

483 **A1 Optical field simulation**

484 Analyzing how light field propagate through those components are critical to device optimization  
 485 and photonic integrated circuit design. Given a linear isotropic optical component, we will shine  
 486 time-harmonic continuous-wave light on its input ports and analyze the steady-state electromagnetic  
 487 field distributions  $\mathbf{E} = \hat{x}\mathbf{E}_x + \hat{y}\mathbf{E}_y + \hat{z}\mathbf{E}_z$  and  $\mathbf{H} = \hat{x}\mathbf{H}_x + \hat{y}\mathbf{H}_y + \hat{z}\mathbf{H}_z$  in it, each of which includes  
 488 horizontal ( $x$ ), vertical ( $y$ ), and longitudinal ( $z$ ) components. The light field follows the Maxwell  
 489 PDE under certain absorptive boundary conditions [15],

$$\nabla \times \mathbf{E}(\mathbf{r}, t) = \mu_0 \frac{\partial \mathbf{H}(\mathbf{r}, t)}{\partial t} + \mathbf{J}_e(\mathbf{r}, t), \quad \nabla \times \mathbf{H}(\mathbf{r}, t) = -\epsilon_0 \epsilon_r(\mathbf{r}) \frac{\partial \mathbf{E}(\mathbf{r}, t)}{\partial t} + \mathbf{J}_e(\mathbf{r}, t), \quad (7)$$

490 where  $\nabla \times$  is the curl operator of a vector function,  $\mu_0$  is the vacuum magnetic permeability,  $\epsilon_0$  and  
 491  $\epsilon_r$  are the vacuum and relative electric permittivity,  $\mathbf{J}_m$  and  $\mathbf{J}_e$  are the magnetic and electric current  
 492 sources. Since the input light is time-harmonic at a vacuum angular frequency  $\omega$ , the time-domain  
 493 PDE can be transformed to the frequency domain for the steady state as follows,

$$\nabla \times \mathbf{E}(\mathbf{r}) = j\omega \mu_0 \mathbf{H}(\mathbf{r}) + \mathbf{J}_m(\mathbf{r}), \quad \nabla \times \mathbf{H}(\mathbf{r}) = -j\omega \epsilon_0 \epsilon_r(\mathbf{r}) \mathbf{E}(\mathbf{r}) + \mathbf{J}_e(\mathbf{r}). \quad (8)$$

494 A simple variable substitution gives us the *curl-of-curl* Maxwell PDE,

$$((\mu_0^{-1} \nabla \times \nabla \times) - \omega^2 \epsilon_0 \epsilon_r(\mathbf{r})) \mathbf{E}(\mathbf{r}) = j\omega \mathbf{J}_e(\mathbf{r}), \quad (\nabla \times (\epsilon_r^{-1}(\mathbf{r}) \nabla \times) - \omega^2 \mu_0 \epsilon_0) \mathbf{H}(\mathbf{r}) = j\omega \mathbf{J}_m(\mathbf{r}). \quad (9)$$

495 To restrict a unique solution without boundary reflection, complicated boundary conditions will be  
 496 inserted [15]. An artificial material, i.e., coordinate-stretched perfectly matched layer (SC-PML),  
 497 will be padded around the solving domain. Such PML materials have large imaginary parts in  
 498 the permittivities to introduce strong energy absorption and changes the derivative operator to  
 499  $\nabla = (\frac{1}{s_x(x)} \frac{\partial}{\partial x}, \frac{1}{s_y(y)} \frac{\partial}{\partial y}, \frac{1}{s_z(z)} \frac{\partial}{\partial z})$ , where  $s$  is a location-determined complex value. Solving the  
 500 above PDEs will give the steady-state frequency-domain complex magnitude of the optical fields.

501 **A2 Dataset generation**

502 We generate our customized MMI device simulation dataset using an open-source FDFD simulator  
 503 [angler](#) [15]. The tunable MMI dataset has 5.5 K *single-source* training data, 614 validation data, and  
 504 1.5 K multi-source test data. The etched MMI dataset has 12.4 K *single-source* training data, 1.4 K  
 505 validation data, and 1.5 K *multi-source* test data. We summarize how we generate random devices in  
 506 Table A4. We randomly sample the physical dimension of the MMI, input/output waveguide width,  
 507 the width of the perfectly matched layer (PML), device border width away from PML, controlling  
 508 pad sizes, input light source frequencies, etched cavity sizes and ratio (determines the number of  
 cavities in the MMIs), and permittivities in the controlling region.

Table A4: Summary of device design variable's sampling range, distribution, and unit.

Variables	Value/Distribution		Unit
	$ \mathbf{J}  \times  \mathbf{J} $ Tunable MMI	$ \mathbf{J}  \times  \mathbf{J} $ Etched MMI	
Length	$\mathcal{U}(20, 30)$	$\mathcal{U}(20, 30)$	$\mu m$
Width	$\mathcal{U}(5.5, 7)$	$\mathcal{U}(5.5, 7)$	$\mu m$
Port Length	3	3	$\mu m$
Port Width	$\mathcal{U}(0.8, 1.1)$	$\mathcal{U}(0.8, 1.1)$	$\mu m$
Border Width	0.25	0.25	$\mu m$
PML Width	1.5	1.5	$\mu m$
Pad Length	$\mathcal{U}(0.7, 0.9) \times \text{Length}$	$\mathcal{U}(0.7, 0.9) \times \text{Length}$	$\mu m$
Pad Width	$\mathcal{U}(0.4, 0.65) \times \text{Width}/ \mathbf{J} $	$\mathcal{U}(0.4, 0.65) \times \text{Width}/ \mathbf{J} $	$\mu m$
Wavelengths $\lambda$	$\mathcal{U}(1.53, 1.565)$	$\mathcal{U}(1.53, 1.565)$	$\mu m$
Cavity Ratio	-	$\mathcal{U}(0.05, 0.1)$	-
Cavity Size	-	0.027 Length $\times$ 0.114 Width	$\mu m^2$
Relative Permittivity $\epsilon_r$	$\mathcal{U}(11.9, 12.3)$	{2.07, 12.11}	-

509

510 **A3 Training settings**

511 We implement all models and training logic in PyTorch 1.10.2. All experiments are conducted on a  
 512 machine with Intel Core i7-9700 CPUs and an NVIDIA Quadro RTX 6000 GPU. For training from

513 scratch, we set the number of epochs to 200 with an initial learning rate of 0.002, cosine learning rate  
514 decay, and a mini-batch size of 12. For the tunable MMI dataset, we split all 7,680 examples into  
515 72% training data, 8% validation data, and 20% test data. For the etched MMI dataset, we split all  
516 15,360 examples into 81% training data, 9% validation data, and 10% test data. For device adaptation,  
517 we first perform linear probing for 20 epochs with an initial learning rate of 0.002 and cosine learning  
518 rate decay; then we perform finetuning for 30 epochs with an initial learning rate of 0.0002 and a  
519 cosine learning rate decay. We apply stochastic network depth with a linear scaling strategy and a  
520 maximum drop rate of 0.1.

## 521 A4 Model architectures

522 **UNet.** We construct a 4-level convolutional UNet with a base channel number of 34. The total  
523 parameter count is 3.47 M.

524 **FNO-2d.** For Fourier neural operator (FNO), we use 5 2-D FNO layers with a channel number of 32.  
525 The Fourier modes are set to ( $\#Mode_z=32$ ,  $\#Mode_x=10$ ). The final projection head is  $CONV1 \times 1(256)$ -  
526  $GELU-CONV1 \times 1(2)$ . The total parameter count is 3.29 M.

527 **F-FNO.** For factorized Fourier neural operator (F-FNO), we use 12 F-FNO layers with a channel  
528 number of 48. The Fourier modes are set to ( $\#Mode_z=70$ ,  $\#Mode_x=40$ ). The final projection head is  
529  $CONV1 \times 1(256)$ - $GELU-CONV1 \times 1(2)$ . The total parameter count is 3.16 M.

530 **NeurOLight.** For our proposed NeurOLight, we use 12 F-FNO layers for tunable MMIs and 16  
531 layers for etched MMIs with a base channel number  $C=64$ . The convolution stem is  $BSConv3 \times 3(32)$ -  
532  $BN-ReLU-BSConv3 \times 3(64)$ - $BN-ReLU$ , where  $BSConv$  is blueprint convolution [12]. The Fourier  
533 modes are set to ( $\#Mode_z=70$ ,  $\#Mode_x=40$ ). The channel expansion ratio in the FFN is set to  $s=2$ .  
534 The final projection head is  $CONV1 \times 1(256)$ - $GELU-CONV1 \times 1(2)$ . The total parameter count is  
535 1.58 M.

## 536 A5 Animation of NeurOLight

537 We animate the prediction process of NeurOLight on a  $3 \times 3$  tunable MMI. The prediction  
538 throughput reaches 120 frames per second (FPS) which allows designers to freely tune input fre-  
539 quencies, device sizes, permittivities, and input light sources with *real-time* simulation result feedback.

540

541