# From Hours to Seconds: Towards 100x Faster Quantitative Phase Imaging via Differentiable Microscopy

UDITH HAPUTHANTHRI, [1,2] KITHMINI HERATH, [1,2] RAMITH HETTIARACHCHI, [1,2] HASINDU KARIYAWASAM, [1,2] AZEEM AHMAD, [3] BALPREET S. AHLUWALIA, [3] GANESH ACHARYA, [4] CHAMIRA U. S. EDUSSOORIYA, [2,*] AND DUSHAN N. WADDUWAGE [1,*]

[1] *Center for Advanced Imaging, Faculty of Arts and Sciences, Harvard University, Cambridge, MA 02138, USA*
[2] *Department of Electronic and Telecommunication Engineering, University of Moratuwa, Sri Lanka*
[3] *Department of Physics and Technology, UiT The Arctic University of Norway, Tromsø 9037, Norway*
[4] *Division of Obstetrics and Gynecology, Department of Clinical Science, Intervention and Technology, Karolinska Institute, Stockholm, Sweden*
[*] *wadduwage@fas.harvard.edu; chamira@uom.lk*

**Abstract:** With applications ranging from metabolomics to histopathology, quantitative phase microscopy (QPM) is a powerful label-free imaging modality. Despite significant advances in fast multiplexed imaging sensors and deep-learning-based inverse solvers, the throughput of QPM is currently limited by the speed of electronic hardware. Complementarily, to improve throughput further, here we propose to acquire images in a compressed form such that more information can be transferred beyond the existing electronic hardware bottleneck. To this end, we present a learnable optical compression-decompression framework that learns content-specific features. The proposed differentiable quantitative phase microscopy ($\partial\mu$) first uses learnable optical feature extractors as image compressors. The intensity representation produced by these networks is then captured by the imaging sensor. Finally, a reconstruction network running on electronic hardware decompresses the QPM images. In numerical experiments, the proposed system achieves compression of $\times 64$ while maintaining the SSIM of $\sim 0.90$ and PSNR of $\sim 30$ dB on cells. The results demonstrated by our experiments open up a new pathway for achieving end-to-end optimized (i.e., optics and electronic) compact QPM systems that may provide unprecedented throughput improvements.

## 1. Introduction

Among the label-free imaging modalities, quantitative phase microscopy (QPM) is a simple but powerful approach, providing important biophysical information by quantifying optical phase differences [1, 2]. From the phase map, one can further yield quantitative information about the morphology and dynamics of the examined specimens [3, 4]. In addition to morphology, the measured phase maps can be converted to dry mass of the cells with accuracy that is of the order of femtograms per square microns [5, 6]. QPM has found many important applications in biomedicine [7] including pathogen screening [8], cancer cell classification [9], and label-free analysis of histopathology specimens [10, 11]. Importantly, label-free histopathology based on QPM has been shown to capture subtle, nanoscale morphological properties of tissues that could lead to early detection of cancer [12]. This technique also preserves the tissue sample for further molecular-specific pathological analysis [13]. Moreover, recently quantitative phase imaging has even been extended to image the structures of thick biological systems such as zebrafish larval [14]. QPM has also been demonstrated for stain-free quantification of chromosomal dry mass in living cells [15], quantification of different growth phases of chondrocytes [16], identification of biophysical markers of sickle cell drug responses [17] and quantification of

nuclear mechanical properties that may play roles in cancer metastasis [18].

The first phase imaging mechanism was introduced by Zernike in his phase contrast microscopy [19]. Here the phase shifts due to the refractive indices and depth differences in the specimen are converted into detectable intensity variations. Zernike's original design consisted of a phase filter which directly displays phase information by interfering scattered portion of light from an image, with its unscattered portion. Even though the work improved with several extensions [20, 21], due to the non-linear dependency between phase and intensity, direct phase contrast techniques are incapable of quantitative phase measurements. QPM techniques overcome this problem by computational inverse reconstruction [7]. A typical quantitative phase microscope consists of an optical system (forward model) and a computational phase retrieval algorithm (inverse model) [22]. The forward optical system converts undetectable phase information into detectable fringe patterns; from the fringe patterns, the inverse reconstruction algorithm retrieves phase and intensity maps of the specimen. Recent developments in QPM have mostly been focused on improving the inverse reconstruction using GPU acceleration [23–25], deep-learning-based inverse solvers [26–31], and illumination patterns optimization [32, 33].

Orthogonal to the current developments, the main bottleneck of QPM-based imaging cytometry is the image acquisition speed, which is fundamentally governed by the pixel rate of image sensors. Currently, the pixel rate of a state-of-the-art camera is around $1 \times 10^{10}$ pixels/sec. However, the pixel throughput of the front-end optics is virtually unlimited. An image passes through optics at the speed of light and has been the rationale for developing optical signal processing technologies [34]. Here we propose to exploit this property to optically compress an image in order to measure the compressed form of the image using a high-speed light detector (such as a high-speed camera). Thus the pixel throughput of the original image would be increased at a rate proportional to the degree of compression. Compressive imaging of biological specimens has been demonstrated before, using random sampling of the linearly projected image space [35]. Better compression, however, may be achieved through learning dataset-specific features of images. To this end, here we propose differentiable microscopy ($\partial \mu$) to identify important image features for compression, learned through machine learning approaches [36]. First, we modeled the compression and de-compression as a single auto-encoder neural network. Afterward, with a large set of target images, we trained the neural network to find a low-dimensional compressed representation of the images. Once trained, the decoder part of the network acts as a decompression algorithm. We then model the encoder part of the network as a learnable optical system to be used as an optical phase feature extractor. More specifically, the optical phase feature extractor encodes phase information at the compressed latent intensity field. The intensity field is then acquired by a photodetector array. The detected intensity feature map is decoded as the phase image by an electronic phase reconstruction network.

Our optical phase feature extractor is motivated by the previous work on all-optically retrieving the phase information through learnable optics [37]. This work modeled phase retrieval as an optimization problem where the optical network's physical parameters were optimized by minimizing the pixel distance between the input phase and the output intensity. The main limitation of this previous work is the lack of *non-linearity* of the all-optical method. In contrast, here we treat the all-optical linear network as *a feature extractor* which can extract features from both phase and amplitude of the input field. The optical model only has to learn a faithful representation that contains sufficient information to reconstruct the original phase. Combining this all-optical phase feature extractor with the non-linear phase reconstruction network, we realize an end-to-end-learnable non-linear function for the phase retrieval task. In the following sections, we first introduce the proposed end-to-end differentiable compressive QPM (section 2.1). Second, we assess the feasibility of phase retrieval using linear compression and non-linear decompression functions (section 2.2). Third, we demonstrate *Compressed QPM* using an optical encoder and an electronic decoder. We use learnable Fourier filters (LFFs) [37] or PhaseD2NNs [37] as the
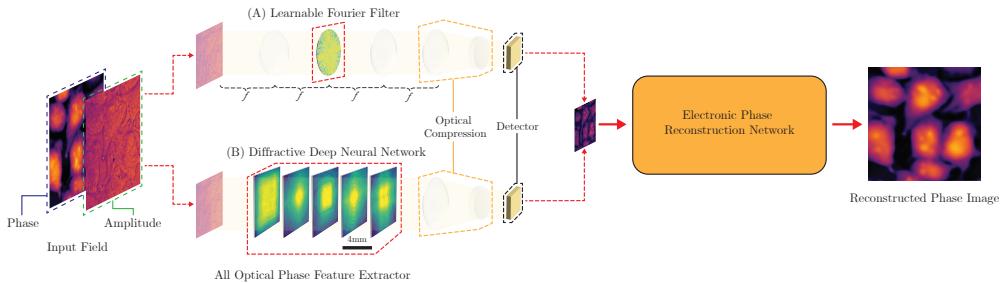
Fig. 1. **End-to-end pipeline**: Input field with high contrast phase information is fed into the proposed optical phase feature extraction network. The resultant compressed output intensity field which contains the phase features is captured by the detector. The electronic phase reconstruction network utilizes these features to reconstruct phase information.

optical feature extractors. We use SwinIR [38], a state-of-the-art super-resolution pipeline, as the electronic reconstruction network. Our experiments with experimental datasets suggest that the proposed method can perform orders of magnitude faster QPM on cells than the current state of the art through $\times 64 - \times 256$ compression.

## 2. Results

### 2.1. End-to-end Differentiable Compressive Quantitative Phase Microscopy

We model the QPM system as a combination of an optical and electronic network that learns the task of reconstructing phase information of the input light field at the output of the electronic network. The optical network converts useful morphological features to an intensity field which is then captured by the detector array placed at the output plane of this network. Then, the electronic network constructs the phase map of the input light field from the captured features. This task is embedded in a loss function, and the entire network is parameterized in a differentiable manner. The parameters of the network are optimized to reduce the loss for a particular dataset. We follow a 3-stage optimization criteria for the improved stability of the end-to-end optimization; 1) optimize the all-optical phase extraction network; 2) optimize the electronic phase reconstruction network; 3) end-to-end fine-tuning.

**Optimize the all-optical phase extraction network.** Here the optical network is optimized to reconstruct the phase at its output intensity. For an input optical field $x_{in} = A_{in}e^{j\phi_{in}}$ we train an optical model $H_O$ through which the input field is propagated to produce the output field $x_{out} = A_{out}e^{j\phi_{out}} = H_O(x_{in})$. The **phase reconstruction loss**, $\mathcal{L}_\phi$ introduced in previous work [37] is utilized here as,

$$\mathcal{L}_\phi = \mathbb{E}_{x_{in} \sim P_X} \left[ L1(|A_{out}|^2, \phi_{in}/(2\pi)) \right], \tag{1}$$

where, $P_X$ and $L1(.)$ respectively represent the probability distribution of phase objects and the L1 loss.

**Optimize the electronic phase reconstruction network.** At this stage, we consider the end-to-end network, however, only the weights of the electronic network are optimized. The pretrained optical network discussed earlier is utilized as a feature extractor to encode the input phase. We demagnify the output field of the optical network to compress the intensity representation.

The electronic super-resolution network reconstructs the input phase from the compressed intensity representation. The reconstructed phase information is given by $\hat{\phi} = H_E(D(|A_{out}|^2))$. Here $H_E(.)$ and $D(.)$ represent the electronic phase reconstruction network and the optical demagnification layer, respectively. $D(.)$, the optical demagnification layer is implemented through a stack of $2 \times 2$ average pooling operations [39]. Similar to previous work [38], we consider $\mathcal{L}_{swin}$, a combination of loss functions for this optimization,

$$\mathcal{L}_{swin} = \mathbb{E}_{x_{in} \sim P_X} \left[ L1(\hat{\phi}, \phi_{in}/(2\pi)) + \mathcal{L}_{perceptual}(\hat{\phi}, \phi_{in}/(2\pi)) + \mathcal{L}_{adversarial}(\hat{\phi}, \phi_{in}/(2\pi)) \right],$$

$$(2)$$

where, $\mathcal{L}_{perceptual}$ and $\mathcal{L}_{adversarial}$ represent the perceptual loss [38] and adversarial loss [38], respectively.

**End-to-end fine-tuning.** As the final stage, we finetune the entire optical-electronic network to reconstruct the phase at the output of the network. To improve the reconstruction in terms of capturing fine cell structures, we incorporate the negative structural similarity index measure (SSIM) [40] as the loss function.

$$\mathcal{L}_{SSIM} = \mathbb{E}_{x_{in} \sim P_X} - \frac{1}{M} \sum_{j=1}^{M} \frac{(2\mu_{X_j}\mu_{Y_j} + C_1)(2\sigma_{X_j Y_j} + C_2)}{\left(\mu_{X_j}^2 + \mu_{Y_j}^2 + C_1\right)\left(\sigma_{X_j}^2 + \sigma_{Y_j}^2 + C_2\right)},$$

$$(3)$$

$X_j$ and $Y_j$ represent equal-sized windows from a normalized input phase image ($\phi_{in}/2\pi$) and the corresponding reconstructed phase output ($\hat{\phi}$), respectively, for $M$ number of windows for an image. $P_X$ represents the probability distribution of input phase objects. $\mu_{X_j}, \mu_{Y_j}, \sigma_{X_j}, \sigma_{Y_j}, \sigma_{X_j Y_j}$ are the means, variances and the covariance of the $X_j$ and $Y_j$ windows, respectively. $C_1 = (k_1 \times L)^2$ and $C_2 = (k_2 \times L)^2$ are regularization parameters with $L = 1.0$, $k_1 = 0.01$ and $k_2 = 0.03$.

In our formulation, since the input to the model is a complex-valued optical field, we only require input optical fields as training data. For the experiments, we utilize optical fields that are experimentally measured from phase objects. We consider a HeLa Cell dataset as our primary dataset. The optical fields in this dataset (i.e. full fields of view (full FoVs)) are $250\mu m \times 250\mu m$. Each full FoV is $789 \times 789$ pixel grid where each pixel is 316.4 nm $\times$316.4 nm. The dataset with full FoVs is divided into train and test sets. We cropped full FoVs into $256 \times 256$ sized patches (i.e. patch FoVs) for the end-to-end training of the proposed methods. Further details of the dataset are presented in the methods (section 5). Using the above loss functions and training data, we train optical-electronic networks that consist of either an LFF or a PhaseD2NN (see section 2.3) as the optical network. To further improve the realisticity of the proposed method, we also conduct experiments with detector noise. In the next sections, we discuss these models and their results.

## 2.2. Linear Encoding Does not Degrade Compressibility

The optical feature extractor is a linear system. We therefore first established the feasibility of linear compression in comparison to nonlinear compressor models.

**Linear Encoding and Non-linear Decoding Allow Compression.** Prior to experimentation with an optical-electronic neural network, we first investigated the computational capabilities of the optics-based encoder and electronics-based decoder. Due to the linear nature of the optical encoder, we experimented on an autoencoder (AE) network [41] with a *linear encoder* followed by a *non-linear decoder*. The reconstruction results obtained from this network were compared with a fully linear autoencoder and a fully nonlinear autoencoder. The qualitative results in Fig. 2 show that the autoencoder network with a *linear encoder* and *non-linear decoder* performs on par with the fully nonlinear autoencoder.
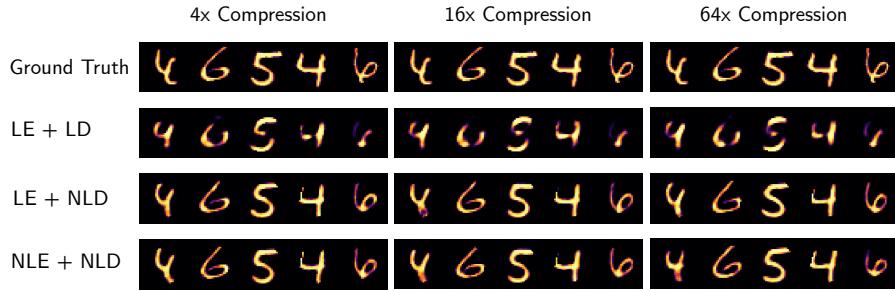
Fig. 2. Compressibility of MNIST images using autoencoders (AE) with linear (L) and nonlinear (NL) encoder (E)/ decoder (D). LE, LD, NLE, and NLD represent linear encoder, linear decoder, non-linear encoder, and non-linear decoder respectively.

**Complex-valued Linear Encoding and Non-linear Decoding Allow Compression of Phase Information.** While we assessed the computational feasibility of a linear encoder followed by a non-linear decoder to perform reconstruction, in QPM, another main hurdle is that information of interest is in the phase of the light field. Therefore, we further assessed the ability of an autoencoder network (complex-valued linear encoder + non-linear decoder) to extract, compress, and reconstruct information encoded in the phase. Similar to previous results, Fig. 3 shows that a complex-valued linear encoder with a nonlinear decoder achieves similar qualitative performance as the complex-valued nonlinear encoder and nonlinear decoder.

These results suggest the computational feasibility of a linear optical network (encoder) followed by a nonlinear electronic network (decoder) to reconstruct information in the phase of the light field.
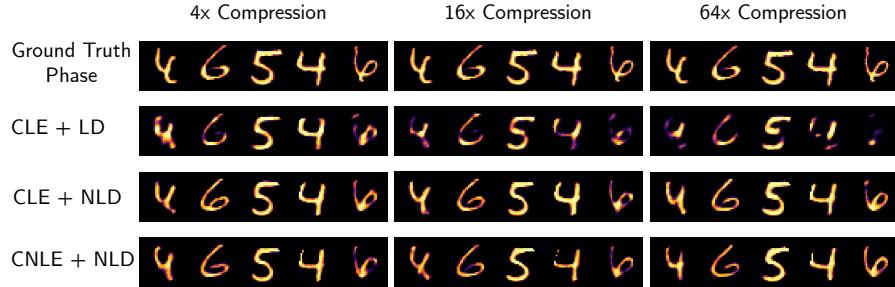


Fig. 3. Phase to intensity conversion and compressibility of MNIST images using linear (L) and nonlinear (NL) encoder (E)/ decoder (D). Both the encoders are complex-valued hence denoted as CLE and CNLE.

## 2.3. Optical Encoding and Electronic Decoding Enable Compressed QPM

Our results in section 2.2 show that an autoencoder with a linear encoder and a non-linear decoder (*AE:LE+NLD*) can reconstruct images as good as a fully nonlinear model. In this section, we propose our model with an all-optical feature extractor (linear encoder) and an electronic image reconstruction network (non-linear decoder) for QPM.
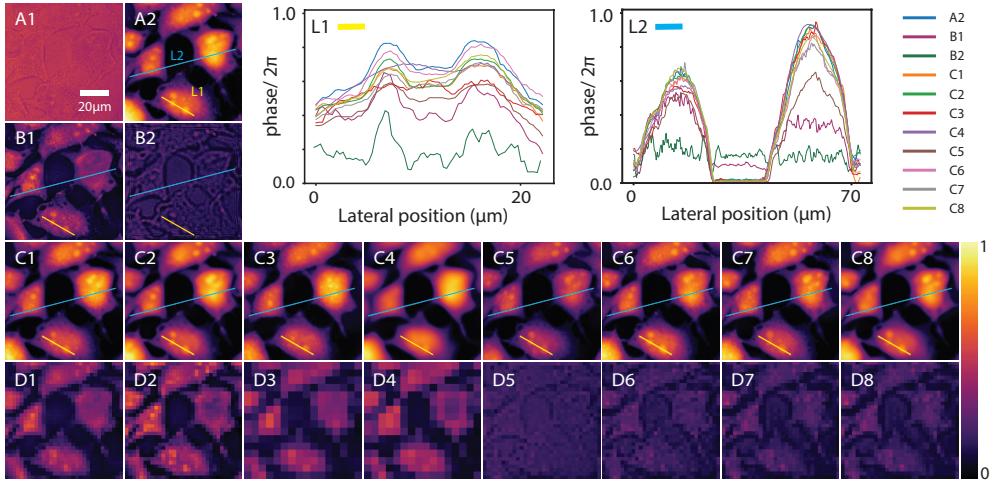
Fig. 4. **Qualitative performance comparison**: Amplitude of input field (A1), phase of the input field (A2) (i.e. patch FoV) from the test set, all-optical phase to intensity conversion results [37] (i.e. baselines) using LFF (B1), PhaseD2NN (B2), Phase reconstructions from our approach 1: LFF + SwinIR with ×64 compression without fine-tuning (C1), with fine-tuning (C2), LFF with ×256 compression without fine-tuning (C3), with fine-tuning (C4), Phase reconstructions from our approach 2: PhaseD2NN + SwinIR with ×64 compression without fine-tuning + 1 optical layer (C5), without fine-tuning + 3 optical layers (C6), without fine-tuning + 5 optical layers (C7), with fine-tuning + 5 optical layers (C8), Corresponding compressed output intensity fields of optical feature extractor (D1-8). Phase values along the L1 and L2 lines show the local and global resolving power of the proposed methods.

**Learnable Fourier Filter (LFF) + SwinIR.** Based on previous work [37], we first used a Learnable Fourier Filter (an LFF) as the optical feature extraction network. The LFF contained an optical 4-$f$ system with a learnable circular Fourier filter. Similar to previous work [37], the transmission coefficients of the circular Fourier filter were treated to be learnable. The input and output fields were $256 \times 256$ squared aperture grids. The circular Fourier filter had a radius of 128 grid points. The coefficients of the filter were randomly initialized. We used SwinIR [38], a state-of-the-art super-resolution network, as the electronic reconstruction network. We observed that directly training the end-to-end model (optical and electronic) was not ideal as the gradient flow between the optical and electronic networks was weak. Therefore, we employed the 3-stage criteria for the optimization of the end-to-end model (as discussed in section 2.1). We tested compression levels ×64 and ×256 for the compressed optical output intensity field in our experiments.

Table 1 shows the performances at ×64, ×256 compression levels for the tested HeLa dataset (section 5). For each compression level, performances are reported with and without the fine-tuning step. The corresponding qualitative results are shown in Figs. 4 and 5. All proposed methods outperformed all-optical baselines (B1, B2) with a significant margin in terms of SSIM (structural similarity index) and PSNR (peak signal-to-noise ratio) [42]. End-to-end fine-tuning showed a considerable improvement in the performance for all the cases. Our best method achieved PSNR= 29.76 dB, SSIM= 0.90 performance at ×64 compression, indicating that the proposed method is suitable for high-throughput QPM. Even at ×256 compression, the proposed method outperformed all-optical baselines by a considerable margin with PSNR= 27.62 dB and SSIM= 0.83. We further tested our approach by including a noise model with Poisson noise and
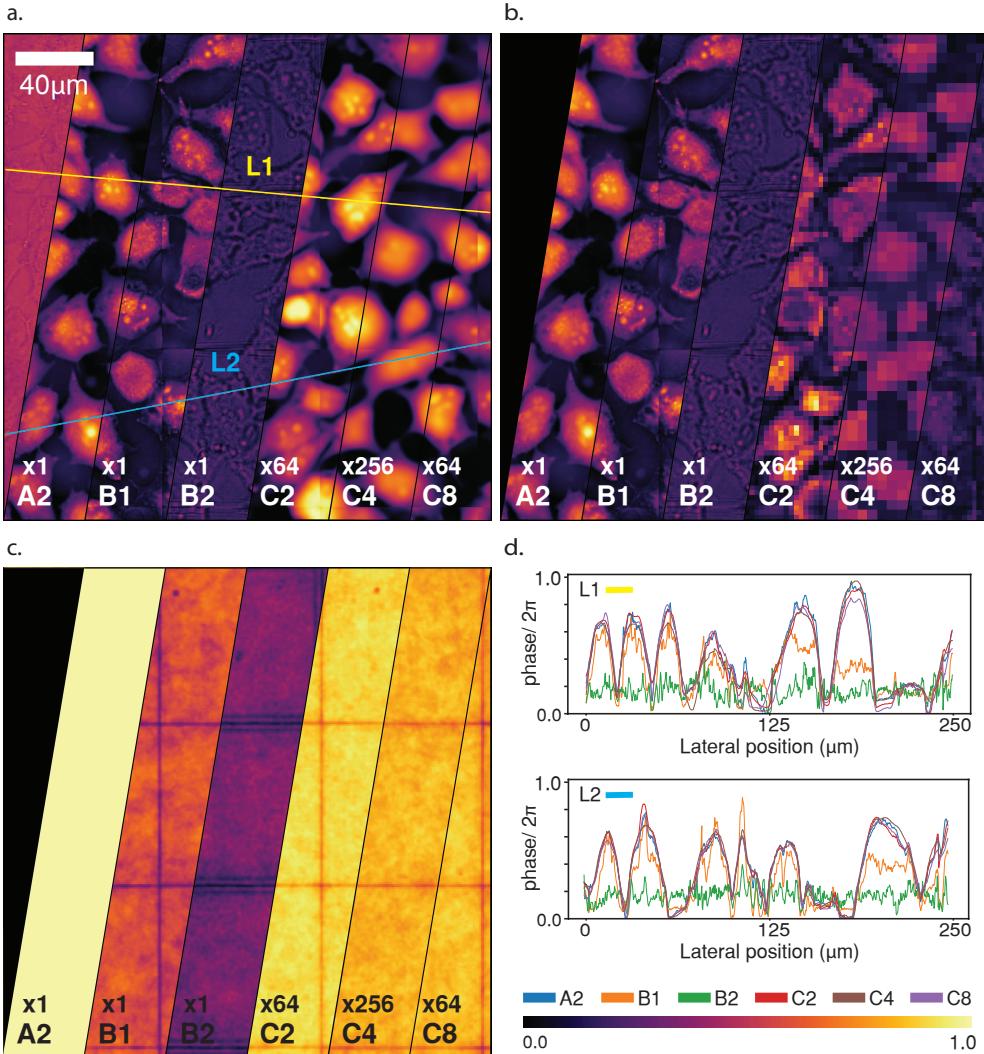
Fig. 5. **Performance comparison of best methods**: a) The phase reconstructions. b) Compressed intensity fields at the detector. c) SSIM maps of reconstructions. d) The resolving power of the phase reconstructions is shown. Phase of the input field (A2) of a full FoV from the test set, all-optical phase to intensity conversion results [37] (i.e. baselines) using LFF (B1), PhaseD2NN (B2), Phase reconstructions from our approach 1: LFF + SwinIR with ×64 compression, with fine-tuning (C2), LFF with ×256 compression with fine-tuning (C4), phase reconstructions from our approach 2: PhaseD2NN + SwinIR with ×64 compression, with fine-tuning + 5 optical layers (C8).

read noise [36]. We fine-tuned the best model (C2) with noise. A read noise with a standard deviation of 6.0 and a detector maximum photon count of 10000 were used. The proposed method with detector noise (E1) performed on par with the best model indicating that our LFF + SwinIR QPM is robust to real-world noise conditions. We further discuss the effect of the detector noise in the discussion (see section 3).

**PhaseD2NN + SwinIR.** Second, we tested a PhaseD2NN [37] as the optical network in the proposed end-to-end framework. Similar to the previous section, the SwinIR super-resolution network was used for reconstruction. We selected the operating range of the PhaseD2NN as the visible wavelength ($\lambda = 632.8$ nm). on the same HeLa cell dataset (see section 5).

The optical network consisted of 5 optical layers each having $256 \times 256$ optical neural grid. The size of each neuron was $\frac{\lambda}{2} \times \frac{\lambda}{2}$ (316.4 nm $\times$316.4 nm). Therefore, the size of the optical layer was 80.9984$\mu$m $\times$80.9984$\mu$m. Optical layers were separated with 3.373$\mu$m distance between each other. The distance between the input plane and the first optical layer was 3.373$\mu$m while the distance between the last optical layer and the detector plane was 5.904$\mu$m. Since the pixel size matched the PhaseD2NN neuron size, we could train the end-to-end network directly on the patch FoVs from the dataset. We followed the optimization criteria presented in section 2.1 for the end-to-end training. Notably, we observed that in step 1, PhaseD2NN training was not stable due to the large number of physical parameters with a larger grid size (e.g., $256 \times 256$). To increase the stability and gradient flow of this optimization step, we used a sub-optimization-schedule (shown in Supplementary algorithm S1). We compressed the output intensity from the optical model $\times 64$ to obtain a higher throughput.

Table 1 shows the performances for $\times 64$ compression level. We report the performances while selecting different layers of the PhaseD2NN as the output layer. The corresponding learned power value was given by the algorithm S1 (Supplementary). The final model with 5 layers was fine-tuned according to the proposed optimization steps. Similar to section 2.3, fine-tuning improved the performance. We explored different numbers of diffractive layers for the PhaseD2NN without the fine-tuning step and the results are presented in Table 1.

We performed further experiments with the 5 layer PhaseD2NN (C8 and E2). Our method achieved the best performance of PSNR= 27.24 dB, SSIM= 0.86 with $\times 64$ compression which was considerably higher than the all-optical baselines. Similar to the previous section, we injected detector noise with similar specifications of a maximum photon count of 10000 and detector read noise standard deviation of 6.0. The resultant performance with the detector noise (E2) was on par with the best model without the noise (C8). This indicates that our PhaseD2NN + SwinIR QPM is robust to real-world noise conditions.

## 3. Discussion

**Overall Comparison.** Fig. 5 presents the qualitative results for best-performing models. Fig. 5(d) shows that the proposed differentiable optical-electronic QPM systems have a higher resolving capability compared to all-optical baselines. Fig. 5(c) SSIM maps show how our methods perform for different regions of full field-of-view (FoV). Low SSIM in edges indicates that there is room to improve the proposed QPM just by refining the edges of generated patches. We observed that, even though both proposed methods: i.e, LFF with SwinIR; and PhaseD2NN with SwinIR outperformed the all-optical baselines by a considerable margin, the LFF-based method performed better than the PhaseD2NN-based one. Further studies are needed to investigate the reason for this behavior.

**Stability of PhaseD2NN Training.** We observed that the optimization step 1, i.e., all-optical reconstruction (see section 2.1), is not stable for the PhaseD2NN. We suspect that the reason for this instability is the large FoV (of $256 \times 256$) resulting in a large number of learnable parameters. To overcome this, we used a sub-optimization-schedule for the PhaseD2NN training motivated by progressive growing learning principles [43] (see algorithm S1 in the supplementary). Instead of training the PhaseD2NN in an end-to-end fashion, here we optimize the PhaseD2NN layer by layer progressively with the phase reconstruction loss. With this schedule, we could efficiently train the optical network. Even though one can argue that the proposed schedule leads to a sub-optimal solution, we achieved a sufficient performance for QPM [31] with this schedule.

Nevertheless, an interesting future direction is to explore more efficient methods to train large D2NNs.

**Resolving Power of Differentiable Optical-Electronic QPM.** We analyzed the resolving power of the proposed differentiable optical-electronic QPM in Fig. 4(d). The plots demonstrate the phase value variations along lines L1 and L2 on the patch FoV. Phase variations along L1 and L2 further show the superior resolving ability of the proposed methods for local and global features respectively. Furthermore, fine-tuning improved the resolving capability for most of the models.

**Effect of Photodetector Noise.** To further evaluate the behavior of the proposed method with detector noise, we evaluated the method with maximum photon counts of 100 and 10000, and read noise standard deviations of 4.0 and 6.0. Table 2 shows that differentiable optical-electronic QPM is robust to noise when the maximum photon count is 10000 (for most QPM applications high light conditions can be used). Even though the proposed method performs worse for lower photon counts (e.g., 100), it still performs far better than all-optical baselines. We further note that the performance degradation due to higher read noise is insignificant for larger photon counts (e.g., 10000).

**Compressibility limitations.** Last, we tested our LFF-based approach on a QPM dataset of tissue with much more complex features (see section 5). The goal of this experiment was to investigate the limitations of our approach at high compression levels. We observed that our method failed to reconstruct high-resolution features at both ×64 and ×256 (see supplementary Fig. S1). There could be two potential reasons for the subpar performance. It could be the case that the optical compressor cannot efficiently convert phase information to the latent intensity field at the detector. Alternatively, it could be the case that the reconstruction network is not capable of reconstructing highly compressed information from images with complex features. To investigate the latter we tested our reconstruction network on a simple resolution enhancement task on the same tissue dataset. As shown in supplementary Fig. S1, here too the reconstruction network failed. Thus we conclude that in our method, the compressibility is limited in the presence of complex features. Further studies are required to establish the compressibility bounds for data distributions of interest.

## 4. Conclusion

Quantitative phase microscopy (QPM) is an emerging label-free imaging modality with a wide range of biological and clinical applications. Recent advances in QPM are focused on developing fast instruments through better detectors and fast deep-learning-based inverse solvers. However, currently, the QPM throughput is fundamentally limited by the pixel throughput of the imaging detectors. Orthogonal to current advances, to improve QPM throughput beyond the hardware bottleneck, here we propose to use content-aware compressive data acquisition. Specifically, we utilize learnable optical front-ends to extract compressed phase features. A state-of-the-art transformer deep network then decodes the captured information to quantitatively reconstruct the phase image. The proposed pipeline inherently improves the imaging speed while achieving high-quality reconstructions. Moreover, the advances presented in this work can lead to similar developments in a wide range of label-free coherent imaging modalities such as photothermal, coherent anti-Stokes Raman scattering (CARS), and stimulated Raman scattering (SRS).

## 5. Methods

### 5.1. Datasets

In our numerical experiments, we used two datasets.

**HeLa Cell Dataset :**  We used a HeLa cell dataset [37] as the primary dataset for our experiments. We followed the sample preparation procedure explained in previous work [37]. The initial dataset contained 501 complex-valued images (i.e. detected FoVs). Each detected FoV was obtained by a camera with a $2304 \times 2304$ pixel grid where the pixel size was $6.5\ \mu$m $\times 6.5\ \mu$m. The light field from the specimen was magnified $\times 60$ before imaging onto the detector.

To pre-process the dataset, we first calculated the side length of the light fields before the magnification (= $\frac{2304\,pixels \times 6.5\ \mu\text{m}/pixel}{60}$ = 249.6 $\mu$m). Second, we calculated the number of 316.4 nm $\times$316.4 nm sized pixels in these light fields (= $round(\frac{249.6\ \mu\text{m}}{316.4\ \text{nm}/pixel}$ = $789 pixels$)). Finally, we resized the detected FoVs (i.e. $2304 \times 2304$ pixel grids) into $789 \times 789$ pixel grids. This resulted in the light field before the magnification with a pixel size of 316.4 nm $\times$316.4 nm. We refer to these FoVs as full FoVs. We obtained train and test sets by dividing the full FoV dataset into 401 and 100 sets. For the training of the proposed networks, we used $256 \times 256$ cropped patches (i.e. patch FoVs) from the full FoVs.

**Tissue Dataset:**  We also acquired a tissue dataset to further validate our observations and to derive an empirical upper bound for the results. We followed preparation, acquisition, and preprocessing procedures similar to HeLa cells, with a magnification of $\times 20$. There were 470 detected FoVs. Camera had $2367 \times 2367$ pixel grid where the pixel size was $6.5\ \mu$m $\times 6.5\ \mu$m. Side length of the light fields before the magnification was, $\frac{2304\,pixels \times 6.5\ \mu\text{m}/pixel}{20}$ = 748.8 $\mu$m). Number of 316.4 nm $\times$316.4 nm sized pixels in these light field was = $round(\frac{748.8\ \mu\text{m}}{316.4\ \text{nm}/pixel}$ = $2367 pixels$). We resized the detected FoVs (i.e. $2304 \times 2304$ pixel grids) into $2367 \times 2367$ pixel grids to match the pixel sizes of the light fields and the algorithm (full FoVs). The full FoV dataset was divided into 470 and 117, train and test sets.

### 5.2. Implementation Details

We implemented the proposed optical-electronic networks with Python version 3.6.13. We used auto differentiation in PyTorch [44] framework version 1.8.0 for the joint optimization/ training of the proposed optical-electronic networks. All experiments were conducted on a server with 12 Intel(R) Xeon(R) Platinum 8358 (2.60 GHz) CPU Cores, an NVIDIA A100 Graphics Processing Unit with 40 GB memory running on the Centos operating system.

We used batch size 32, learning rates of 0.1, 0.001 respectively for LFF and PhaseD2NN in the optimization stage 1. LFF was trained for 1500 epochs with multi-step learning rate scheduler [44] (milestones : [50, 400, 650, 1000, 1400], $\gamma = 0.1$). PhaseD2NN was trained for 1500 epochs after each optimizer initialization step in algorithm S1. For joint multi-layer optimizations in algorithm S1, a learning rate of 0.00005 was used for better stability. For the optimization stage 2, we followed similar training configurations used in SwinIR [38]. Lastly, for the final optimization stage (i.e. end-to-end fine-tuning), we fine-tuned the LFF + SwinIR and PhaseD2NN + SwinIR for 24000 and 3000 epochs with a learning rate of $5 \times 10^{-6}$, respectively. We used Adam [45] as the optimizer for all optimizations.

## 6. Backmatter

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

**Data availability.** Data underlying the results presented in this paper can be obtained from the authors upon reasonable request.

**Supplemental document.** See the supplemental document for supporting content.

## References

1. G. Popescu, T. Ikeda, R. R. Dasari, and M. S. Feld, "Diffraction phase microscopy for quantifying cell structure and dynamics," Opt. Lett. **31**, 775–777 (2006).
2. Y. Park, G. Popescu, K. Badizadegan, R. R. Dasari, and M. S. Feld, "Diffraction phase and fluorescence microscopy," Opt. Express **14**, 8263–8268 (2006).
3. C. Fang-Yen, S. Oh, Y. Park, W. Choi, S. Song, H. S. Seung, R. R. Dasari, and M. S. Feld, "Imaging voltage-dependent cell motions with heterodyne mach-zehnder phase microscopy," Opt. Lett. **32**, 1572–1574 (2007).
4. M. S. Amin, Y. Park, N. Lue, R. R. Dasari, K. Badizadegan, M. S. Feld, and G. Popescu, "Microrheology of red blood cell membranes using dynamic scattering microscopy," Opt. Express **15**, 17001–17009 (2007).
5. Y. Sung, N. Lue, B. Hamza, J. Martel, D. Irimia, R. R. Dasari, W. Choi, Z. Yaqoob, and P. So, "Three-dimensional holographic refractive-index measurement of continuously flowing cells in a microfluidic channel," Phys. Rev. Appl. **1**, 014002 (2014).
6. W. Choi, C. Fang-Yen, K. Badizadegan, S. Oh, N. Lue, R. R. Dasari, and M. S. Feld, "Tomographic phase microscopy," Nat Methods **4**, 717–719 (2007).
7. Y. K. Park, C. Depeursinge, and G. Popescu, "Quantitative phase imaging in biomedicine," Nat. Photonics **12**, 578–589 (2018).
8. Y. Jo, S. Park, J. Jung, J. Yoon, H. Joo, M.-H. Kim, S.-J. Kang, M. C. Choi, S. Y. Lee, and Y. Park, "Holographic deep learning for rapid optical screening of anthrax spores," Sci Adv **3**, e1700606 (2017).
9. D. Roitshtain, L. Wolbromsky, E. Bal, H. Greenspan, L. L. Satterwhite, and N. T. Shaked, "Quantitative phase microscopy spatial signatures of cancer cells," Cytom. Part A **91**, 482–493 (2017).
10. H. Majeed, A. Keikhosravi, M. Kandel, T. Nguyen, Y. Liu, A. Kajdacsy-Balla, K. Tangella, K. Eliceiri, and G. Popescu, "Quantitative histopathology of stained tissues using color spatial light interference microscopy (cslim)," Sci. Reports **9** (2019).
11. Y. Rivenson, T. Liu, Z. Wei, Y. Zhang, K. de Haan, and A. Ozcan, "PhaseStain: The digital staining of label-free quantitative phase microscopy images using deep learning," Light. Sci. Appl. **8** (2019).
12. Z. Wang, K. Tangella, A. Balla, and G. Popescu, "Tissue refractive index as marker of disease," J Biomed Opt **16**, 116017 (2011).
13. I. A. Cree, Z. Deans, M. J. L. Ligtenberg, N. Normanno, A. Edsjö, E. Rouleau, F. Solé, E. Thunnissen, W. Timens, E. Schuuring, S. Dequeker, S. Murray, M. Dietel, P. Groenen, J. H. Van Krieken, European Society of Pathology Task Force on Quality Assurance in Molecular Pathology, and Royal College of Pathologists, "Guidance for laboratories performing molecular pathology for cancer patients," J Clin Pathol **67**, 923–931 (2014).
14. M. Kandel, C. Hu, G. Naseri Kouzehgarani, E. Min, K. Sullivan, H. Kong, J. Li, D. Robson, M. Gillette, C. Best-Popescu, and G. Popescu, "Epi-illumination gradient light interference microscopy for imaging opaque structures," Nat. Commun. **10** (2019).
15. Y. Sung, W. Choi, N. Lue, R. R. Dasari, and Z. Yaqoob, "Stain-free quantification of chromosomes in live cells using regularized tomographic phase microscopy," PLoS One **7**, e49502 (2012).
16. Y. Sung, A. Tzur, S. Oh, W. Choi, V. Li, R. R. Dasari, Z. Yaqoob, and M. W. Kirschner, "Size homeostasis in adherent cells studied by synthetic phase microscopy," Proc. National Acad. Sci. United States Am. **110** (2013).
17. P. Hosseini, S. Z. Abidi, E. Du, D. P. Papageorgiou, Y. Choi, Y. Park, J. M. Higgins, G. J. Kato, S. Suresh, M. Dao, Z. Yaqoob, and P. T. C. So, "Cellular normoxic biophysical markers of hydroxyurea treatment in sickle cell disease," Proc Natl Acad Sci U S A **113**, 9527–9532 (2016).
18. V. Singh, Y. A. Yang, H. Yu, R. Kamm, Z. Yaqoob, and P. So, "Studying nucleic envelope and plasma membrane mechanics of eukaryotic cells using confocal reflectance interferometric microscopy," Nat. Commun. **10** (2019).
19. F. Zernike, "Observation of transparent objects," Physica pp. 974–986 (1942).
20. J. Glückstad, "Phase contrast image synthesis," Opt. Commun. **130**, 225–230 (1996).
21. N. Shibata, S. D. Findlay, Y. Kohno, H. Sawada, Y. Kondo, and Y. Ikuhara, "Differential phase-contrast microscopy at atomic resolution," Nat. Phys. **8**, 611–615 (2012).
22. Y. J. Jo, H. Cho, S. Y. Lee, G. Choi, G. Kim, H. S. Min, and Y. K. Park, "Quantitative phase imaging and artificial intelligence: A review," IEEE J. Sel. Top. Quantum Electron. **25** (2018).
23. K. Kim, K. S. Kim, H. Park, J. C. Ye, and Y. Park, "Real-time visualization of 3-d dynamic microscopic objects using optical diffraction tomography," Opt. Express **21**, 32269–32278 (2013).

24. J. Lim, K. Lee, K. H. Jin, S. Shin, S. Lee, Y. Park, and J. C. Ye, "Comparative study of iterative reconstruction algorithms for missing cone problems in optical diffraction tomography," Opt. Express **23**, 16933–16948 (2015).

25. Y. Sung, W. Choi, C. Fang-Yen, K. Badizadegan, R. R. Dasari, and M. S. Feld, "Optical diffraction tomography for high resolution live cell imaging," Opt. InfoBase Conf. Pap. **17**, 1977–1979 (2009).

26. T. Nguyen, V. Bui, and G. Nehmetallah, "Computational optical tomography using 3d deep convolutional neural networks (3d-dcnns)," Opt. Eng. **57** (2018).

27. J. Di, K. Wang, Y. Li, and J. Zhao, "Deep learning-based holographic reconstruction in digital holography," in *Imaging and Applied Optics Congress,* (Optica Publishing Group, 2020), p. HTu4B.2.

28. Y. Zhu, C. H. Yeung, and E. Y. Lam, "Digital holographic imaging and classification of microplastics using deep transfer learning," Appl. Opt. **60**, A38–A47 (2021).

29. K. Wang, Q. Kemao, J. Di, and J. Zhao, "Y4-net: A deep learning solution to one-shot dual-wavelength digital holographic reconstruction," Opt. Lett. **45**, 4220–4223 (2020).

30. H. Wang, M. Lyu, and G. Situ, "eholonet: A learning-based end-to-end approach for in-line digital holographic reconstruction," Opt. Express **26**, 22603–22614 (2018).

31. K. Wang, J. Dou, Q. Kemao, J. Di, and J. Zhao, "Y-net: A one-to-two deep learning framework for digital holographic reconstruction," Opt. Lett. **44**, 4765–4768 (2019).

32. M. R. Kellman, E. Bostan, N. A. Repina, and L. Waller, "Physics-based learned design: Optimized coded-illumination for quantitative phase imaging," IEEE Trans. on Comput. Imaging **5** (2019).

33. A. Matlock and L. Tian, "High-throughput, volumetric quantitative phase imaging with multiplexed intensity diffraction tomography," Biomed. Opt. Express **10**, 6432–6448 (2019).

34. P. Yeh and C. Gu, *Landmark Papers on Photorefractive Nonlinear Optics* (World Scientific, 1995).

35. V. Studer, J. Bobin, M. Chahid, H. S. Mousavi, E. Candes, and M. Dahan, "Compressive fluorescence microscopy for biological and hyperspectral imaging," Proc. National Acad. Sci. **109**, E1679–E1687 (2012).

36. U. Haputhanthri, A. Seeber, and D. Wadduwage, "Differentiable microscopy for content and task aware compressive fluorescence imaging," (2022).

37. K. Herath, U. Haputhanthri, R. Hettiarachchi, H. Kariyawasam, A. Ahmad, B. S. Ahluwalia, C. U. S. Edussooriya, and D. Wadduwage, "Differentiable microscopy designs an all optical quantitative phase microscope," (2022).

38. J. Liang, J. Cao, G. Sun, K. Zhang, L. V. Gool, and R. Timofte, "Swinir: Image restoration using swin transformer," 2021 IEEE/CVF Int. Conf. on Comput. Vis. Work. (ICCVW) pp. 1833–1844 (2021).

39. H. Gholamalinezhad and H. Khosravi, "Pooling methods in deep neural networks, a review," ArXiv **abs/2009.07485** (2020).

40. Z. Wang, A. Bovik, H. Sheikh, and E. Simoncelli, "Image quality assessment: From error visibility to structural similarity," IEEE Trans. on Image Process. **13**, 600–612 (2004).

41. Y. Wang, H. Yao, and S. Zhao, "Auto-encoder based dimensionality reduction," Neurocomputing **184**, 232–242 (2016). RoLoD: Robust Local Descriptors for Computer Vision 2014.

42. A. Horé and D. Ziou, "Image quality metrics: PSNR vs. SSIM," in *2010 20th International Conference on Pattern Recognition,* (2010), pp. 2366–2369.

43. T. Karras, T. Aila, S. Laine, and J. Lehtinen, "Progressive growing of gans for improved quality, stability, and variation," ArXiv **abs/1710.10196** (2018).

44. A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, "Pytorch: An imperative style, high-performance deep learning library," in *Advances in Neural Information Processing Systems 32,* H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, eds. (Curran Associates, Inc., 2019), pp. 8024–8035.

45. D. Kingma and J. Ba, "Adam: A method for stochastic optimization," Int. Conf. on Learn. Represent. (2014).

| Experiment | Detector Noise | Optical Net | Compression | Finetune | # layers | full FoV/ patch FoV Metrics PSNR | SSIM |
|---|---|---|---|---|---|---|---|
| B1 |  | All optical LFF |  |  |  | 16.1565/ 16.9761 | 0.5880/ 0.6008 |
| B2 |  | All optical PhaseD2NN |  |  |  | 12.3730/ 12.6631 | 0.3163/ 0.3320 |
| C1 | ✗ | LFF | 64 | ✗ | - | 23.8267/ 25.7840 | 0.8225/ 0.8278 |
| C2 |  |  |  | ✓ | - | **27.2579/ 29.7608** | **0.8967/ 0.9031** |
| C3 |  |  | 256 | ✗ | - | 22.5457/ 23.9536 | 0.7470/ 0.7548 |
| C4 |  |  |  | ✓ | - | 26.0003/ 27.6129 | 0.8223/ 0.8302 |
| C5 |  | Phase-D2NN | 64 | ✗ | 1 | 22.6495/ 23.8566 | 0.7808/ 0.7889 |
| C6 |  |  |  |  | 3 | 24.7560/ 26.0716 | 0.8224/ 0.8313 |
| C7 |  |  |  |  | 5 | 24.8015/ 26.0551 | 0.8107/ 0.8185 |
| C8 |  |  |  | ✓ | 5 | **25.8617/ 27.2449** | **0.8519/ 0.8602** |
| E1 | ✓ | LFF | 64 | ✓ | - | 27.3794/ 29.8110 | 0.8935/ 0.8998 |
| E2 |  | Phase-D2NN | 64 | ✓ | 5 | 25.7665/ 27.0651 | 0.8477/ 0.8558 |

Table 1. **Performance comparison**: Best results for optical feature extraction networks **LFF**, **PhaseD2NN** are highlighted. These best models are further fine-tuned end-to-end with the detector noise simulation (noise specifications of the detector: read noise standard deviation= 6.0, maximum photon count= 10000) to improve realisticity. We calculate the patch and full FoV metrics on the test patch FoVs and full FoVs respectively. We reconstruct the full FoVs by tiling the reconstructed patch FoVs.

| Optical Net | Noise Specifications | | full FoV/ patch FoV metrics | |
|---|---|---|---|---|
| | max. photon count | $\sigma_{read}$ | **PSNR** | **SSIM** |
| LFF | 100 | 4 | 23.4928/ 25.0266 | 0.777/ 0.7834 |
| | | 6 | 22.5004/ 24.1487 | 0.7606/ 0.7663 |
| | 10000 | 4 | **27.4122/ 29.8110** | **0.8935/ 0.8997** |
| | | 6 | **27.3794/ 29.8110** | **0.8935/ 0.8998** |
| Phase-D2NN | 100 | 4 | 17.8942/ 18.7526 | 0.6441/ 0.6502 |
| | | 6 | 16.8953/ 17.7111 | 0.6094/ 0.6158 |
| | 10000 | 4 | **25.7193/ 27.0456** | **0.8478/ 0.8559** |
| | | 6 | **25.7665/ 27.0651** | **0.8477/ 0.8558** |

Table 2. **Performance of our method for different detector noise conditions**: Our best models (C2, C8 in table 1) are further fine-tuned with the corresponding noise specifications.

# From Hours to Seconds: Towards 100x Faster Quantitative Phase Imaging via Differentiable Microscopy: Supplemental Document

## 1. SUB-OPTIMIZATION-SCHEDULE FOR STABLE ALL-OPTICAL PHASED2NN TRAINING

We observe that using PhaseD2NN as the optical network hinders the convergence of the optimization step-1 (section 2.1 in the manuscript). To improve the stability of this step, we consider a sub-optimization-schedule in algorithm S1.

**Algorithm S1.** Sub-optimization-schedule for stable all-optical PhaseD2NN training

---

**Data:** $P_X$
**Result:** $H_O^*(.)$ ;                    /* Optimal PhaseD2NN all-optical model */
$H_O(.) \leftarrow (f_n \circ f_{n-1} \circ ... \circ f_1)(.)$ ;           /* Define PhaseD2NN. $f_n$ is the $n^{th}$ layer */
$P_1, P_2, ., P_n \leftarrow 1.0$ ;          /* Define learnable power values for each layer */
**for** $i \in [1, n]$ **do**
  Initializing optimizers, schedulers
    **for** *each epoch* **do**
      **for** *each data batch* $x_{in} \sim P_X$ **do**
        $A_{out} = ||P_i \times (f_i \circ f_{i-1} \circ ... \circ f_1)(x_{in})||$
        $\mathcal{L}_\phi$ calculated according to section 2.1
        $f_i^*, P_i^* = \underset{f_i, P_i}{\mathrm{argmin}}(\mathcal{L}_\phi)$
      **end**
    **end**
  **if** $i \neq 1$ **then**
    Initializing optimizers, schedulers
      **for** *each epoch* **do**
        **for** *each data batch* $x_{in} \sim P_X$ **do**
          $A_{out} = ||P_i \times (f_i \circ f_{i-1} \circ ... \circ f_1)(x_{in})||$
          $\mathcal{L}_\phi$ calculated according to section 2.1
          $f_i^*, f_{i-1}^*, ..., f_1^*, P_i^* = \underset{f_i, f_{i-1}, ..., the f_1, P_i}{\mathrm{argmin}}(\mathcal{L}_\phi)$
        **end**
      **end**
    **end**
**end**

---

## 2. COMPARISON WITH EMPIRICAL UPPER BOUND

We further tested our method on a tissue dataset (section 5). All the hyper-parameters were similar to the HeLa experiments. We compared our method with a baseline we call *PhaseSR*. PhaseSR is a hypothetical phase reconstruction method that assumes perfect phase-to-intensity transformation. This removes the bottleneck of converting phase to intensity optically. Therefore only bottleneck the end-to-end pipeline has to handle is reconstructing missing information due to the downsampling of the field. The input to the super-resolution network was a downscaled phase distribution where the phase values were distributed in $[0, 2\pi]$. We considered $\times 64$ and $\times 256$ downscaling and compared the qualitative (Fig. S1) and quantitative (Table S1) results with our method. Due to the ideal phase-to-intensity transformation in PhaseSR, we considered it as the empirical upper bound for the results.

Poor performances in all-optical baselines indicate that phase retrieval of the tissue dataset is more challenging than that of the HeLa cell dataset. The proposed methods also exhibited lower performance on this dataset. Furthermore, as anticipated, PhaseSR demonstrated higher quantitative performance compared to the proposed method. However, the qualitative performance similarity between the proposed method and PhaseSR suggests that compression plays a more significant role in degrading performance than phase-to-intensity conversion.

| Experiment | Optical Net | Compression | Finetune | full FoV/ patch FoV | |
|---|---|---|---|---|---|
| | | | | **PSNR** | **SSIM** |
| B | All optical LFF | | | 11.5446/ 12.3976 | 0.0640/ 0.0632 |
| C1 | LFF | 64 | ✓ | 22.1981/ 24.2548 | 0.7134/ 0.7147 |
| C2 | | 256 | ✓ | 19.5154/ 21.1397 | 0.5555/ 0.5563 |
| D1 | PhaseSR | 64 | - | 24.8808/ 27.6686 | 0.7759/ 0.7735 |
| D2 | | 256 | - | 22.539/ 23.7128 | 0.6043/ 0.6036 |

**Table S1. Performance comparison on tissue dataset**: We compare the performance of our method with PhaseSR, for ×64 and ×256 compressions. Since PhaseSR assumes a perfect all-optical phase-to-intensity transformation, we treat it as the empirical upper bound for our results. Refer to Fig. S1 for qualitative comparison
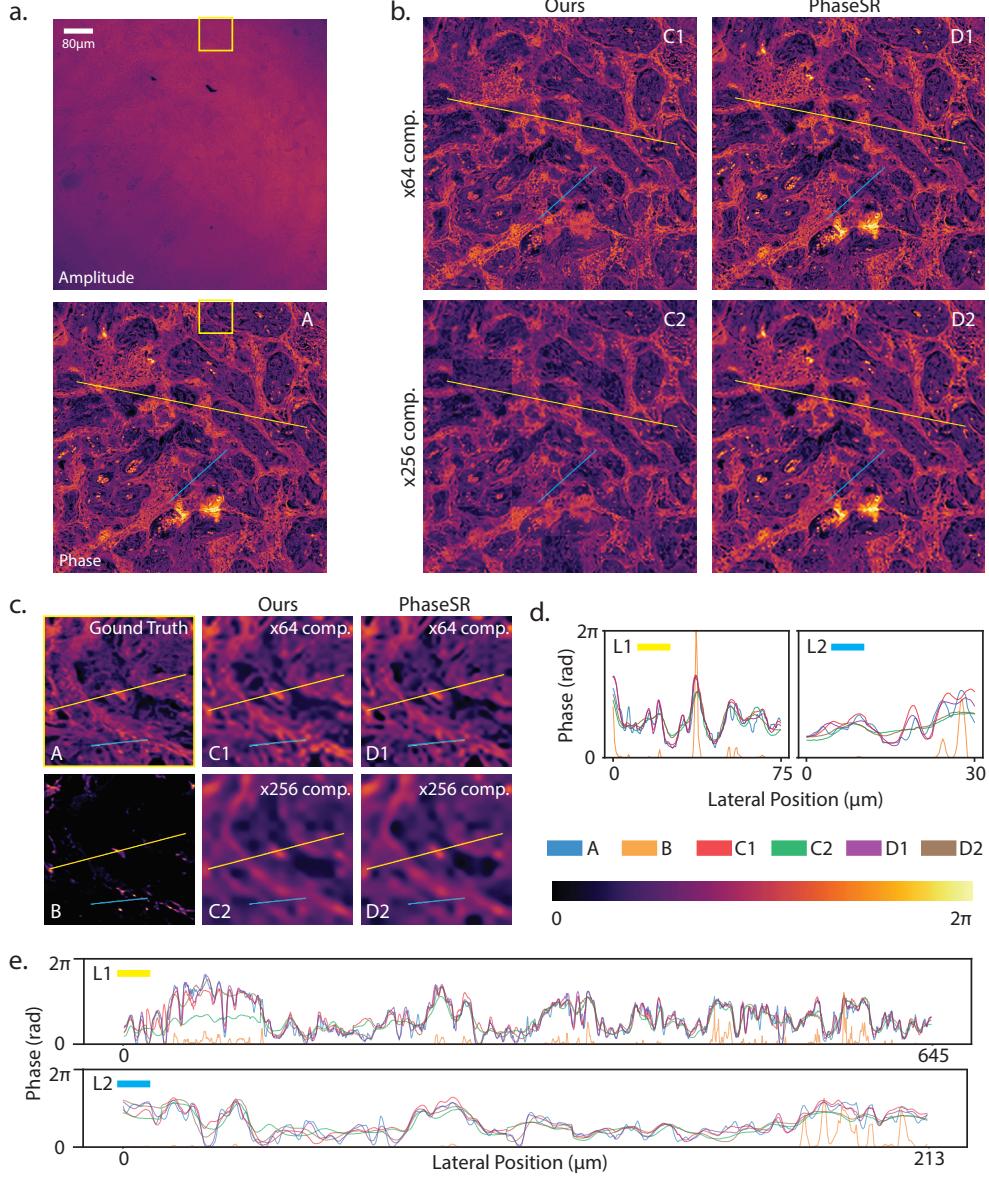
**Fig. S1. Performance comparison on QPM tissue dataset**: a) A sample from test dataset. b) Reconstructions from our method and PhaseSR for ×64 and ×256 compressions. PhaseSR assumes a perfect all-optical phase to intensity transformation. The reconstruction network takes a (×64 or ×256) downscaled ground truth phase as the input. Due to the perfect all-optical phase-to-intensity assumption, we treat PhaseSR as the empirical upper bound for our results. c) Magnified ground truth and corresponding results. Figure B shows the all-optical LFF reconstruction. d) Phase fluctuations along two lines on small field of views (C). e) Phase fluctuations along two lines on large field of views (B). Colors A-D2 are ground truth, all-optical LFF, ours with ×64 compression, ours with ×256 compression, PhaseSR with ×64 compression, PhaseSR with ×256 compression (please refer to table S1).