

Nonlinear Optimization Algorithm for Partially Coherent Phase Retrieval and Source Recovery

Jingshan Zhong, Lei Tian, Paroma Varma and Laura Waller

Abstract—We propose a new algorithm for recovering both complex field (phase and amplitude) and source distribution (illumination spatial coherence) from a stack of intensity images captured through focus. The joint recovery is formulated as a nonlinear least-square-error optimization problem, which is solved iteratively by a modified Gauss-Newton method. We derive the gradient and Hessian of the cost function and show that our second-order optimization approach outperforms previously proposed phase retrieval algorithms, for datasets taken with both coherent and partially coherent illumination. The method is validated experimentally in a commercial microscope with both Köhler illumination and a programmable LED dome.

Index Terms—Phase retrieval, Image reconstruction.

I. INTRODUCTION

Many algorithms have been proposed for reconstructing quantitative phase from through-focus intensity images [1]–[23]. Owing to their experimental simplicity, such methods have been widely adopted for phase imaging in existing imaging systems (e.g. optical microscopes [17]–[20], X-ray [5] and EUV Lithography [8]). This paper contributes both an improved phase retrieval algorithm and a method for jointly recovering the source pattern (spatial coherence of the illumination) and complex field, with high resolution in both.

Our modified Gauss-Newton nonlinear optimization method builds upon algorithms used in coherent phase retrieval [1], [3], [4], [24] and Ptychography [25]–[28]. Previous algorithms mainly employ gradient descent (e.g. Gerchberg-Saxton [21]), which is a *first-order* optimization method [25]. First-order methods use only the first derivative of the cost function for the complex field update at each iteration, whereas *second-order* optimization (e.g. Newton’s method) uses both the first and second derivatives for a better update and faster convergence [29]. Here, we derive the analytic expression for the second-order derivative (Hessian matrix) of our cost function and implement a second-order optimization approach.

In practice, computing the Hessian is prohibitively computational. Its size is the number of pixels squared (too large to store for multi-Megapixel images) and computational complexity scales with the cubic of the number of pixels. Newton’s method and some Quasi-Newton methods (e.g. Davidon-Fletcher-Powell (DFP), Broyden-Fletcher-Goldfarb-Shanno (BFGS)) [29] are computationally impractical since

J. Zhong, L. Tian and L. Waller are with the Department of Electrical Engineering and Computer Sciences at the University of California, Berkeley, USA.

P. Varma is with the Department of Electrical Engineering at Stanford University, USA.

E-mail: lwaller@alum.mit.edu
Website: www.laurawaller.com

they require storage or computation of the Hessian matrix or its inverse. The limited-memory Broyden-Fletcher-Goldfarb-Shanno method (L-BFGS) [30] circumvents these problems by approximating the Hessian from a few first-order gradients. Here, we propose a better second-order estimate by showing how the necessary quantity - the product of the Hessian with a vector - can be computed efficiently using FFTs and a conjugate gradient method [29]. The result is a stable, accurate and fast phase retrieval algorithm for coherent or partially coherent illumination and known or unknown sources.

Partially coherent illumination is particularly useful in brightfield microscopes and lithography projection imaging systems, providing advantages of reduced speckle noise, greater light throughput and up to $2\times$ better resolution [6], [31]. Many phase reconstruction algorithms assume a coherent imaging model (e.g. illumination by an on-axis plane wave) in which the complex field at focus is the sample’s transmittance function. With partially coherent systems, the defocused intensity will also depend on the illumination coherence. Using a coherent model with partially coherent data will incur model mis-match errors that can significantly degrade the result [12], [32]. Only recently have phase retrieval algorithms been developed to explicitly account for partial coherence [9]–[12]. Interestingly, the mathematical formulation for partially coherent phase imaging is very similar to that of joint estimation of object and aberration parameters [14], [23], [33].

The first step is to account for partially coherent illumination in the forward model. A full description of partial coherence requires a higher-dimensional function (e.g. mutual intensity) [34], which is impractically large. However, partial coherence generated by an incoherent 2D extended source can be represented by a set of mutually incoherent point sources. Here, we consider two types of partially coherent illumination: Köhler illumination (used in most microscopes) and a domed LED array [35]. We approximate the sum over source points as a convolution of the coherent intensity with the source pattern scaled by defocus distance [36], [37], such that our joint phase and source recovery algorithm becomes practical and efficient. Partially coherent phase recovery has been demonstrated already for TIE [9], [38], iterative methods [11] and the Kalman filter method [10], [12] using a similar model.

The complex field and the source affect through-focus intensity in different ways, which is why it is possible to solve for both. Recovering three 2D distributions (phase, amplitude and source) from a single focus stack leads to a large-scale nonlinear and non-convex optimization problem that is difficult to solve directly. Thus, we take an alternating minimization approach by splitting the problem into two parts:

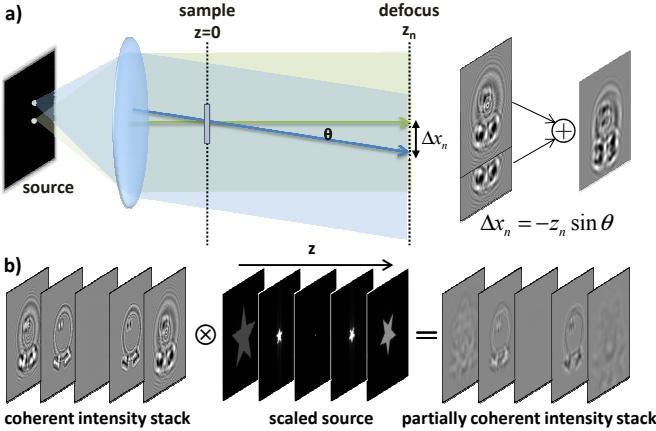


Fig. 1. Convolution model for defocused partially coherent intensity in K\"{o}hler geometry. a) Each point on the source generates a plane wave at the sample arriving at a particular angle. The on-axis source point's coherent diffracted intensity after propagation (green) is the same as the off-axis point source's (blue), but laterally shifted by an amount determined by the defocus distance and illumination angle. Measured intensity is the sum over all source points, which can be expressed as a convolution. b) Each image in a partially coherent focus stack is thus approximated as the 2D convolution of the coherent intensity image with a scaled version of the source pattern.

phase recovery and *source recovery*. Our algorithm is not guaranteed to find the global minimum, but converges well in practice when appropriate defocus distances are used. It works for *coherent* or *partially coherent* illumination where the source is either *known* or *unknown*.

II. THEORY

Our goal is to recover both the complex field of the sample and the source pattern (which defines illumination spatial coherence) from a stack of intensity images captured through-focus. We first introduce a *forward model*, which mathematically describes the physical process. Next, we formulate the *inverse problem* and propose a nonlinear optimization method that computationally estimates phase, amplitude and the source pattern from the measurements.

A. Forward model

We start by mathematically describing the forward model. Figure 1(a) shows a diagram of the experimental setup for K\"{o}hler geometry: an incoherent 2D source is placed in the Fourier plane of the sample, then intensity images are measured at various defocus distances. In our microscope, the sample's complex field is magnified and relayed to the detector plane by a 4f system (not shown).

The source is assumed to be incoherent, so it can be decomposed into a set of mutually incoherent point sources. In both K\"{o}hler geometry and our domed LED illuminator, each source point illuminates the sample with a plane wave at a different tilt angle. The measured intensity at a given defocus distance is thus the incoherent sum of the defocused intensity patterns from each source point. We define the source pattern, $S(\sin(\theta_x), \sin(\theta_y))$, in terms of illumination angles which correspond to numerical aperture (NA) in the Fourier plane. Through focus, the diffracted field from an off-axis

source point passing through a thin complex object will shift laterally across the field of view by Δx , which is proportional to the tilt angle and defocus distance z (see Fig. 1):

$$\Delta x = -z \sin(\theta_x). \quad (1)$$

This shift is entirely predicted by ray optics, though our model includes wave optics via coherent propagation effects.

Given that the defocused partially coherent intensity is the sum of shifted versions of the coherent intensity, the result can be written as a convolution between the scaled source pattern and the coherent diffracted intensity image. From Eq. (1), an off-axis point source at $\sin(\theta_x)$ causes a shift amount of $-z \sin(\theta_x)$ at defocus distance z . Therefore, the scaled source in the convolution model $S_z(x, y)$ is obtained from the original source distribution function $S(\sin(\theta_x), \sin(\theta_y))$ by replacing $\sin(\theta_x)$ with $-\frac{x}{z}$ and $\sin(\theta_y)$ with $-\frac{y}{z}$:

$$S_z(x, y) = \frac{1}{z^2} S\left(-\frac{x}{z}, -\frac{y}{z}\right), \quad (2)$$

where $\frac{1}{z^2}$ is a normalization factor that keeps total energy constant. The source $S_z(x, y)$ scales with defocus distance z . Mathematically, the partially coherent intensity $I^{PC}(x, y; z)$ at defocus distance z is a 2D convolution of the coherent intensity $I^C(x, y; z)$ with the scaled source pattern [36], [37]:

$$I^{PC}(x, y; z) = I^C(x, y; z) \otimes S_z(x, y), \quad (3)$$

where \otimes denotes convolution. The coherent intensity $I^C(x, y; z)$ is the stack of images that would be produced by illuminating the sample with an on-axis plane wave, ($\sin(\theta_x) = 0, \sin(\theta_y) = 0$). For small defocus distances the scaled source is small, so partial coherence effects will be minimal [37]. For large defocus distances the scaled source is large, so the intensity images will be blurred, as expected.

The partially coherent forward model in Eq. (3) accurately predicts through-focus intensity, but ignores the finite numerical aperture of the microscope's 4f system. In reality, each illumination angle (source point) causes a shift of the sample spectrum in the Fourier (pupil) plane of the microscope. Hence, each source point effectively sees a different shifted pupil, breaking the convolution model. It is straightforward to incorporate this differential pupil clipping in our forward model by implementing a shifted pupil cutoff for each source point. However, it is computationally intensive, so we use the convolution approximation to speed up computation.

We first introduce some notation to describe our forward model with matrices. Column vectors denoted in bold by \mathbf{I}_n^{PC} , \mathbf{I}_n^C and \mathbf{s}_n are raster-scanned from their 2D counterparts, $I^{PC}(x, y; z_n)$, $I^C(x, y; z_n)$ and $S_{z_n}(x, y)$. n is the index of the image measured at defocus distance z_n . The sample's complex field $A(x, y, z = 0)$ is solved for via its Fourier transform, denoted by vector \mathbf{b} . Defocused coherent intensity can thus be written as the square modulus of the propagated field:

$$\mathbf{I}_n^C = |\mathbf{K}^H \mathbf{H}_n \mathbf{b}|^2. \quad (4)$$

\mathbf{H}_n is a diagonal matrix describing wave propagation by distance z_n , the superscript H (in \mathbf{K}^H) represents matrix Hermitian and \mathbf{K} is the discrete Fourier transform matrix (with its Hermitian \mathbf{K}^H being the inverse Fourier transform matrix).

$|\cdot|^2$ takes absolute square for each element of the vector inside. The diagonal entries of \mathbf{H}_n are raster-scanned from analytical values given by either the Fresnel or the angular spectrum propagation kernel in frequency domain.

We define a vector \mathbf{s} as a raster-scanned column vector of the source, scaled for the largest defocus distance, z_{max} . The vector \mathbf{s}_n is related to \mathbf{s} by scaling it by a factor of $\frac{z_{max}}{z_n}$ (Eq. (3)). In computation, the scaling of the source is realized by four steps: zero-padding in real-space, doing a Fourier transform to the zero-padded source, cropping the central part of spectrum in frequency domain and performing an inverse Fourier transform back to real-space. This scaling procedure can be computed efficiently by FFT:

$$\mathbf{s}_n = \mathbf{K}^H \mathbf{O}_n \mathbf{K} \mathbf{Z}_n \mathbf{s}, \quad (5)$$

where \mathbf{Z}_n represents zero-padding, \mathbf{O}_n cropping in frequency domain, and \mathbf{K} and \mathbf{K}^H are the Fourier transform and its inverse, respectively. By grouping the first three steps of the scaling procedure, we have:

$$\mathbf{s}_n = \mathbf{K}^H \mathbf{P}_n \mathbf{s}, \quad (6)$$

where $\mathbf{P}_n = \mathbf{O}_n \mathbf{K} \mathbf{Z}_n$. Therefore, $\mathbf{P}_n \mathbf{s}$ is the Fourier transform of the scaled source \mathbf{s}_n . Since the scaling here is performed with discrete computation, it may have numerical errors that manifest as Gibbs ringing at the edge of the source due to the cropping in frequency domain.

The forward model of the joint recovery problem is thus a discretized version of Eq. (3):

$$\begin{aligned} \mathbf{I}_n^{PC} &= \mathbf{I}_n^C \otimes \mathbf{s}_n \\ &= |\mathbf{K}^H \mathbf{H}_n \mathbf{b}|^2 \otimes (\mathbf{K}^H \mathbf{P}_n \mathbf{s}). \end{aligned} \quad (7)$$

B. Inverse problem

Our forward model gives insight into the joint complex field and source recovery problem. Consider the case of a *known complex field*, which means that the coherent intensity stack $I^C(x, y; z)$ is also known. In this case, the source retrieval problem is akin to a deconvolution problem in which we have multiple known images convolved by scaled versions of the same source – a solvable and well-determined problem. Alternatively, consider the case of a *known source*, whose scaled pattern can be deconvolved from each defocused intensity image to recover the coherent intensity stack, resulting in the traditional (coherent) phase retrieval problem. Here, we describe how to solve each of these problems efficiently, as well as the joint recovery problem.

We formulate the joint inverse problem as a nonlinear least-square optimization which aims to minimize the cost function in terms of the complex field, \mathbf{b} , and the source, \mathbf{s} :

$$\begin{aligned} \min_{\mathbf{b}, \mathbf{s}} Err(\mathbf{b}, \mathbf{s}) &= \frac{1}{2} \sum_{n=1}^{N_{im}} \left\| \mathbf{I}_n^{PC} - \mathbf{I}_n^C \otimes \mathbf{s}_n \right\|_2^2 \\ &= \frac{1}{2} \sum_{n=1}^{N_{im}} \left\| \mathbf{I}_n^{PC} - |\mathbf{K}^H \mathbf{H}_n \mathbf{b}|^2 \otimes (\mathbf{K}^H \mathbf{P}_n \mathbf{s}) \right\|_2^2, \end{aligned} \quad (8)$$

where $\|\cdot\|$ represents ℓ_2 norm, N_{im} is the total number of measured images, \mathbf{I}_n^{PC} is the measured partially coherent intensity and $|\mathbf{K}^H \mathbf{H}_n \mathbf{b}|^2 \otimes (\mathbf{K}^H \mathbf{P}_n \mathbf{s})$ is the estimated intensity predicted from the current estimates of \mathbf{b} and \mathbf{s} . The goal is to find an estimate of \mathbf{b} and \mathbf{s} that minimizes the mean-square error between the measured and predicted intensity stacks.

It is difficult to obtain a direct solution to the joint recovery problem in Eq. (8), since the optimization problem is nonlinear, non-convex and large-scale. In the following sections, we break the problem into two smaller sub-problems. First, we consider the *phase recovery problem* when the source is known. Next, we consider the *source recovery problem* when the complex field is known. Finally, we combine these two to recover both \mathbf{b} and \mathbf{s} using an alternating update scheme.

C. Phase recovery with known source

We start with the simplified *known source* case, in which the illumination coherence (source pattern) is known or measured *a priori*. The optimization in Eq. (8) becomes a problem of phase recovery from a partially coherent intensity stack. If the source \mathbf{s} is known, then $\mathbf{K}^H \mathbf{P}_n \mathbf{s}$ (equivalent to \mathbf{s}_n) is also known and Eq. (8) can be reformulated as:

$$\min_{\mathbf{b}} F(\mathbf{b}) = \frac{1}{2} \sum_{n=1}^{N_{im}} \left\| \mathbf{I}_n - \mathbf{C}_n |\mathbf{K}^H \mathbf{H}_n \mathbf{b}|^2 \right\|_2^2, \quad (9)$$

where \mathbf{C}_n is a matrix describing the convolution of the scaled source \mathbf{s}_n . By imposing periodic boundary conditions, we have

$$\mathbf{C}_n = \mathbf{K}^H \text{diag}(\mathbf{P}_n \mathbf{s}) \mathbf{K}, \quad (10)$$

where $\mathbf{P}_n \mathbf{s}$ is the Fourier transform of \mathbf{s}_n (see Eq. (6)), and $\text{diag}(\cdot)$ gives a square diagonal matrix with the elements of the vector on the main diagonal. In the coherent case, the matrix \mathbf{C}_n is simply the identity matrix. Note that there is no need to compute \mathbf{C}_n explicitly; rather, we use multiplication in Fourier space to implement the convolution using FFTs.

For comparison, we evaluate various nonlinear optimization algorithms (both *first-order* and *second-order*) for solving Eq. (9). First, we calculate the gradient and Hessian of the cost function in Eq. (9) using the theory of \mathbb{CR} calculus [39], [40]. In \mathbb{CR} calculus, the complex variable \mathbf{b} is augmented as $[\mathbf{b} \ \bar{\mathbf{b}}]^T$ for the convenience of calculating the gradient and Hessian, in which $\bar{\mathbf{b}}$ is the complex conjugate of \mathbf{b} . In Appendix A, we derive the gradient, $\nabla F(\mathbf{b}, \bar{\mathbf{b}})$ in Eq. (29), and an approximation of the Hessian, $\nabla^2 F(\mathbf{b}, \bar{\mathbf{b}})$ in Eq. (30).

To start the iteration procedure, we initialize the complex field \mathbf{b}^0 with the measured intensity at focus. At the $k+1$ iteration, the general equation for updating the estimate is:

$$\mathbf{b}^{k+1} = \mathbf{b}^k + \alpha^k \mathbf{d}^k, \quad (11)$$

where α^k is a learning step size found by a backtracking line search algorithm [29], [41] and \mathbf{d}^k is the search direction. The line search algorithm gives an appropriate step size to guarantee that the cost function decreases after each iteration.

Next, we evaluate methods to choose the search direction, which is crucial for determining convergence speed. In the (*first-order*) gradient descent method, the search direction is

simply the gradient of the cost function. In *second-order* methods, the Hessian matrix is used to improve convergence speed. In (*second-order*) quasi-Newton methods or inexact Newton methods [29], the search direction is obtained by solving a linear system of the Hessian and the gradient to obtain the search direction \mathbf{d} :

$$\nabla^2 F(\mathbf{b}, \bar{\mathbf{b}}) \begin{bmatrix} \mathbf{d} \\ \bar{\mathbf{d}} \end{bmatrix} = -\nabla F(\mathbf{b}, \bar{\mathbf{b}}). \quad (12)$$

The block components of the Hessian matrix are Hermitian to each other (see Eq. (30)). From Eq. (29), the top half and bottom half of the gradient are complex conjugate. Therefore, we can express the solution of the linear equations as $[\mathbf{d}^T \bar{\mathbf{d}}^T]^T$, of which the top and bottom half are complex conjugate.

Since the Hessian matrix in our problem is large (size of N^2 , where N is the number of pixels in one image), it is difficult to solve Eq. (12) directly. In the L-BFGS [30] method, the linear equations are solved with an approximate inverse of the Hessian matrix that is built from the estimates and gradients of previous iterations. The L-BFGS method is computationally low-cost, with the optimal number of estimates and gradients used in the approximation being problem-dependent [29]. Despite being a crude approximation to the Hessian, we will show that L-BFGS can reduce the number of iterations needed and speed up runtime, compared to gradient descent.

Here, we propose a modified Gauss-Newton method (or line search Newton-CG [29]) for iteratively solving the linear Eq. (12) with a conjugate gradient method (Chapter 5 of [29]). We compute the multiplication of the matrix with a vector (estimates) at each iteration, which acts as the main computational overhead. The matrix $\nabla^2 F(\mathbf{b}, \bar{\mathbf{b}})$ consists of the DFT matrix \mathbf{K} and circulant matrix \mathbf{C}_n which denotes convolution (see Eq. (30)). So the multiplication of the matrix $\nabla^2 F(\mathbf{b}, \bar{\mathbf{b}})$ with a vector is computed by FFT, which has good computational complexity of $\mathcal{O}(N \log N)$.

In order to obtain the search direction \mathbf{d}^k for the estimate \mathbf{b}^k in the update procedure of the complex field, the linear equations with $\nabla^2 F(\mathbf{b}, \bar{\mathbf{b}})$ and $\nabla F(\mathbf{b}, \bar{\mathbf{b}})$ are solved at the current estimate \mathbf{b}^k . The resulting \mathbf{d}^k (top half of the linear system solution) provides the search direction in Eq. (11). Since the accuracy of this sub-procedure depends on the number of iterations used, there is a tradeoff between computational time and convergence speed of the overall algorithm. By directly using the analytic expression of the Hessian, our algorithm outperforms gradient descent and L-BFGS in terms of *both* number of iterations *and* runtime. Hence, it presents an improved phase retrieval method for coherent phase-from-defocus problems, as well as partially coherent.

D. Source estimation with known phase

Next, we consider another simplified problem: recovering an unknown source \mathbf{s} given a priori knowledge of the complex field \mathbf{b} . The coherent intensity stack \mathbf{I}_n^C in Eq. (8) is fully defined at all defocus steps by propagating \mathbf{b} (Eq. (4)). Hence, solving for the source is essentially an over-constrained deconvolution problem, with the source scaling in a known

way at each defocus distance. Reformulating the optimization problem in Eq. (8) into a source estimation problem gives:

$$\min_{\mathbf{s}} E(\mathbf{s}) = \frac{1}{2} \sum_{n=1}^{N_{im}} \left\| \mathbf{I}_n^{PC} - \mathbf{I}_n^C \otimes (\mathbf{K}^H \mathbf{P}_n \mathbf{s}) \right\|_2^2. \quad (13)$$

Exploiting the Fourier convolution property, Eq. (13) becomes:

$$\min_{\mathbf{s}} E(\mathbf{s}) = \frac{1}{2} \sum_{n=1}^{N_{im}} \left\| \mathbf{I}_n^{PC} - \mathbf{K}^H \text{diag}(\tilde{\mathbf{I}}_n^C) \mathbf{P}_n \mathbf{s} \right\|_2^2, \quad (14)$$

where $\tilde{\mathbf{I}}_n^C$ is the Fourier transform of \mathbf{I}_n^C .

Minimizing the error function $E(\mathbf{s})$ over \mathbf{s} is a quadratic optimization problem. By setting the first derivative to 0, the optimal solution is solved via the linear equations:

$$\left(\sum_{n=1}^N \mathbf{P}_n^H \text{diag}(|\tilde{\mathbf{I}}_n^C|^2) \mathbf{P}_n \right) \mathbf{s} = \sum_{n=1}^N \mathbf{P}_n^H \text{diag}(\tilde{\mathbf{I}}_n^C)^H \mathbf{K} \mathbf{I}_n^{PC}. \quad (15)$$

The matrix on the left side of Eq. (15) has size $N \times N$, where N is the total number of pixels in each image. Directly inverting this large matrix is computationally expensive, given large N (e.g. $N = 10^6$ for our images). However, multiplication of the large matrix with a vector can be efficiently computed by FFT. By exploiting fast matrix-vector multiplication, we solve the linear equations iteratively using the linear conjugate gradient method (Chapter 5 of [29]).

E. Joint estimation of both phase and source

Finally, we combine the methods in Sections II-C and II-D to solve the joint optimization problem (Eq. (8)). We propose an iterative scheme to alternately update \mathbf{b} and \mathbf{s} , as shown in Fig. 2. After initializing \mathbf{b}^0 and \mathbf{s}^0 , we fix our estimate of the complex field \mathbf{b}^k and solve for the source \mathbf{s}^{k+1} according to Section II-D, by minimizing the cost function $Err(\mathbf{b} = \mathbf{b}^k, \mathbf{s})$ over \mathbf{s} . Then, we take that estimate of the source \mathbf{s}^{k+1} and solve for the complex field \mathbf{b}^{k+1} according to Section II-C, by minimizing the cost function $Err(\mathbf{b}, \mathbf{s} = \mathbf{s}^{k+1})$. These new estimates of the complex field and source are then used as the starting point for another iteration of the algorithm. It can be shown that the error will reduce or remain unchanged with each iteration:

$$Err(\mathbf{b}^k, \mathbf{s}^k) \geq Err(\mathbf{b}^k, \mathbf{s}^{k+1}) \geq Err(\mathbf{b}^{k+1}, \mathbf{s}^{k+1}). \quad (16)$$

Once the cost function reaches a minimal value, the iterations stop. Although we cannot prove that our method will always achieve a global minimum (it is possible to get stuck in a local minimum), we show in simulation that the method provides a good estimation of both source and complex field (Section IV-A). Practically, the *second-order* optimization method for complex field recovery is important – particularly for joint recovery and large-scale nonlinear problems, *first-order* methods tend to be slow and get stuck in local minima [42].

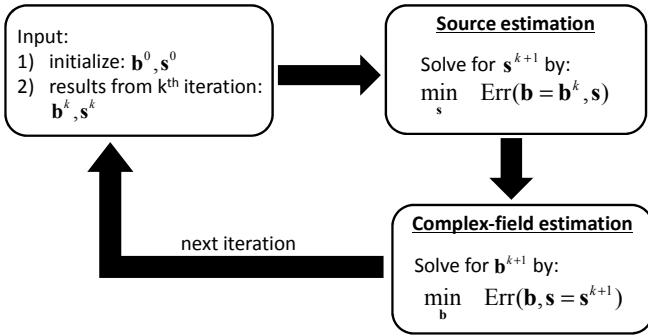


Fig. 2. The joint complex field and source recovery problem is solved iteratively using an alternating update scheme.

F. Choice of defocus distances

When jointly estimating many parameters from a single image stack, one must make sure that the data contains sufficient *redundancy* and *diversity*. We aim to recover three 2D distributions (amplitude, phase and source pattern), each having millions of pixels, so we must capture a minimum of 3 images [43]. Of course, more images will improve the result – in practice we find that 5-7 focus steps gives good results.

The choice of defocus steps (z planes) is important. For (coherent) phase retrieval, the ideal measurement planes depend on the object itself [44]. A weak-object linearization of the problem leads to an object-independent scheme that is approximately ideal: exponentially-spaced focus steps, with more images near focus than at further distances [19]. High spatial frequency phase information is contained in the images at small z distances and low spatial frequency phase information comes from large z distances.

For source recovery, the situation is reversed – high spatial frequency source information comes from large z distances and low spatial frequency source information comes from small z distances. This can be understood from the scaling of the source with defocus in Eq. (13). To quantify, suppose u_{max} is the highest frequency measured in the intensity images. The maximum spatial frequency of the source that is captured by the intensity stack is then zu_{max} ; the measured frequency range $[-zu_{max}, zu_{max}]$ shrinks linearly as z goes zero. In our optimization, we only attempt to recover the frequencies of the source which are measured by our focus stack.

The end result is that the phase and source recovery problems have competing tradeoffs in the choice of defocus distances. Because our images are generally captured relatively close to focus, high spatial frequencies of the source and low spatial frequencies of the phase are more sensitive to noise. We choose our z distances to have a wide range of defocus values, spaced exponentially near focus for good phase retrieval and linearly at large defocus distances for good source retrieval.

III. RESULTS FOR COHERENT PHASE RETRIEVAL

To demonstrate the advantages of our *second-order* algorithm for phase retrieval, we first look at the simplified case of complex field retrieval in a coherent microscope (the source is known to be a point source). We compare multiple algorithms:

- Gerchberg-Saxton (GS): global and sequential
- gradient descent (GD): small step size and large step size
- sparse Kalman filter (Sparse KF)
- L-BFGS
- our proposed modified Gauss-Newton method

Gerchberg-Saxton [21] is the most widely applied phase retrieval method due to its simplicity of implementation. One iteratively propagates the complex field estimate between defocused planes and updates the amplitude estimate with the square root of the measured intensity at each plane. There are two kinds of GS algorithm: sequential and global. The sequential GS algorithm [16], [45] imposes the constraints of the measured intensity images one-by-one as it propagates the complex field estimate back and forth through focus. The global GS algorithm instead updates the in-focus complex field estimate by taking an average of all of the estimates that result from imposing intensity constraints at each plane individually [25], [26]. The Sparse KF method [20], [46] takes a statistical approach, using the covariance matrix in order to implement pixel-wise noise filtering. Hence, it performs well in high noise, which is not the case here. Gradient descent, L-BFGS and Gauss-Newton all attempt to minimize the cost function in Eq. (9) for the case of C_n being the identity matrix (coherent phase retrieval). To classify the algorithms: GS and GD are *first-order* nonlinear optimization methods; L-BGFS and Gauss-Newton are *second-order* nonlinear optimization methods; and Sparse KF is a statistical method.

Though the performance of some algorithms may be tweaked by various parameters, we have attempted to provide a fair comparison. Since the convergence of GD is sensitive to the initialization of the step size in the line search algorithm, we consider two initialization cases for the GD algorithm: small step size (1) and large step size (the number of pixels in one image). The second-order methods, L-BFGS and modified Gauss Newton, are initialized with step size of 1 (default). L-BFGS stores gradients of the previous 100 estimates for the Hessian approximation, and our modified Gauss-Newton method uses 50 iterations to solve the Hessian linear equations.

A. Simulations

We start with simulations, where the ground truth solution is known. The complex field to be recovered is a starfish amplitude distribution and a flower phase distribution, with a phase range of π radians (Fig. 3(a)). We simulate a coherent stack of 7 defocused images by propagating the complex field symmetrically through focus to z distances: $0, \pm 38.3, \pm 80.2, \pm 167.7 \mu m$ (Fig. 3(b)). The illumination has a wavelength of $632.8 nm$, one single image has 301×301 pixels (pixel size $1 \mu m$), and all intensity images are polluted by Gaussian noise of variance 0.002 (the mean square value of the intensity stack is 0.71). Each algorithm uses the same dataset and the same initialization (square root of the measured intensity at focus).

The recovered phase from each algorithm is shown in Fig. 3. Failed results for the sequential GS method (which does not converge correctly) and the GD method with small step size (which updates overly slowly) are not shown. Thus, the GS

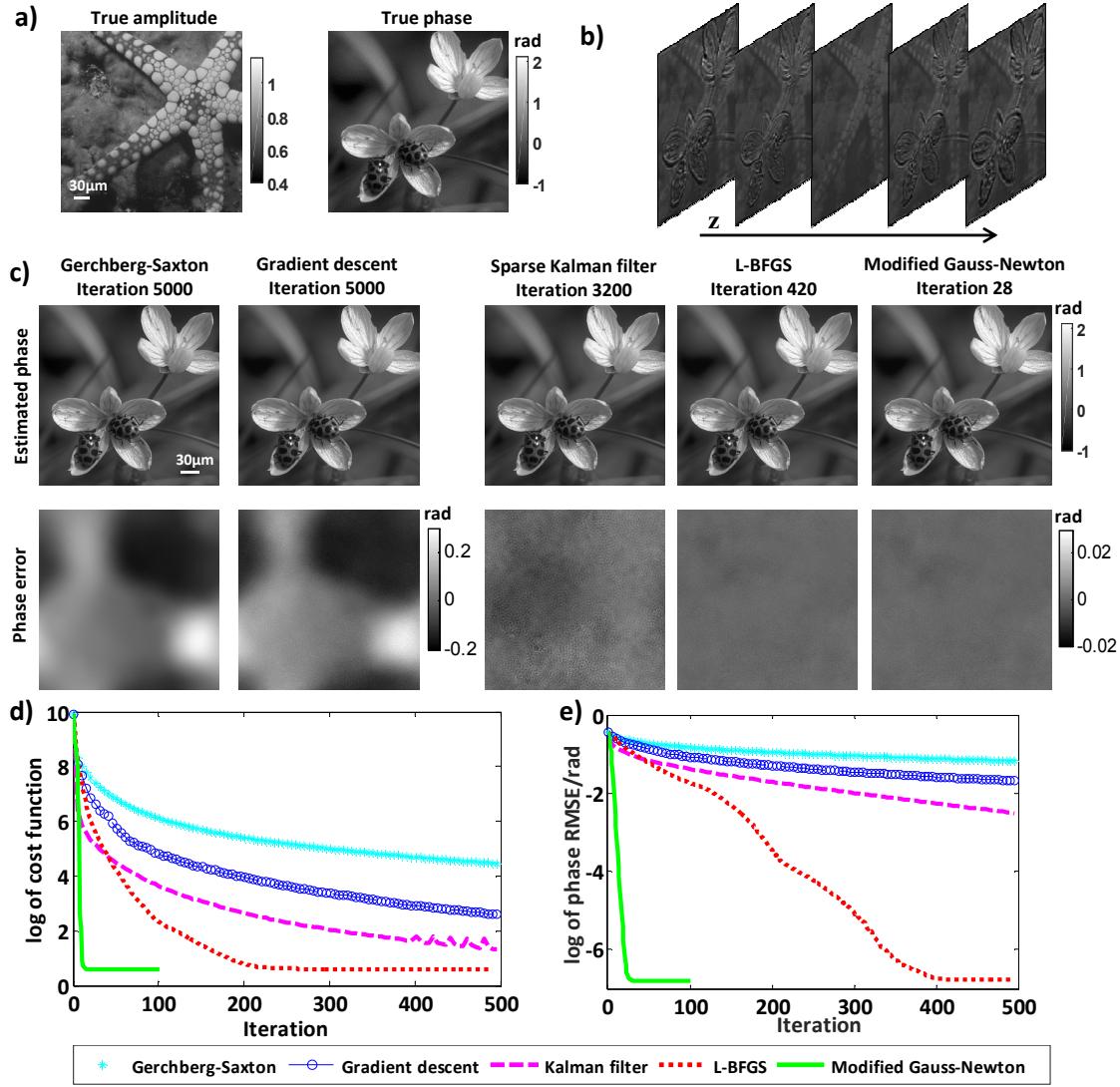


Fig. 3. Simulations comparing phase retrieval algorithms for a coherent intensity focus stack. a) Amplitude and phase of the simulated true complex field. b) The stack of intensity images captured through-focus. c) Recovered phase from each of 5 algorithms and their error maps. The *first-order* algorithms (GS and GD) show strong low frequency noise, even after 5000 iterations, while the Sparse KF, L-BFGS, and our Gauss-Newton algorithm incur less error with fewer iterations. d) The cost function (mean square error between the measured and predicted intensity) decreases with more iterations, with our Gauss-Newton algorithm converging in the fewest iterations. e) The root-mean-square error (RMSE) describing the difference between the recovered and true phase.

method in Fig. 3 refers to the global GS method, and the GD method is for a large step size. One can see from Fig. 3(c) that the *first-order* optimization methods (GS and GD) result in larger phase error, particularly for low spatial frequencies, even after 5000 iterations. The second-order optimization methods (L-BFGS and Gauss-Newton) give the best results and converge in fewer iterations, with the modified Gauss-Newton method using the fewest (~ 15).

Comparing the plots in Fig. 3(d) and (e), we see that the convergence metric (cost function) decreases fastest for our Gauss-Newton method, as does the actual error (Phase MSE) which is computed from the ground truth. Interestingly, for all but our Gauss-Newton method, the cost function decreases faster than the actual error, indicating that the cost function may be providing false confidence in the result. For fair comparison, the results shown in Fig. 3(c) allow the GS, GD and Sparse KF algorithms to run for many iterations to show

that the low frequency error persists for GS and GD results, while the recovered phase of Sparse KF method eventually almost converges to the true phase.

Of course, the reduction in number of iterations for the *second-order* methods comes at a cost of additional computation for each iteration. To fairly compare speed of convergence, we look at computational complexity and runtime. The computational complexity of all methods scales as $\mathcal{O}(N \log N)$, due to the FFT implementation. However, the L-BFGS takes two loops through the stored gradients to compute the search direction [29]. The modified Gauss-Newton method also takes extra steps in each iteration to solve the linear equation of the Hessian matrix with a conjugate gradient method. Despite the extra computational cost in using *second-order* approaches, the reduction in the number of iterations still results in shorter runtime due to the greater accuracy of each update step.

To quantify this tradeoff, we compare the accuracy and

runtime for all methods in Table I. The *first-order* methods (GS, GD) and Sparse KF take significantly more runtime and result in larger error than the *second-order* methods (L-BFGS, Gauss-Newton). Our proposed Gauss-Newton algorithm is the fastest, taking approximately half the runtime of the L-BFGS method in this case. Our method also has advantages because it does not need to make preparation for storing many gradients as in the Hessian approximation of L-BFGS.

TABLE I
COMPARISON OF ACCURACY AND RUNTIME FOR PHASE RETRIEVAL ALGORITHMS. GAUSS-NEWTON PERFORMS BEST IN THIS CASE.

	GS	GD	Sparse KF	L-BFGS	Gauss-Newton
Iteration	5000	5000	3200	420	15
log of cost function	1.39	0.61	0.60	0.59	0.59
log of phase RMSE/rad	-2.33	-4.60	-5.19	-6.78	-6.79
Runtime/s	260.0	684.4	237.63	62.2	26.2

Figure 4 shows the cost function and phase error when initializing the complex field with the square root of the measured in-focus intensity and a random phase. Each phase pattern is randomly generated by i.i.d. random Gaussian distribution of zero mean and variance 1. Both the cost functions and phase errors almost converge to the result with zero phase initialization (Table I), so our method is stable with phase initialization.

Next, we show the effect of varying the number of iterations in the conjugate gradient step within each iteration of the Gauss-Newton method. More iterations will result in more accurate solutions to the Hessian equations (enabling fewer iterations of the Gauss-Newton method), at a cost of longer compute time for each Gauss-Newton iteration. Thus, there will be an optimal choice that balances the trade-off for the fastest overall compute time. To demonstrate, we run the Gauss-Newton method and stop iterating it when the cost function converges to the value in Table I, for each case from 10 to 100 iterations in the conjugate gradient step. Figure. 5(a) shows that the number of iterations needed for the Gauss-Newton method decreases as the number of iterations in the conjugate gradient step increases, as expected. Figure. 5(b) shows that the corresponding overall runtime is optimized with approximately 50 iterations in the conjugate gradient step.

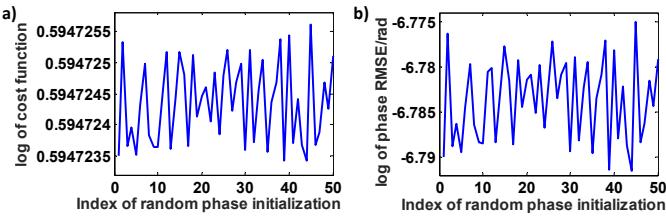


Fig. 4. Reconstruction errors in the complex field are relatively insensitive to different random phase initializations. a) The cost function and b) the phase error for different random phase initializations.

B. Experiments

To demonstrate that our findings also hold for experimental data, we compare algorithms for coherent phase retrieval using a measured focus stack from a human cheek cell. The

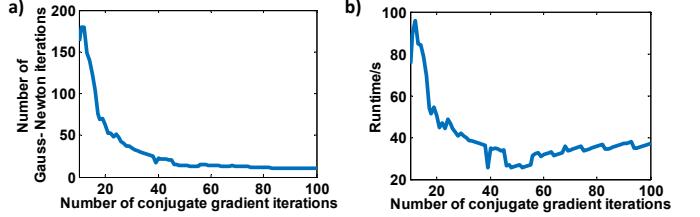


Fig. 5. The effect of changing the number of iterations in the conjugate gradient step. a) The Gauss-Newton method converges in fewer iterations as the number of iterations in the conjugate gradient step increases. b) Overall runtime as the number of iterations in the conjugate gradient step increases.

illumination is white light filtered with a center wavelength of 650 nm. The sample is placed at the focal plane and 13 intensity images are measured, spaced exponentially [19] at both sides of focus, from $6\mu\text{m}$ to $40\mu\text{m}$. Each image has 380×370 pixels of size $0.312\mu\text{m} \times 0.312\mu\text{m}$.

Figure 6 shows the recovered phase using all of the methods discussed. Since the *first-order* methods (GS and GD) converge slowly, we use more iterations. The recovered phase of sequential GS suffers from low frequency noise, since the algorithm is not noise stable. The global GS method does better, since it averages the constraints from all measurements. However, much of the low spatial frequency information is missing, giving the cells a hollowed-out look. The global GS method is equivalent to gradient descent with a constant step size for all iterations [25], so the quality of its recovered phase is in between that of the small and large step size GD methods. The *second-order* optimization methods, L-BFGS and our modified Gauss-Newton method, provide the best recovered phase and the smallest fitting error to the measured intensity images. A phase cross-section in Fig. 6(b) highlights the low spatial frequency nature of the error. All of the *first-order* algorithms recover a phase result that is lower in value than it should be (except the unstable sequential GS), indicating missing low frequency information, which is important for characterizing the cell's volume and shape.

IV. RESULTS FOR PARTIALLY COHERENT JOINT RECOVERY OF COMPLEX FIELD AND SOURCE

A. Simulations

To verify the joint complex field and source recovery algorithm that was proposed in Section II-E, we begin with simulations. A stack of through-focus intensity images representing a partially coherent imaging system (in Köhler geometry) is generated from the same complex field as in Section III. We consider three different source shapes, shown in Fig. 7(a): a circle, a ring, and an arbitrary 2D pattern. The z planes are chosen to be the same for each dataset. The first 5 near focus are exponentially spaced from $3\mu\text{m}$ to $23.4\mu\text{m}$ by a growing rate of 1.5, and the next 10 are equally spaced by $15\mu\text{m}$. Each intensity image is polluted with Gaussian white noise of variance 0.002. This data is used as input for the modified Gauss-Newton algorithm described in Section II-E.

The recovered source pattern, amplitude and phase for each source pattern are shown in Fig. 7(c), along with the error maps for each. Notice that the source error is mainly in the

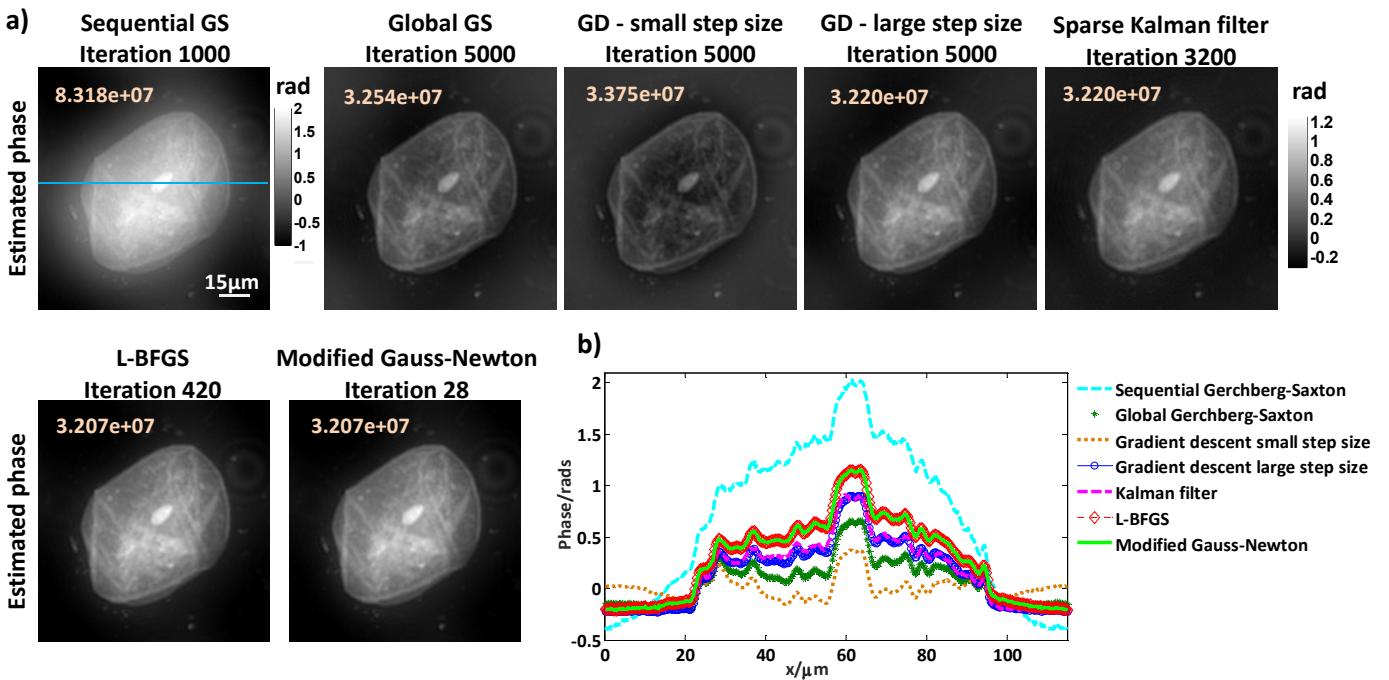


Fig. 6. Comparison of coherent phase recovery methods with experimental data. a) The recovered phase using methods: Gerchberg-Saxton (global and sequential), gradient descent (small step size and large step size), sparse Kalman filter, L-BFGS and our modified Gauss-Newton. The mean-square errors between the measured and predicted intensity stacks are shown at the top-left corner of each. The L-BFGS and our modified Gauss-Newton methods give the best recovered phase. b) The cross-section of all of the recovered phases, taken at the blue line shown in the sequential GS result.

high spatial frequencies, while the phase error is mainly in the low spatial frequencies, as expected. Figure 7(d) shows that both the source and the phase error converge towards the correct solution. Interestingly, the source recovery seems to converge within approximately 10 iterations, while the phase error requires approximately 40 iterations in each case.

B. Experiments

Next, we test our joint recovery algorithm using experimental data from MCF10A cells measured on a Nikon TE300 microscope (Fig. 8(a)). Illumination is centered at a wavelength of 550nm with 10nm bandwidth (assumed spatially incoherent and monochromatic). To generate different source patterns, we use three different condenser apertures (Fig. 8(c)). For each, a focus stack containing 29 through-focus images is collected: the first near-focus 7 images are exponentially spaced from $0.5\mu\text{m}$ to $32\mu\text{m}$, and the next 7 images are equally spaced by $10\mu\text{m}$. Each intensity image has 1000×1000 pixels of size $0.16\mu\text{m} \times 0.16\mu\text{m}$. Figure 8(b) shows measured intensity stacks for two different source sizes. The larger source produces more blur, as expected.

We consider both the *known* and *unknown* source situations. An image of the source is captured at the microscope's side port. Using this as the *known* source, we recover phase with the algorithm in Section II-C (Fig. 8(c)). Next, we recover both the phase and the source (assuming now that it is *unknown*) from the same focus stack, using the algorithm in Section II-E. The recovered source closely matches the measured one (Fig. 8(c)). However, it is slightly off-axis, which is likely caused by a misalignment of the microscope or the camera. This offset

also explains the asymmetric low frequency artifacts in the recovered phase from the measured source, and suggests that our measurements were shifted from the actual source position. The low frequency noise is alleviated in the joint phase and source recovery (Fig. 8(d)), which implicitly corrects the misalignment. This provides a demonstration of the self-calibration capabilities of the joint recovery method.

To demonstrate our algorithm with more complicated source patterns, we use a different type of source - a programmable LED array dome [35], shown in Fig. 9(a). The LED dome was originally designed for high-resolution imaging [47], [48]. Its radius is sufficiently large that we can approximate the light from each LED as a tilted plane wave at the sample, and each LED is incoherent with others, much like the Köhler geometry. We choose three illumination patterns, measure a focus stack for each (Fig. 9(b)) and use these as input to our joint recovery algorithm. A human cheek cell sample was placed on the z stage and the same defocus steps as in Fig. 8 were used. Each intensity has 864×1018 pixels of size $0.16\mu\text{m} \times 0.16\mu\text{m}$.

Figure 9(d) shows the phase recovered by our algorithm when the illumination is incorrectly assumed to be coherent. Because partially coherent effects are ignored, the result is very blurred, exacerbated by the relatively large focus range in this dataset. Large focus range is required for high-resolution source recovery (and low-frequency phase information). Figure 9(d) shows the results of our Gauss-Newton optimization method without any assumption on the source. The algorithm correctly reconstructs the patterns and resolves the LEDs, which is useful for system calibration of the LED dome.

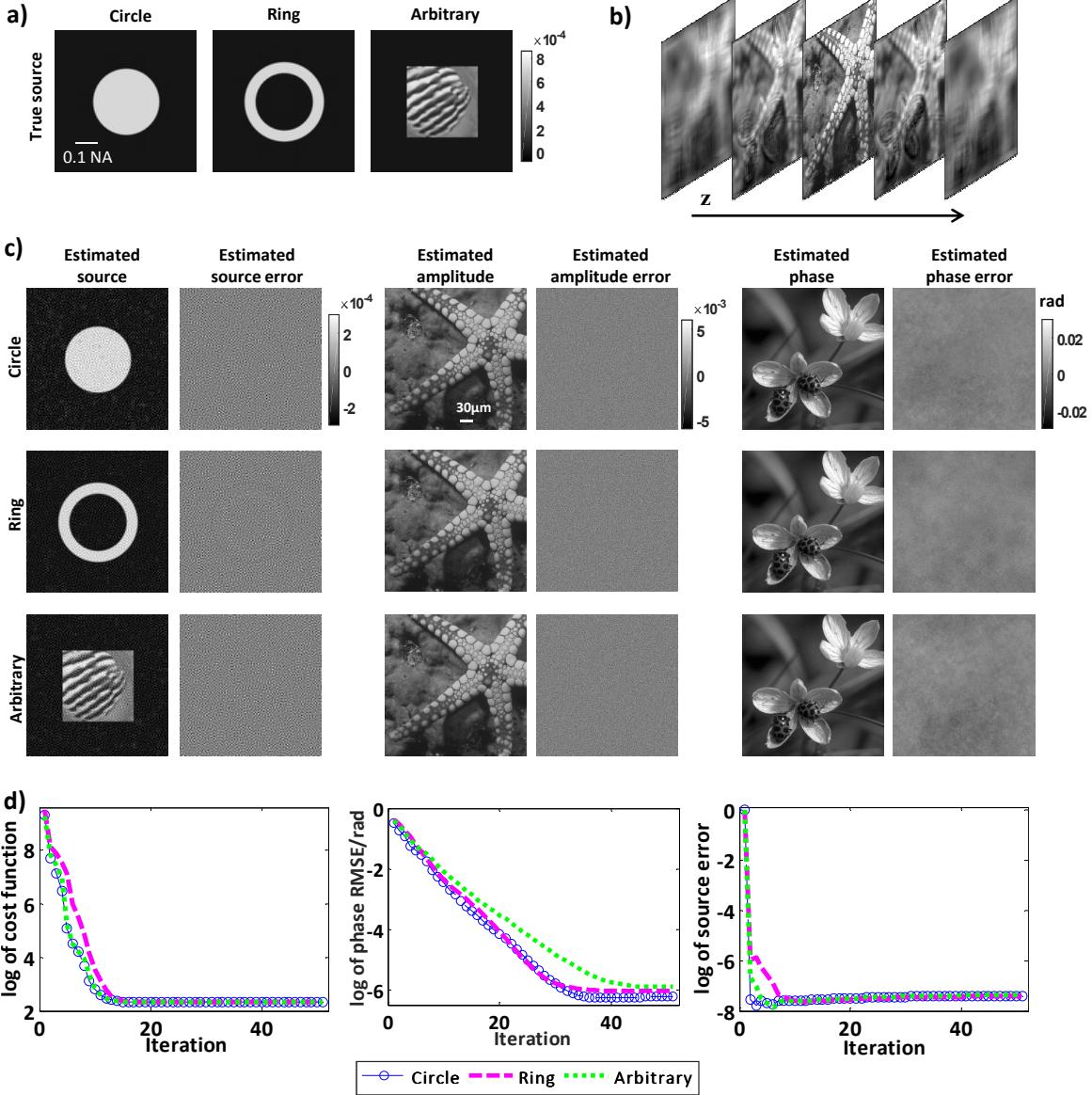


Fig. 7. Simulated joint recovery of both complex field and source using our modified Gauss-Newton algorithm. a) 3 different source patterns are tested: a circle, a ring, and a square with an arbitrary pattern inside. b) The simulated through-focus intensity dataset for the arbitrary source case. c) Recovered source, amplitude, phase, and their errors compared to the true values. The recovered amplitude and phase images share the same colorbar as Fig. 3(a). d) Convergence plots of the cost function, source error and phase error (radians) with iteration number.

V. CONCLUSION

We have proposed a robust and accurate nonlinear optimization method for recovering both the sample's complex field and an arbitrary source pattern (illumination coherence) from a single partially coherent defocused intensity stack. We derived the gradient and Hessian of the cost function with complex variables. In the simplified case of coherent imaging, we showed that our second-order Gauss-Newton optimization method performs better than existing phase retrieval algorithms including Gerchberg-Saxton, gradient descent, Kalman filter and the L-BFGS method. In general, both second-order methods (Gauss-Newton, L-BFGS) converged significantly faster than first-order methods. While our modified Gauss-Newton method performed best in this case (defocus stacks),

the L-BFGS method is more easily generalized for different phase recovery situations, so may find use.

We then formulated the partially coherent optimization problem based on a simple convolution model for defocused intensity. After solving the phase and source recovery inverse problems separately, the joint optimization problem was solved iteratively using an alternating minimization approach.

We validated our joint recovery of complex field and source in simulations and with experimental data measured on a commercial microscope. Our modified Gauss-Newton algorithm can be used for both *coherent* and *partially coherent* phase retrieval applications with either a *known* or *unknown* source pattern. Potential applications beyond that of biological microscopy include surface metrology, lithography inspection systems and X-ray phase imaging.

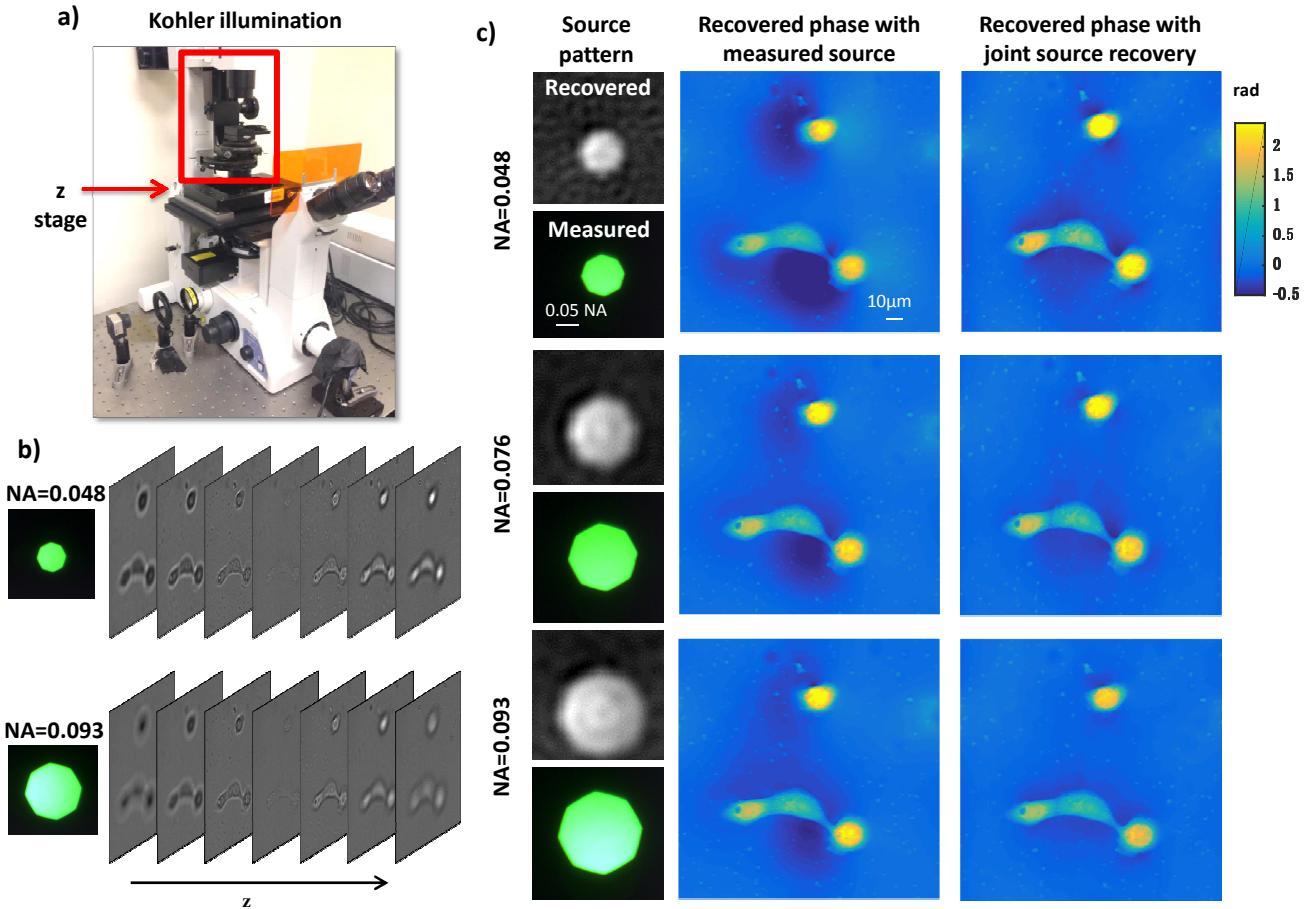


Fig. 8. Experimental results using: a) a Nikon TE300 microscope equipped with a piezo z stage for capturing partially coherent intensity stacks under different condenser apertures (source patterns). b) Two focus stacks with sources of different sizes. c) Left: measured (green) and reconstructed (grayscale) source patterns for apertures of three different sizes (NA: 0.048, 0.076 and 0.093). Middle: recovered phase using the measured source. Right: recovered phase with joint estimation of phase and source, exhibiting reduced low frequency artifacts, due to self-correction of source misalignment. Amplitude reconstructions are not shown because the samples are transparent.

Source code and sample datasets can be found at <http://www.laurawaller.com/opensource/>.

APPENDIX A DERIVING THE FIRST-ORDER AND SECOND-ORDER DERIVATIVES OF THE COST FUNCTION WITH \mathbb{CR} CALCULUS

The cost function in Eq. (9) with the complex variable \mathbf{b} is non-analytic, so it is not complex-differentiable. Instead, we take the derivatives with respect to the real and imaginary part of the complex variable. However, it is not straightforward to do so [49]. The \mathbb{CR} calculus theory [39], [40] provides a convenient method to take derivatives for a cost function with complex variable. The \mathbb{CR} calculus method rewrites the cost function with an augmented complex variable $[\mathbf{b} \ \bar{\mathbf{b}}]^T$, and then treats both \mathbf{b} and $\bar{\mathbf{b}}$ as independent variables when taking the derivatives. The derivatives with respect to the real and imaginary parts of the complex variable are linear transforms for the derivatives with respect to the augmented variable. More importantly, \mathbb{CR} calculus shows that nonlinear optimization methods can be used to update the augmented variable, which avoids flipping backward and forward between the real and complex domain [40].

We write a general form of cost function for phase recovery problem from the square error cost function in Eq. (9):

$$\min_{\mathbf{b}} F = \frac{1}{2} \sum_{n=1}^{N_{im}} \left(\mathbf{I}_n - \mathbf{C}_n |\mathbf{A}_n \mathbf{b}|^2 \right)^H \left(\mathbf{I}_n - \mathbf{C}_n |\mathbf{A}_n \mathbf{b}|^2 \right), \quad (17)$$

where \mathbf{b} is the unknown and the matrix \mathbf{A}_n denotes the transfer function of the system. In the cost function of Eq. (9), the matrices $\mathbf{A}_n = \mathbf{K}^H \mathbf{H}_n$.

First, we express the cost function as a summation of error from each intensity measurement. We define the difference between measurement and prediction for image n as a vector $\mathbf{f}_n(\mathbf{b}, \bar{\mathbf{b}})$, and the error for image n as $F_n(\mathbf{b}, \bar{\mathbf{b}})$:

$$\mathbf{f}_n(\mathbf{b}, \bar{\mathbf{b}}) = \mathbf{I}_n - \mathbf{C}_n |\mathbf{A}_n \mathbf{b}|^2, \quad (18)$$

$$\begin{aligned} F_n(\mathbf{b}, \bar{\mathbf{b}}) &= \frac{1}{2} \mathbf{f}_n(\mathbf{b}, \bar{\mathbf{b}})^H \mathbf{f}_n(\mathbf{b}, \bar{\mathbf{b}}) \\ &= \frac{1}{2} \sum_{m=1}^M \mathbf{f}_{nm}(\mathbf{b}, \bar{\mathbf{b}}) \bar{\mathbf{f}}_{nm}(\mathbf{b}, \bar{\mathbf{b}}), \end{aligned} \quad (19)$$

where M is the length of the vector $\mathbf{f}_n(\mathbf{b}, \bar{\mathbf{b}})$, and $\mathbf{f}_{nm}(\mathbf{b}, \bar{\mathbf{b}})$ is the m^{th} entry in the vector. The cost function can then be

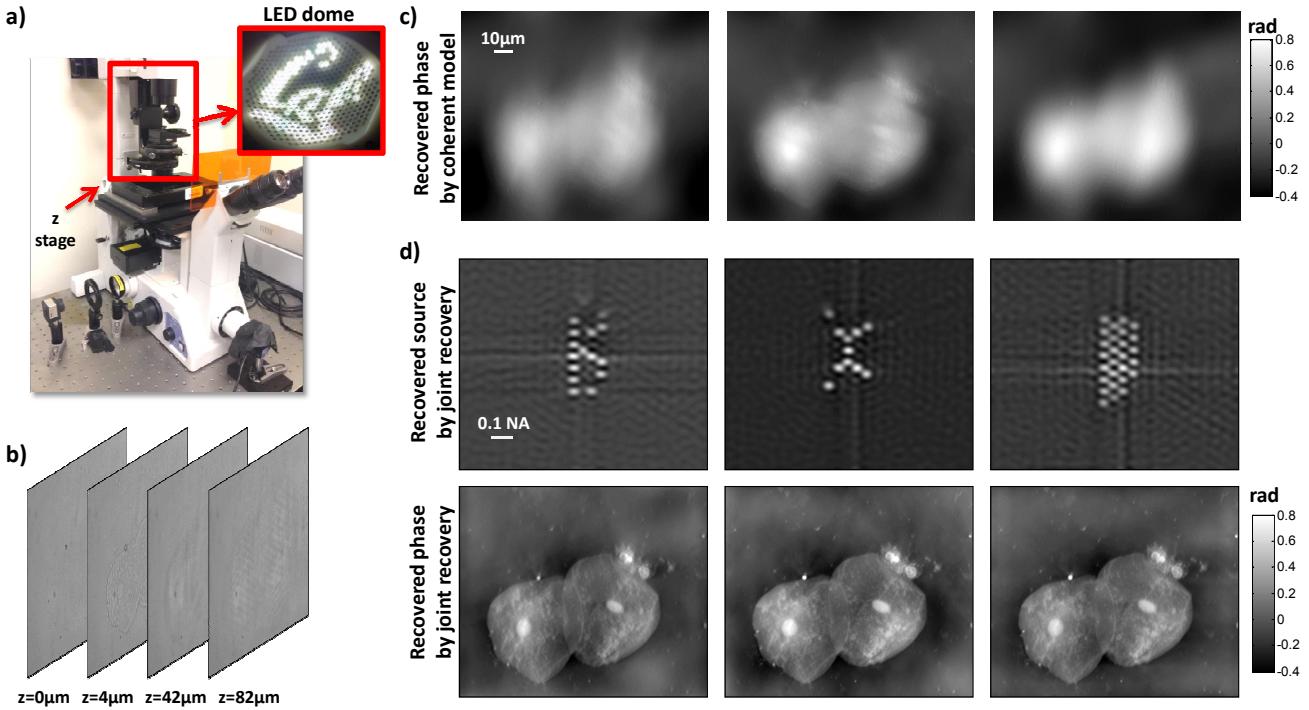


Fig. 9. Experimental results using coded illumination angles in our LED dome. a) The illumination unit of the microscope is replaced with the LED dome. b) An example of the captured intensity images through-focus for one illumination pattern. c) The reconstructed phase using a coherent phase retrieval algorithm. d) Top row: three reconstructed source patterns from our joint recovery algorithm. Bottom row: recovered phase images.

expressed as a summation the error $F_n(\mathbf{b}, \bar{\mathbf{b}})$ for all images:

$$\begin{aligned} \min_{\mathbf{b}} F &= \frac{1}{2} \sum_{n=1}^{N_{im}} \left(\mathbf{I}_n - \mathbf{C}_n |\mathbf{A}_n \mathbf{b}|^2 \right)^H \left(\mathbf{I}_n - \mathbf{C}_n |\mathbf{A}_n \mathbf{b}|^2 \right) \\ &= \sum_{n=1}^{N_{im}} F_n(\mathbf{b}, \bar{\mathbf{b}}). \end{aligned} \quad (20)$$

Next, we derive the first-order and second-order derivatives (gradient and Hessian) for each component $F_n(\mathbf{b}, \bar{\mathbf{b}})$. From [29], [40], the gradient and Hessian for $F_n(\mathbf{b}, \bar{\mathbf{b}})$ can be obtained from the chain rule:

$$\nabla F_n(\mathbf{b}, \bar{\mathbf{b}}) = \mathbf{J}_n(\mathbf{b}, \bar{\mathbf{b}})^H \mathbf{f}_n(\mathbf{b}, \bar{\mathbf{b}}), \quad (21)$$

$$\begin{aligned} \nabla^2 F_n(\mathbf{b}, \bar{\mathbf{b}}) &= \mathbf{J}_n(\mathbf{b}, \bar{\mathbf{b}})^H \mathbf{J}_n(\mathbf{b}, \bar{\mathbf{b}}) \\ &\quad + \sum_{m=1}^M \mathbf{f}_{nm}(\mathbf{b}, \bar{\mathbf{b}}) \mathbf{f}_{nm}^H(\mathbf{b}, \bar{\mathbf{b}}), \end{aligned} \quad (22)$$

where $\mathbf{J}_n(\mathbf{b}, \bar{\mathbf{b}})^H$ is Jacobian of the function vector $\mathbf{f}_n(\mathbf{b}, \bar{\mathbf{b}})$. From \mathbb{CR} calculus, we have

$$\begin{aligned} \mathbf{J}_n(\mathbf{b}, \bar{\mathbf{b}}) &= \left[\begin{array}{cc} \frac{\partial \mathbf{f}_n}{\partial \mathbf{b}} & \frac{\partial \mathbf{f}_n}{\partial \bar{\mathbf{b}}} \end{array} \right] \\ &= \left[\begin{array}{cc} -\mathbf{C}_n \text{diag}(\bar{\mathbf{A}}_n \bar{\mathbf{b}}) \mathbf{A}_n & -\mathbf{C}_n \text{diag}(\mathbf{A}_n \mathbf{b}) \bar{\mathbf{A}}_n \end{array} \right]. \end{aligned} \quad (23)$$

We follow the Gauss-Newton method (Chapter 10 of [29]) to approximate the Hessian matrix with only the Jacobian matrix by omitting the second term $\sum_{m=1}^M \mathbf{f}_{nm}(\mathbf{b}, \bar{\mathbf{b}}) \mathbf{f}_{nm}^H(\mathbf{b}, \bar{\mathbf{b}})$:

$$\nabla^2 F_n(\mathbf{b}, \bar{\mathbf{b}}) \approx \mathbf{J}_n(\mathbf{b}, \bar{\mathbf{b}})^H \mathbf{J}_n(\mathbf{b}, \bar{\mathbf{b}}). \quad (24)$$

Then, the gradient and Hessian (and its approximation) of the cost function could be expressed as summation of the

gradient and Hessian of each component $F_n(\mathbf{b}, \bar{\mathbf{b}})$. From Eqs. (21)(23)(24), we have

$$\begin{aligned} \nabla F(\mathbf{b}, \bar{\mathbf{b}}) &= \sum_{n=1}^{N_{im}} \nabla F_n(\mathbf{b}, \bar{\mathbf{b}}) \\ &= \sum_{n=1}^{N_{im}} \left[\begin{array}{c} -\mathbf{A}_n^H \text{diag}(\mathbf{A}_n \mathbf{b}) \mathbf{C}_n^H \\ -\mathbf{A}_n^T \text{diag}(\bar{\mathbf{A}}_n \bar{\mathbf{b}}) \mathbf{C}_n^H \end{array} \right] \left[\mathbf{I}_n - \mathbf{C}_n |\mathbf{A}_n \mathbf{b}|^2 \right], \end{aligned} \quad (25)$$

$$\begin{aligned} \nabla^2 F(\mathbf{b}, \bar{\mathbf{b}}) &= \sum_{n=1}^{N_{im}} \nabla^2 F_n(\mathbf{b}, \bar{\mathbf{b}}) = \sum_{n=1}^{N_{im}} \mathbf{J}_n(\mathbf{b}, \bar{\mathbf{b}})^H \mathbf{J}_n(\mathbf{b}, \bar{\mathbf{b}}) \\ &= \sum_{n=1}^{N_{im}} \left[\begin{array}{cc} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{array} \right], \end{aligned} \quad (26)$$

where

$$\begin{aligned} \mathbf{B}_{11} &= \mathbf{A}_n^H \text{diag}(\mathbf{A}_n \mathbf{b}) \mathbf{C}_n^H \mathbf{C}_n \text{diag}(\bar{\mathbf{A}}_n \bar{\mathbf{b}}) \mathbf{A}_n \\ \mathbf{B}_{12} &= \mathbf{A}_n^H \text{diag}(\mathbf{A}_n \mathbf{b}) \mathbf{C}_n^H \mathbf{C}_n \text{diag}(\mathbf{A}_n \mathbf{b}) \bar{\mathbf{A}}_n \\ \mathbf{B}_{21} &= \mathbf{A}_n^T \text{diag}(\bar{\mathbf{A}}_n \bar{\mathbf{b}}) \mathbf{C}_n^H \mathbf{C}_n \text{diag}(\bar{\mathbf{A}}_n \bar{\mathbf{b}}) \mathbf{A}_n \\ \mathbf{B}_{22} &= \mathbf{A}_n^T \text{diag}(\bar{\mathbf{A}}_n \bar{\mathbf{b}}) \mathbf{C}_n^H \mathbf{C}_n \text{diag}(\mathbf{A}_n \mathbf{b}) \bar{\mathbf{A}}_n. \end{aligned} \quad (27)$$

Last, we derive the gradient and Hessian of Eq. (9) from Equations (25)(26) by replacing \mathbf{A}_n with $\mathbf{K}^H \mathbf{H}_n$. Since \mathbf{K} is the DFT matrix and \mathbf{H}_n is a diagonal matrix, we have

$$\begin{aligned} \mathbf{K}^T &= \mathbf{K}, \quad (\mathbf{K}^H)^T = \mathbf{K}^H, \quad \bar{\mathbf{K}} = \mathbf{K}^H; \\ \mathbf{H}_n^T &= \mathbf{H}_n, \quad \mathbf{H}_n^H = \bar{\mathbf{H}}_n; \\ \mathbf{A}_n^H &= \bar{\mathbf{H}}_n \mathbf{K}, \quad \mathbf{A}_n^T = \mathbf{H}_n \mathbf{K}^H, \quad \bar{\mathbf{A}}_n = \mathbf{K} \bar{\mathbf{H}}_n. \end{aligned} \quad (28)$$

From Eqs. (25)(26), we obtain

$$\nabla F(\mathbf{b}, \bar{\mathbf{b}}) = \sum_{n=1}^{N_{im}} \begin{bmatrix} -\bar{\mathbf{H}}_n \mathbf{K} \text{diag}(\mathbf{K}^H \mathbf{H}_n \mathbf{b}) \mathbf{C}_n^H \\ -\mathbf{H}_n \mathbf{K}^H \text{diag}(\mathbf{K} \bar{\mathbf{H}}_n \bar{\mathbf{b}}) \mathbf{C}_n^H \end{bmatrix} \left[\mathbf{I}_n - \mathbf{C}_n |\mathbf{K}^H \mathbf{H}_n \mathbf{b}|^2 \right], \quad (29)$$

$$\nabla^2 F(\mathbf{b}, \bar{\mathbf{b}}) \approx \sum_{n=1}^{N_{im}} \begin{bmatrix} \mathbf{M}_{11} & \mathbf{M}_{12} \\ \mathbf{M}_{21} & \mathbf{M}_{22} \end{bmatrix}, \quad (30)$$

where

$$\mathbf{M}_{11} = \bar{\mathbf{H}}_n \mathbf{K} \text{diag}(\mathbf{K}^H \mathbf{H}_n \mathbf{b}) \mathbf{C}_n^H \mathbf{C}_n \text{diag}(\mathbf{K} \bar{\mathbf{H}}_n \bar{\mathbf{b}}) \mathbf{K}^H \mathbf{H}_n$$

$$\mathbf{M}_{12} = \bar{\mathbf{H}}_n \mathbf{K} \text{diag}(\mathbf{K}^H \mathbf{H}_n \mathbf{b}) \mathbf{C}_n^H \mathbf{C}_n \text{diag}(\mathbf{K}^H \mathbf{H}_n \mathbf{b}) \mathbf{K} \bar{\mathbf{H}}_n$$

$$\mathbf{M}_{21} = \mathbf{H}_n \mathbf{K}^H \text{diag}(\mathbf{K} \bar{\mathbf{H}}_n \bar{\mathbf{b}}) \mathbf{C}_n^H \mathbf{C}_n \text{diag}(\mathbf{K} \bar{\mathbf{H}}_n \bar{\mathbf{b}}) \mathbf{K}^H \mathbf{H}_n$$

$$\mathbf{M}_{22} = \mathbf{H}_n \mathbf{K}^H \text{diag}(\mathbf{K} \bar{\mathbf{H}}_n \bar{\mathbf{b}}) \mathbf{C}_n^H \mathbf{C}_n \text{diag}(\mathbf{K}^H \mathbf{H}_n \mathbf{b}) \mathbf{K} \bar{\mathbf{H}}_n.$$

ACKNOWLEDGMENTS

The authors thank Gongguo Tang, Jingzhao Zhang, and Zachary Phillips for discussions and help with coding. This work was funded by the Gordon and Betty Moore Foundation's Data-Driven Discovery Initiative Grant GBMF4562 to Laura Waller (UC Berkeley), and by the Office of Naval Research Grant N00014-14-1-0083.

REFERENCES

- [1] J. R. Fienup, "Phase retrieval algorithms: a comparison," *Appl. Opt.*, vol. 21, no. 15, pp. 2758–2769, 1982.
- [2] M. Teague, "Deterministic phase retrieval: a Green's function solution," *J. Opt. Soc. Am.*, vol. 73, no. 11, pp. 1434–1441, Nov 1983.
- [3] J. R. Fienup, "Reconstruction of an object from the modulus of its Fourier transform," *Opt. Lett.*, vol. 3, no. 1, pp. 27–29, Jul 1978.
- [4] G. R. Brady and J. R. Fienup, "Nonlinear optimization algorithm for retrieving the full complex pupil function," *Opt. Express*, vol. 14, no. 2, pp. 474–486, 2006.
- [5] K. Nugent, T. Gureyev, D. Cookson, D. Paganin, and Z. Barnea, "Quantitative phase imaging using hard X rays," *Phys. Rev. Lett.*, vol. 77, no. 14, pp. 2961–2964, 1996.
- [6] C. J. R. Sheppard, "Defocused transfer function for a partially coherent microscope and application to phase retrieval," *J. Opt. Soc. Am. A*, vol. 21, no. 5, pp. 828–831, May 2004.
- [7] D. Paganin and K. A. Nugent, "Noninterferometric phase imaging with partially coherent light," *Phys. Rev. Lett.*, vol. 80, pp. 2586–2589, Mar 1998.
- [8] R. A. Claus, Y.-G. Wang, A. Wojdyla, M. P. Benk, K. A. Goldberg, A. R. Neureuther, P. P. Naulleau, and L. Waller, "Phase measurements of EUV mask defects," in *SPIE Advanced Lithography*. International Society for Optics and Photonics, 2015, pp. 942.217–942.217.
- [9] J. C. Petruccielli, L. Tian, and G. Barbastathis, "The transport of intensity equation for optical path length recovery using partially coherent illumination," *Opt. Express*, vol. 21, no. 12, pp. 14 430–14 441, Jun 2013.
- [10] Z. Jingshan, L. Tian, R. A. Claus, J. Dauwels, and L. Waller, "Partially coherent phase recovery by Kalman filtering," in *Frontiers in Optics*. Optical Society of America, 2013, pp. FW6A–9.
- [11] J. A. Rodrigo and T. Alieva, "Rapid quantitative phase imaging for partially coherent light microscopy," *Opt. Express*, vol. 22, no. 11, pp. 13 472–13 483, Jun 2014.
- [12] Z. Jingshan, L. Tian, J. Dauwels, and L. Waller, "Partially coherent phase imaging with simultaneous source recovery," *Biomed. Opt. Express*, vol. 6, no. 1, pp. 257–265, Jan 2015.
- [13] D. Paganin, A. Barty, P. McMahon, and K. Nugent, "Quantitative phase-amplitude microscopy. III. the effects of noise," *J. Microsc.*, vol. 214, no. 1, pp. 51–61, 2004.
- [14] R. G. Paxman, T. J. Schulz, and J. R. Fienup, "Joint estimation of object and aberrations by using phase diversity," *JOSA A*, vol. 9, no. 7, pp. 1072–1085, 1992.
- [15] M. Soto and E. Acosta, "Improved phase imaging from intensity measurements in multiple planes," *Appl. Opt.*, vol. 46, no. 33, pp. 7978–7981, Nov 2007.
- [16] L. Allen and M. Oxley, "Phase retrieval from series of images obtained by defocus variation," *Opt. Commun.*, vol. 199, no. 1, pp. 65–75, 2001.
- [17] N. Streibl, "Phase imaging by the transport equation of intensity," *Opt. Commun.*, vol. 49, no. 1, pp. 6–10, 1984.
- [18] L. Waller, L. Tian, and G. Barbastathis, "Transport of intensity phase-amplitude imaging with higher order intensity derivatives," *Opt. Express*, vol. 18, no. 12, pp. 12 552–12 561, Jun 2010.
- [19] Z. Jingshan, R. A. Claus, J. Dauwels, L. Tian, and L. Waller, "Transport of intensity phase imaging by intensity spectrum fitting of exponentially spaced defocus planes," *Opt. Express*, vol. 22, no. 9, pp. 10 661–10 674, May 2014.
- [20] L. Waller, M. Tsang, S. Ponda, S. Yang, and G. Barbastathis, "Phase and amplitude imaging from noisy images by Kalman filtering," *Opt. Express*, vol. 19, no. 3, pp. 2805–2814, Jan 2011.
- [21] R. W. Gerchberg and W. O. Saxton, "A practical algorithm for the determination of phase from image and diffraction plane pictures," *Optik*, vol. 35, p. 237, 1972.
- [22] D. C. Youla and H. Webb, "Image restoration by the method of convex projections: Part 1, theory," *IEEE Transactions on Medical Imaging*, vol. 1, no. 2, pp. 81–94, 1982.
- [23] R. A. Gonsalves and R. Chidlaw, "Wavefront sensing by phase retrieval," in *23rd Annual Technical Symposium*. International Society for Optics and Photonics, 1979, pp. 32–39.
- [24] D. S. Weller, A. Pnueli, G. Divon, O. Radzyner, Y. C. Eldar, and J. A. Fessler, "Undersampled phase retrieval with outliers," *IEEE Transactions on Computational Imaging*, vol. 1, no. 4, pp. 247–258, 2015.
- [25] C. Yang, J. Qian, A. Schirozek, F. Maia, and S. Marchesini, "Iterative algorithms for ptychographic phase retrieval," Ernest Orlando Lawrence Berkeley National Laboratory, Berkeley, CA (US), Tech. Rep., 2011.
- [26] J. Qian, C. Yang, A. Schirozek, F. Maia, and S. Marchesini, "Efficient algorithms for ptychographic phase retrieval," *Inverse Problems and Applications*, vol. 615, p. 261, 2014.
- [27] J. Rodenburg, "Ptychography and related diffractive imaging methods," *Advances in Imaging and Electron Physics*, vol. 150, pp. 87–184, 2008.
- [28] L.-H. Yeh, J. Dong, J. Zhong, L. Tian, M. Chen, G. Tang, M. Soltanolkotabi, and L. Waller, "Experimental robustness of Fourier ptychography phase retrieval algorithms," *Opt. Express*, vol. 23, no. 26, pp. 33 214–33 240, 2015.
- [29] J. Nocedal and S. Wright, *Numerical Optimization*, ser. Springer series in operations research and financial engineering. Springer, 1999.
- [30] D. C. Liu and J. Nocedal, "On the limited memory BFGS method for large scale optimization," *Mathematical programming*, vol. 45, no. 1-3, pp. 503–528, 1989.
- [31] C. Mack, *Fundamental principles of optical lithography: the science of microfabrication*. John Wiley & Sons, 2008.
- [32] G. J. Williams, H. M. Quiney, A. G. Peele, and K. A. Nugent, "Coherent diffractive imaging and partial coherence," *Phys. Rev. B*, vol. 75, p. 104102, Mar 2007.
- [33] J. R. Valenzuela, J. A. Fessler, and R. G. Paxman, "Joint estimation of stokes images and aberrations from phase-diverse polarimetric measurements," *JOSA A*, vol. 27, no. 5, pp. 1185–1193, 2010.
- [34] L. Mandel and E. Wolf, *Optical coherence and quantum optics*. Cambridge university press, 1995.
- [35] Z. F. Phillips, M. V. D'Ambrosio, L. Tian, J. J. Rulison, H. S. Patel, N. Sadras, A. V. Gande, N. A. Switz, D. A. Fletcher, and L. Waller, "Multi-contrast imaging and digital refocusing on a mobile microscope with a domed LED array," *PLoS one*, vol. 10, no. 5, 2015.
- [36] K. Nugent, "A generalization of Schell's theorem," *Opt. Commun.*, vol. 79, no. 5, pp. 267 – 269, 1990.
- [37] A. Martin, F.-R. Chen, W.-K. Hsieh, J.-J. Kai, S. Findlay, and L. Allen, "Spatial incoherence in phase retrieval based on focus variation," *Ultramicroscopy*, vol. 106, no. 10, pp. 914 – 924, 2006.
- [38] M. H. Jenkins, J. M. Long, and T. K. Gaylord, "Multifilter phase imaging with partially coherent light," *Appl. Opt.*, vol. 53, pp. D29–D39, 2014.
- [39] D. Brandwood, "A complex gradient operator and its application in adaptive array theory," in *IEE Proceedings F (Communications, Radar and Signal Processing)*, vol. 130, no. 1. IET, 1983, pp. 11–16.
- [40] K. Kreutz-Delgado, "The complex gradient operator and the CR-calculus," *arXiv preprint arXiv:0906.4835*, 2009.
- [41] L. Armijo, "Minimization of functions having lipschitz continuous first partial derivatives," *Pacific Journal of mathematics*, vol. 16, no. 1, pp. 1–3, 1966.
- [42] J. Fienup and C. Wackerman, "Phase-retrieval stagnation problems and solutions," *JOSA A*, vol. 3, no. 11, pp. 1897–1907, 1986.

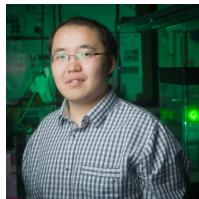
- [43] S. T. Thurman, R. T. DeRosa, and J. R. Fienup, "Amplitude metrics for field retrieval with hard-edged and uniformly illuminated apertures," *JOSA A*, vol. 26, no. 3, pp. 700–709, 2009.
- [44] D. Lee, M. Roggemann, and B. Welsh, "Cramér-Rao analysis of phase-diverse wave-front sensing," *J. Opt. Soc. Am. A*, vol. 16, no. 5, pp. 1005–1015, May 1999.
- [45] G. Pedrini, W. Osten, and Y. Zhang, "Wave-front reconstruction from a sequence of interferograms recorded at different planes," *Opt. Lett.*, vol. 30, no. 8, pp. 833–835, 2005.
- [46] Z. Jingshan, J. Dauwels, M. A. Vázquez, and L. Waller, "Sparse ACEKF for phase reconstruction," *Opt. Express*, vol. 21, no. 15, pp. 18 125–18 137, Jul 2013.
- [47] G. Zheng, R. Horstmeyer, and C. Yang, "Wide-field, high-resolution Fourier ptychographic microscopy," *Nature photonics*, vol. 7, no. 9, pp. 739–745, 2013.
- [48] L. Tian, X. Li, K. Ramchandran, and L. Waller, "Multiplexed coded illumination for Fourier ptychography with an LED array microscope," *Biomedical optics express*, vol. 5, no. 7, pp. 2376–2389, 2014.
- [49] A. S. Jurling and J. R. Fienup, "Applications of algorithmic differentiation to phase retrieval algorithms," *J. Opt. Soc. Am. A*, vol. 31, no. 7, pp. 1348–1359, 2014.



Laura Waller is an Assistant Professor in Electrical Engineering and Computer Sciences at UC Berkeley, where she leads the Computational Imaging Lab. Laura was a Postdoctoral Research Associate in Electrical Engineering and Lecturer of Physics at Princeton University from 2010-2012 and received B.S., M.Eng., and Ph.D. degrees in EECS from the Massachusetts Institute of Technology (MIT) in 2004, 2005, and 2010, respectively. She is a Packard Fellow, Baker Fellow, Moore Foundation Data-Driven Investigator and NSF CAREER Awardee.



Jingshan Zhong is a Postdoctoral Associate in the department of Electrical Engineering and Computer Sciences at UC Berkeley. He received a B. S. in Electronic Science and Technology from the University of Science and Technology of China in 2010. He completed his Ph.D. in Electrical and Electronic Engineering from Nanyang Technological University in 2015. His research interests include phase imaging, high resolution imaging, imaging through scattering by phase space measurements, signal processing, numerical optimization and machine learning.



Lei Tian is a Postdoctoral Associate in the department of Electrical Engineering and Computer Sciences (EECS) at UC Berkeley. Prior to that, he completed an M.S. (2010) and Ph.D. (2013) in the department of Mechanical Engineering at the Massachusetts Institute of Technology (MIT). His research is in computational imaging, which jointly designs signal processing and optics. He has developed compressive imaging techniques for visible and X-ray phase imaging, as well as computational microscopes for 3D, phase and gigapixel imaging.



Paroma Varma is a graduate student at Stanford University. She completed a B.S. in Electrical Engineering and Computer Sciences at UC Berkeley in 2015. Her research focuses on signal processing for computational imaging.