary medium.

In the main text, we showed the confusion matrix for the case of the 3D NNM. On one axis, this matrix shows the ground truth labels for each input image, and on the other, the labels produced by the NNM for that image. Therefore, the diagonal elements show the number of correct classifications by the NNM, while the rest of the elements on each row show the wrong

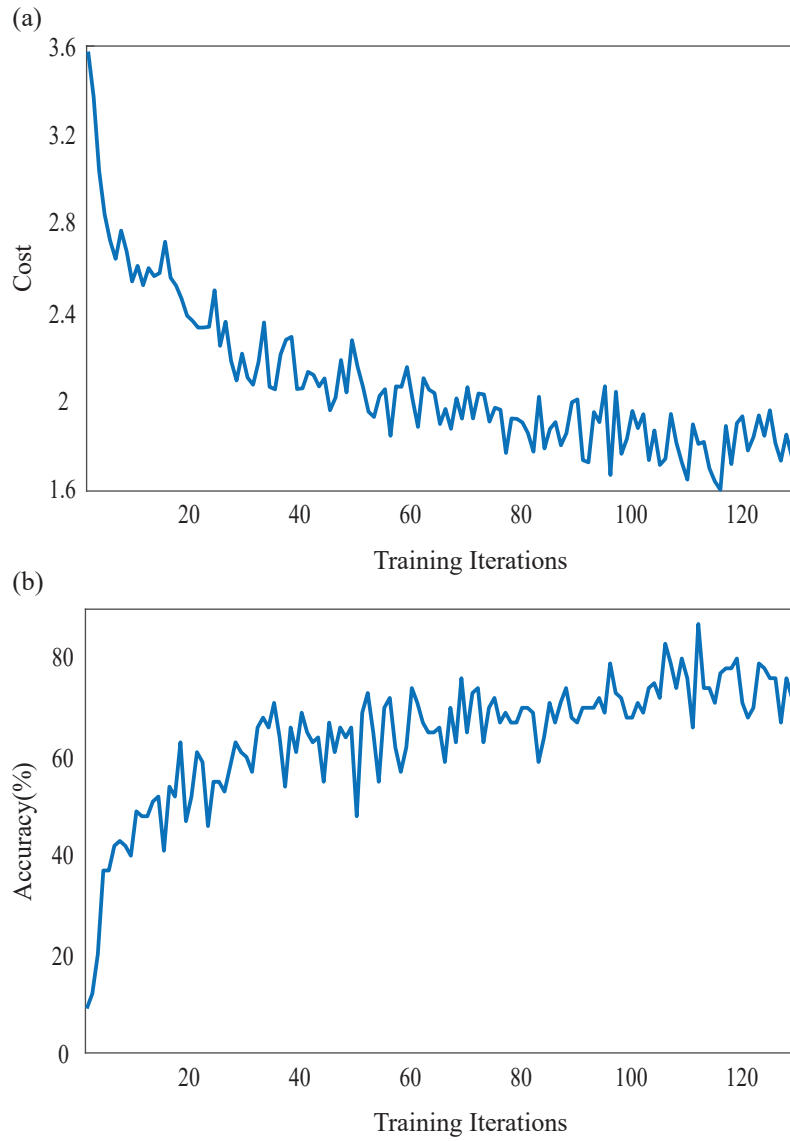


Fig. S6. (a) The training curve for 2D NNM. (b) Evolution of the accuracy with each training iteration on batches of test set equal in size to the training batch for the same medium.

classifications. In Fig.S5, the confusion matrix for the 2D NNM has been depicted where 10 samples of each digit have been presented for the medium.

2D

Ground truth labels	0	10	0	0	0	0	0	0	0	0	
	1	0	10	0	0	0	0	0	0	0	
	2	0	0	7	1	0	1	0	0	1	
	3	0	0	1	8	0	0	0	0	1	
	4	0	0	0	0	7	0	0	0	3	
	5	1	0	0	3	0	6	0	0	0	
	6	0	1	0	0	0	0	8	0	1	
	7	0	0	0	0	0	1	0	8	0	
	8	0	0	1	1	1	0	0	0	7	
	9	0	0	0	0	4	0	0	0	0	
		0	1	2	3	4	5	6	7	8	9

Predicted labels

Fig. S7. Confusion matrix for the 2D PNCM. The diagonal elements show the number of correct classifications among the 10 input digits presented as the input, while the other columns show what the wrong classes that each digit was categorized as are.

### Funding

The work was financially supported by DARPA Young Faculty Award program.

### Acknowledgments

The authors thank W. Shin for his help on improving the computational speed of the implementation of this method.

### Disclosures

The authors declare that there are no conflicts of interest related to this article.

### Conclusion

Here we show that the wave dynamics in the Maxwell's equations is capable of performing highly sophisticated computing. There is intricate connection between differential equations that govern many physical phenomena and neural computing (see more discussion in supplementary), which could be further explored. From the perspective of optics, the functions of most nanophotonic devices can be described as mode mapping [27]. In traditional nanophotonic devices, mode mapping mostly occurs between eigenmodes. For example, a polarization beam splitter [13] maps each polarization eigenmode to a spatial eigenmode. Here, we introduce a class of nanophotonic media that can perform complex and nonlinear mode mapping equivalent to artificial neural computing. The neural computing media shown here has an appearance of disorder media. It would be also interesting to see how disorder media, which support rich physics such as Anderson localization, could provide a new platform for neural computing. In comparison, today's optical neural computing mostly follows layered structures. While highly efficient for digital computing, layered structures could be counter-productive in optical analog computing.

For practical applications that routinely use millions of connections, it can be challenging to implement deep and dense layers of optical components in a compact form. Nevertheless, the concept of NNM shows that any nanostructures could be optimized to perform neural computing without the rigid constraint of layer structures. Combined with ultra-high computing density, NNM could be used in a wide range of information devices as the analog preprocessing unit.

## References

1. Y. Shen, N. C. Harris, S. Skirlo, M. Prabhu, T. Baehr-Jones, M. Hochberg, X. Sun, S. Zhao, H. Larochelle, D. Englund, and M. Soljacic, "Deep learning with coherent nanophotonic circuits," *Nat. Photonics* **11**, 441 (2017).
2. M. Hermans, M. Burm, T. Van Vaerenbergh, J. Dambre, and P. Bienstman, "Trainable hardware for dynamical computing using error backpropagation through physical media," *Nat. communications* **6**, 6729 (2015).
3. T. W. Hughes, M. Minkov, Y. Shi, and S. Fan, "Training of photonic neural networks through in situ backpropagation and gradient measurement," *Optica* **5**, 864–871 (2018).
4. S. R. Skinner, E. C. Behrman, A. A. Cruz-Cabrera, and J. E. Steck, "Neural network implementation using self-lensing media," *Appl. optics* **34**, 4129–4135 (1995).
5. P. R. Prucnal and B. J. Shastri, *Neuromorphic photonics* (CRC Press, 2017).
6. H. G. Chen, S. Jayasuriya, J. Yang, J. Stephen, S. Sivaramakrishnan, A. Veeraraghavan, and A. Molnar, "Asp vision: Optically computing the first layer of convolutional neural networks using angle sensitive pixels," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, (2016), pp. 903–912.
7. J. Bueno, S. Maktoobi, L. Froehly, I. Fischer, M. Jacquot, L. Larger, and D. Brunner, "Reinforcement learning in a large-scale photonic recurrent neural network," *Optica* **5**, 756–760 (2018).
8. 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**, 1004–1008 (2018).
9. M. Hermans and T. Van Vaerenbergh, "Towards trainable media: Using waves for neural network-style training," arXiv preprint arXiv:1510.03776 (2015).
10. H.-G. Park, S.-H. Kim, S.-H. Kwon, Y.-G. Ju, J.-K. Yang, J.-H. Baek, S.-B. Kim, and Y.-H. Lee, "Electrically driven single-cell photonic crystal laser," *Science* **305**, 1444–1447 (2004).
11. J. D. Joannopoulos, S. G. Johnson, J. N. Winn, and R. D. Meade, *Photonic crystals: molding the flow of light* (Princeton university press, 2011).
12. W. Cai and V. Shalaev, *Optical metamaterials: fundamentals and applications* (Springer Science & Business Media, 2009).
13. B. Shen, P. Wang, R. Polson, and R. Menon, "An integrated-nanophotonics polarization beamsplitter with  $2.4 \times 2.4 \mu\text{m}^2$  footprint," *Nat. Photonics* **9**, 378 (2015).
14. A. Y. Piggott, J. Petykiewicz, L. Su, and J. Vučković, "Fabrication-constrained nanophotonic inverse design," *Sci. Reports* **7**, 1786 (2017).
15. L. Su, A. Y. Piggott, N. V. Saprà, J. Petykiewicz, and J. Vuckovic, "Inverse design and demonstration of a compact on-chip narrowband three-channel wavelength demultiplexer," *ACS Photonics* (2017).
16. V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *Proceedings of the 27th international conference on machine learning (ICML-10)*, (2010), pp. 807–814.
17. W. Shin, "MaxwellFDFD Webpage," (2015). <https://github.com/wsshin/maxwelldfd>.
18. H. J. Caulfield, J. Kinser, and S. K. Rogers, "Optical neural networks," *Proc. IEEE* **77**, 1573–1583 (1989).
19. L. Jing, Y. Shen, T. Dubcek, J. Peurifoy, S. Skirlo, Y. LeCun, M. Tegmark, and M. Soljačić, "Tunable efficient unitary neural networks (EUNN) and their application to RNNs," in *Proceedings of the 34th International Conference on Machine Learning*, vol. 70 of *Proceedings of Machine Learning Research* D. Precup and Y. W. Teh, eds. (PMLR, International Convention Centre, Sydney, Australia, 2017), pp. 1733–1741.
20. T. W. Hughes, M. Minkov, I. A. Williamson, and S. Fan, "Adjoint method and inverse design for nonlinear nanophotonic devices," *ACS Photonics* **5**, 4781–4787 (2018).
21. M. Courbariaux, I. Hubara, D. Soudry, R. El-Yaniv, and Y. Bengio, "Binarized neural networks: Training deep neural networks with weights and activations constrained to+ 1 or-1," arXiv preprint arXiv:1602.02830 (2016).
22. C. Li, C. Xu, C. Gui, and M. D. Fox, "Distance regularized level set evolution and its application to image segmentation," *IEEE transactions on image processing* **19**, 3243–3254 (2010).
23. J. Bergstra and Y. Bengio, "Random search for hyper-parameter optimization," *J. Mach. Learn. Res.* **13**, 281–305 (2012).
24. J. Snoek, H. Larochelle, and R. P. Adams, "Practical bayesian optimization of machine learning algorithms," in *Advances in neural information processing systems*, (2012), pp. 2951–2959.
25. S. Saxena and J. Verbeek, "Convolutional neural fabrics," in *Advances in Neural Information Processing Systems*, (2016), pp. 4053–4061.
26. G. D. Marshall, A. Politi, J. C. Matthews, P. Dekker, M. Ams, M. J. Withford, and J. L. O'Brien, "Laser written waveguide photonic quantum circuits," *Opt. express* **17**, 12546–12554 (2009).
27. D. A. Miller, "All linear optical devices are mode converters," *Opt. Express* **20**, 23985–23993 (2012).