



Sparse reconstruction of wavefronts using an over-complete phase dictionary

S. HOWARD,^{1,2} N. WEISSE,²  J. SCHRÖDER,²  C. BARBERO,³ 
B. ALONSO,^{3,4}  Í. SOLA,^{3,4}  P. NORREYS,¹ AND A. DÖPP^{1,2,*} 

¹Department of Physics, Clarendon Laboratory, University of Oxford, Oxford OX1 3PU, United Kingdom

²Ludwig-Maximilians-Universität München, Am Coulombwall 1, 85748 Garching, Germany

³Grupo de Investigación en Aplicaciones del Láser y Fotónica (ALF), Universidad de Salamanca, 37008 Salamanca, Spain

⁴Unidad de Excelencia en Luz y Materia Estructuradas (LUMES), Universidad de Salamanca, Spain

*a.doepp@lmu.de

Abstract: Wavefront reconstruction is a critical component in various optical systems, including adaptive optics, interferometry, and phase contrast imaging. Traditional reconstruction methods often employ either the Cartesian (pixel) basis or the Zernike polynomial basis. While the Cartesian basis is adept at capturing high-frequency features, it is susceptible to overfitting and inefficiencies due to the high number of degrees of freedom. The Zernike basis efficiently represents common optical aberrations but struggles with complex or non-standard wavefronts such as optical vortices, Bessel beams, or wavefronts with sharp discontinuities. This paper introduces a novel approach to wavefront reconstruction using an over-complete phase dictionary combined with sparse representation techniques. By constructing a dictionary that includes a diverse set of basis functions—ranging from Zernike polynomials to specialized functions representing optical vortices and other complex modes—we enable a more flexible and efficient representation of complex wavefronts. Furthermore, a trainable rigid transform is implemented to account for misalignment. Utilizing principles from compressed sensing and sparse coding, we enforce sparsity in the coefficient space to avoid overfitting and enhance robustness to noise.

Published by Optica Publishing Group under the terms of the [Creative Commons Attribution 4.0 License](https://creativecommons.org/licenses/by/4.0/). Further distribution of this work must maintain attribution to the author(s) and the published article's title, journal citation, and DOI.

1. Introduction

Wavefront sensing serves as a pivotal element in numerous domains ranging from computer vision and phase contrast imaging to optical imaging and astrophysics. The wavefront plays a crucial role in shaping the intensity distribution of light as it propagates. This fundamental relationship between phase and intensity forms the basis for various wavefront measurement techniques. As light travels through space or interacts with optical elements, the phase of the wavefront undergoes changes, which manifest as variations in the observed intensity patterns. By carefully analyzing these intensity variations, it becomes possible to infer the underlying phase information.

Wavefront measurement techniques leverage this principle in different ways. Some techniques, such as the Gerchberg-Saxton (GS) algorithm [1] or transport of intensity equation (TIE) methods [2], rely on knowledge of the intensity at a reference plane and compare it with intensities measured at other planes to reconstruct the phase. Others, like interferometric methods [3] or wavefront slope sensors [4], manipulate the wavefront to create intensity patterns that encode the phase information. The distinction between these techniques lies in their assumptions about the reference intensity and how they utilize the available intensity measurements. Some techniques require precise knowledge of the reference intensity, while others can work with relaxed assumptions or infer it from measurements in a different plane.

For any wavefront measurement technique, the wavefront is extracted from the raw signal using a reconstruction process. This process varies depending on the method - for example, it might involve propagation calculations in the Gerchberg-Saxton algorithm or the stitching of gradients in other methods. In this reconstruction step, an important consideration is the basis in which the wavefront is expressed. The raw signal on the sensor is typically measured on a pixel grid arranged in Cartesian coordinates, often referred to as the pixel basis. Each pixel represents a mode in this basis, with its intensity value serving as the coefficient. Conceptually, it might seem straightforward to perform the reconstruction of the wavefront in this same basis, determining the phase value for each individual pixel. This approach, known as zonal reconstruction, is commonly used in various wavefront sensing methods, including Lateral Shearing Interferometry (LSI) and Shack-Hartmann (SH) sensors [5–9].

The pixel basis maintains a direct correspondence between spatial positions in the raw signal and the reconstructed wavefront, making it well-suited for capturing high-frequency features such as edges. This has made it prominent in applications like computer vision and phase contrast imaging. However, the pixel basis has a significant limitation: its proneness to overfitting. Given the high degrees of freedom in the Cartesian basis - typically one value per pixel - there is an increased risk of modeling noise instead of the actual signal, especially in the presence of limited data or high noise levels. These limitations of the pixel basis motivate the exploration of alternative representations for wavefront reconstruction that can leverage prior knowledge about the optical system to mitigate overfitting and provide more robust reconstructions.

With prior knowledge about the system, one can choose a more suitable basis that avoids the problems of overfitting noise; for example, the often used Zernike basis [10–15]. The Zernike basis can be intuitively understood as a polar coordinate adaptation of polynomial regression, optimized for representing wavefront aberrations in circular optical systems. By design it is directly related to common optical surface aberrations like coma and astigmatism, making it highly suitable to describing wavefronts that were generated with optical components. Defined on a unit disk, the Zernike polynomials provide an efficient way to handle rotationally symmetric systems, common in optical applications. The Zernike modes are also somewhat hierarchical and by truncating the coefficients one can reliably calculate standard aberrations. For example, a laser wavefront can typically be described by a number of Zernike modes on the order of 10, which represents many fewer coefficients than required in the pixel basis. However, the efficiency of the Zernike basis is diminished when confronted with more complex or non-standard aberrations, for example those with sharp edges; Fig. 1(a-b) displays the benefit of employing the Cartesian or Zernike basis in certain scenarios.

There are other wavefronts which are not efficiently captured by either the Zernike or Pixel basis, such as optical vortices [16], which possess an azimuthally varying phase that remains radially constant. The phase profile of an optical vortex can be written as $\phi(\theta) = \ell\theta$, where ℓ is the topological charge and θ is the azimuthal angle. Despite its very simple form, this phase profile cannot be efficiently captured by either Cartesian or Zernike polynomials, a fact demonstrated in Fig. 1(c). In addition to optical vortices, myriad other complex wavefront morphologies, such as Bessel beams [17], Laguerre-Gaussian beams, and Airy beams [18], each carrying unique phase and amplitude profiles, are often encountered in optics. Each of these can be captured via concise mathematical descriptions, but require complex expressions in other bases. Therefore, it is apparent that a more versatile and encompassing framework is required for wavefront reconstruction.

Over-complete dictionaries offer a promising solution to this challenge. In contrast to traditional basis sets which are typically complete and orthogonal, over-complete dictionaries contain a surplus of basis functions. While traditional basis sets allow each signal to be represented in one unique way, over-complete dictionaries offer multiple representations for the same signal, providing additional flexibility. The reconstruction process is then designed to find the most

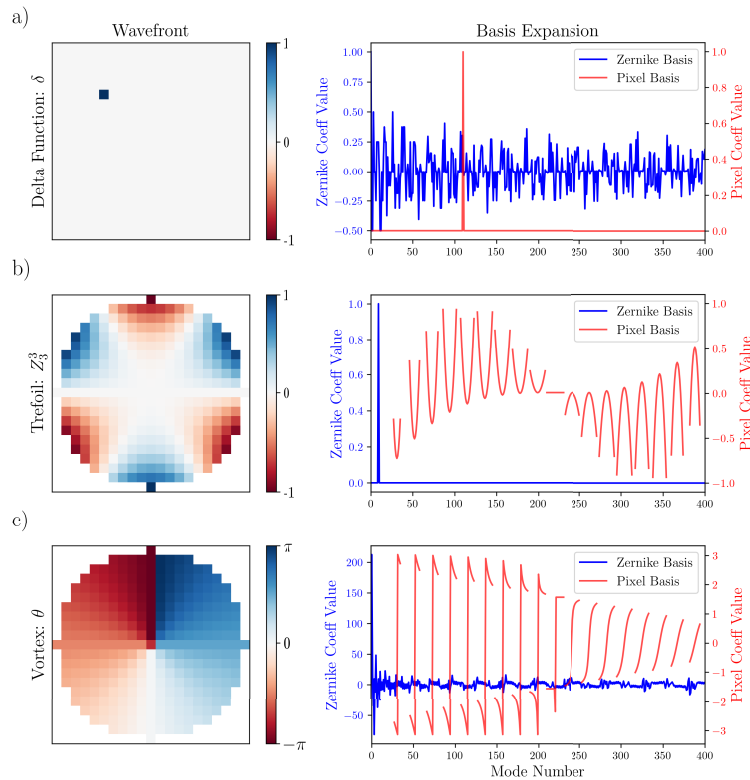


Fig. 1. The expansion of a number of wavefronts in terms of the Pixel and Zernike bases. In order to describe wavefront function, Ψ , the coefficient for mode n in basis Φ was calculated by $c_n = \langle \Psi | \Phi_n \rangle$. a) To express a single pixel requires many Zernike modes. b) A single Zernike mode requires many modes in the Pixel basis. c) There exist wavefronts, such as the optical vortex, which are not efficiently represented in either the Zernike or Pixel basis. The Zernike modes are ordered according to their Noll index [19], and the colorbar applies for all wavefronts.

efficient representation of the signal by adopting methods from compressed sensing and sparse signal representation. Intuitively, this is achieved by favoring solutions that minimize the L0 norm - the number of non-zero modal coefficients. As shall be discussed, the L0 norm is not easy to optimize; fortunately the minimization of the L1 norm results in the same solution in many situations. This approach enables the finding of the most efficient representation of the waveform.

Capitalizing on this theory, we propose the construction of an over-complete phase dictionary, a versatile set of basis functions that extend beyond the capabilities of Cartesian and Zernike polynomials. The proposed dictionary, rich in the variety of basis functions it offers, aims to efficiently handle a wide array of wavefront morphologies.

This paper explores the design of the over-complete phase dictionary, delves into the sparse reconstruction methodology using the dictionary, and its potential implications. This way, we provide a more accurate, interpretable, and computationally efficient representation of wavefronts, pushing the boundaries of traditional wavefront reconstruction methodologies.

2. Theory

Traditional methodologies for wavefront reconstruction have primarily relied on two types of basis sets, known as zonal and modal reconstruction, respectively: the Cartesian and Zernike

basis [20]. While each has its merits, they both possess inherent limitations that render them inadequate for capturing complex wavefronts. The constraints observed in these approaches call for the formulation of a comprehensive, versatile basis set that efficiently handles a broad array of wavefront morphologies.

The following considers the reconstruction of the wavefront from measurements of its gradient, \vec{s} , which is the relevant reconstruction process for techniques such as the Shack-Hartmann sensor and lateral shearing interferometry. Using a finite difference derivative matrix, \mathbf{D} , the forward process can be described using the vectorised (flattened) form of the wavefront, \vec{w} , as, $\vec{s} = \mathbf{D}\vec{w}$. The wavefront can be expressed in any complete basis via, $\vec{w} = \mathbf{B}\vec{c} = \sum_i (\vec{B}_i)c_i$, where \vec{c} and \mathbf{B} are the modal coefficients and basis functions respectively. Combining these equations, and defining a new matrix as the derivatives of a set of basis functions, $\mathbf{D}_B = \mathbf{D}\mathbf{B}$, one arrives at,

$$\vec{s} = \mathbf{D}_B\vec{c}. \quad (1)$$

The reconstruction task is to find the coefficients, \vec{c} , given \vec{s} and \mathbf{D}_B . The key differences in existing approaches centers on the choice of basis functions, \mathbf{B} , as is discussed in the following section.

2.1. Zonal and modal reconstruction

Zonal reconstruction represents the wavefront using a piecewise function, typically in the Cartesian basis. This approach excels at resolving high-frequency spatial features but suffers from issues like overfitting and sensitivity to noise.

$$W(x, y) = \sum_{i=0}^n \sum_{j=0}^m c_{ij} \delta(x - i) \delta(y - j) \quad (2)$$

Here, i, j, x, y describe pixel coordinates, $W(x, y)$ is the wavefront, and c_{ij} are the coefficients.

Modal reconstruction, on the other hand, employs global functions, such as Zernike polynomials, that span the entire aperture. This approach is particularly effective for describing rotationally symmetric systems and is less susceptible to overfitting. A wavefront is expanded in the Zernike polynomials, Z_{nm} , according to,

$$W(R, \theta) = \sum_{n=0}^{\infty} \sum_{m=-n}^n c_n^m Z_{nm}(R, \theta), \quad (3)$$

where R and θ are the radial coordinate and azimuthal angle respectively, where R is inside the pupil ($R \in [0, 1]$), and c_n^m are the corresponding coefficients.

2.2. Limitations of zonal and modal bases

Each basis set has its limitations. The Cartesian basis is often inefficient for low-frequency components due to its high dimensionality. Conversely, Zernike polynomials struggle to capture high-frequency details or non-standard aberrations. As a result, researchers are often forced to compromise between robustness and representational power. In the following two simple examples are discussed, in which the wavefront can be easily parameterized under a specific basis, but is inefficiently encoded in the common zonal and modal approaches.

2.2.1. Optical vortex

Consider an optical vortex, with a wavefront described by $\ell\theta$. Trivially, in the Cartesian (pixel) basis, one requires a number of elements equal to the number of pixels to describe this wavefront.

Less obvious is that the Zernike basis also requires more coefficients. As the Zernike polynomials form a complete basis, one can expand any function in this basis,

$$f(\theta, R) = \sum_{m,n} c_n^m Z_n^m \quad (4)$$

To find the expansion coefficients, one simply calculates the overlap with the given Zernike and the function over the domain which Zernikes are orthogonal, as can be seen below;

$$\int_{-\pi}^{\pi} f(\theta, R) Z_n^m = \int_{-\pi}^{\pi} \left(\sum_{m',n'} c_{n'}^{m'} Z_{n'}^{m'} \right) Z_n^m = c_n^m \quad (5)$$

Setting the function equal to the vortex, $f(\theta, R) = \ell\theta$, one may attempt to expand it in terms of Zernike polynomials:

$$\ell\theta = \sum_{n=0}^{\infty} \sum_{m=-n}^n c_n^m Z_n^m(R, \theta) \quad (6)$$

Detailed analysis (see appendix) shows that one needs the complete set of coefficients c_n^m , each proportional to $1/m$ in magnitude. Thus, the representation of this seemingly simple function is very inefficient. In practice, this problem is amplified by the fact that higher order Zernike modes are numerically unstable. This can only to some extent be mitigated using recursive definitions [21], making it impractical to describe wavefronts with more than Zernike modes.

2.2.2. Misaligned Zernike modes

Consider a wavefront that can be represented by a single Zernike coefficient, e.g., the oblique trefoil mode, $Z_3^3(x, y)$, which is described in polar and cartesian coordinates as,

$$Z_3^3 = R^3 \cos(3\theta) = x^3 - 3xy^2 \quad (7)$$

Now consider a small shift, $x \rightarrow x + \Delta$, so the Zernike mode is offset from the center.

$$Z_3^3(x + \Delta, y) = (x + \Delta)^3 - 3(x + \Delta)y^2 \quad (8)$$

$$= x^3 - 3xy^2 + x^2(-3\Delta) + x(3\Delta^2) - \Delta^3 - 3\Delta y^2 \quad (9)$$

$$= Z_3^3(x, y) + \underbrace{(x^2 - y^2)(-3\Delta)}_{n=2} + \underbrace{x(3\Delta^2)}_{n=1} - \underbrace{\Delta^3}_{n=0} \quad (10)$$

This demonstrates that a simple coordinate shift of a Zernike coefficient in row N of the pyramid will introduce additional terms from rows $[0, \dots, N - 1]$; a basis that was sparse will now be very dense. One also notes for a row, n , the shift factor is proportional to Δ^{N-n} . As $\Delta < 1$ (the Zernike basis is on a unit circle), the higher order modes are added with greater weighting. To generalize this observation, one employs a Taylor expansion of the shifted Zernike mode:

$$Z_N^M(x + \Delta, y) = Z_N^M(x, y) + \Delta \frac{\partial Z_N^M}{\partial x} + \frac{\Delta^2}{2!} \frac{\partial^2 Z_N^M}{\partial x^2} + \dots \quad (11)$$

The derivatives of the Zernike polynomials can be expressed as linear combinations of Zernike polynomials of lower order. Therefore, the first derivative is

$$\frac{\partial Z_N^M}{\partial x} = \sum_{n=N-1}^0 \sum_m a_{nm} Z_n^m(x, y) \quad (12)$$

and similarly, higher-order derivatives involve even lower-order polynomials. Substituting back into the Taylor expansion, one obtains:

$$Z_N^M(x + \Delta, y) = Z_N^M(x, y) + \Delta \sum_{n=N-1}^0 \sum_m a_{nm} Z_n^m(x, y) + \frac{\Delta^2}{2} \sum_{n=N-2}^0 \sum_m b_{nm} Z_n^m(x, y) + \dots \quad (13)$$

This series demonstrates that a small shift Δ introduces contributions from lower-order Zernike modes down to $n = 0$. The coefficients a_{nm} and b_{nm} depend on the specific mode and its derivatives. Even though Δ is small, the impact on the representation is significant because the higher-order terms decay slowly due to the factorial in the denominator being offset by the increasing number of contributing lower-order modes.

Thus, a wavefront that was sparsely represented by a single Zernike coefficient when centered becomes densely represented upon shifting. This density arises from the need to account for the introduced lower-order modes to accurately describe the shifted wavefront. The practical implication is that even minor misalignments in optical systems, or in the detection process, can lead to significant inefficiencies in wavefront representation and interpretation when using Zernike polynomials.

This effect underscores the limitations of the Zernike basis in situations where the wavefront is not perfectly centered or aligned, which is often possible in the case of an imperfectly-circular beam. It highlights the necessity for a more adaptable representation that maintains sparsity despite shifts or misalignments.

2.3. Over-complete dictionaries

Over-complete dictionaries [22,23] provide a solution to the problems outlined above. These dictionaries contain more basis functions than the dimensions in the space they represent, offering flexibility to adapt to a variety of signals while maintaining robustness.

$$W(x, y) = \sum_{i=1}^N c_i \phi_i(x, y) \quad (14)$$

Here, ϕ_i are the basis functions from the over-complete dictionary, and c_i are the coefficients. The overcompleteness of our dictionary presents a fundamental challenge: the existence of multiple valid solutions for a given wavefront. This non-uniqueness might seem to undermine the very purpose of our approach. However, it is precisely this redundancy that allows one to overcome the limitations of single-basis representations that have been previously discussed. The pursuit of sparsity in this context is not merely a mathematical convenience, but a reflection of the underlying physics of wavefronts. Throughout this paper, it has demonstrated that different physical phenomena in optics - be it Zernike-type aberrations, optical vortices, or misaligned beams - naturally lend themselves to sparse representations in appropriate bases. This sparsity is a consequence of the wave equation and the boundary conditions typical in optical systems. The limitations of individual bases in efficiently representing diverse wavefront phenomena stem from the fundamental nature of optical wavefronts. These are in most cases solutions to the paraxial wave equation:

$$\frac{\partial^2 \psi}{\partial x^2} + \frac{\partial^2 \psi}{\partial y^2} + 2ik \frac{\partial \psi}{\partial z} = 0 \quad (15)$$

where $\psi(x, y, z)$ is the complex amplitude of the field, k is the wavenumber, and z is the propagation direction. The solutions to this equation encompass a wide range of phenomena, including

Gaussian beams, Hermite-Gaussian modes, Laguerre-Gaussian modes (which include optical vortices), and Bessel beams, among others.

The pursuit of sparsity in our over-complete dictionary is fundamentally a search for the most appropriate set of solutions to the paraxial wave equation that describe the observed wavefront. This approach can be formalized mathematically as an optimization problem:

$$\min_{\mathbf{c}} |\mathbf{c}|_0 \quad \text{subject to} \quad W = \Phi \mathbf{c} \quad (16)$$

Here, $|\mathbf{c}|_0$ denotes the L0 norm, which quantifies the number of non-zero elements in the coefficient vector \mathbf{c} , and Φ is the over-complete dictionary of modes. This formulation seeks to identify the minimal set of fundamental modes that accurately represent the observed wavefront.

However, L0 minimization is NP-hard, leading us to the more tractable L1 norm minimization, as proposed by Candés et al. [24]:

$$\min_{\mathbf{c}} \|\mathbf{c}\|_1 \quad \text{subject to} \quad W = \Phi \mathbf{c} \quad (17)$$

This relaxation, while computationally necessary, still promotes sparsity and often yields solutions very close to the true L0-sparse solution. Candés and colleagues demonstrated that under certain conditions, L1 minimization can exactly recover the sparsest solution, providing a theoretical foundation for compressed sensing and sparse signal reconstruction [24,25].

Crucially, the paraxial wave equation is defined with respect to an optical axis, which in real optical systems may not perfectly align with the measurement apparatus. As discussed in Sec.2.2.2, this misalignment can lead to apparent complexity in the wavefront when viewed from the perspective of the measurement system. By including the center coordinates, (Δ_x, Δ_y) , as variables in our optimization, one may account for this potential misalignment:

$$\min_{\mathbf{c}, \Delta_x, \Delta_y} |\mathbf{c}|_1 \quad \text{subject to} \quad W = \Phi(x - \Delta_x, y - \Delta_y) \mathbf{c} \quad (18)$$

This formulation allows us to simultaneously determine the most appropriate set of paraxial modes and the true optical axis of the system. It recognizes that apparent complexity in the wavefront may arise from a simple misalignment rather than intrinsic high-order aberrations or complex phase structures. When we introduce the idea of fitting the center during reconstruction, we are implicitly expanding our dictionary to include shifted versions of each basis function. Mathematically, this can be expressed using set notation as:

$$\Phi_{expanded} = \{\phi_i(x - \Delta_x, y - \Delta_y) \mid \phi_i \in \Phi, (\Delta_x, \Delta_y) \in \mathbb{R}^2\} \quad (19)$$

That is, the set of modes in the expanded dictionary contains all possible shifts of each basis function, and is thus infinite-dimensional. However, it's important to note that this expansion doesn't fundamentally change the nature of our problem - it remains an over-complete dictionary, just with an even higher degree of redundancy. The key insight is that by allowing for these shifts, we're enabling our reconstruction to find even sparser representations. A misaligned Zernike mode, which might require many coefficients in the original dictionary, can now be represented by a single coefficient in the expanded dictionary.

This formulation connects seamlessly with the theory of over-complete dictionaries and sparse representations. We're still seeking the sparsest representation, but now in an expanded dictionary that can capture a wider range of physical phenomena efficiently. The concept of dictionary learning in sparse coding literature provides a theoretical framework for understanding this approach. Just as dictionary learning algorithms adapt the basis functions to better represent a class of signals, our method adapts the positioning of the basis functions to better represent the specific wavefront at hand. This adaptive approach not only allows for more accurate

reconstruction of misaligned wavefronts but also provides valuable information about the optical system itself. The optimal (Δ_x, Δ_y) values can indicate misalignments in the optical setup or in the detection process, offering diagnostic capabilities beyond mere wavefront reconstruction.

We refer to our technique by the acronym: PROD (Phase Reconstruction using an Over-complete Dictionary). In the following sections, we will present numerical experiments demonstrating how this expanded dictionary approach, combined with L1-based sparse reconstruction, outperforms static reconstructions across a wide range of wavefront types. We'll explore its robustness to various aberrations and misalignments, showcasing how the pursuit of sparsity in an appropriately designed over-complete dictionary can lead to both more accurate and more interpretable wavefront reconstructions.

3. Numerical simulations

3.1. Construction of the over-complete dictionary

Given a set of wavefront derivatives, with spatial dimensions (N_x, N_y) , the over-complete dictionary is created by the synthesis of a number of modes from different bases. Added first is the pixel basis, with the introduction of $N_x \times N_y$ modes, allowing the technique to represent any discontinuities. Then, to provide an efficient way to represent the common low frequency features, a set number of Zernike modes are added, where the derivatives are found analytically [21]. Finally, with prior knowledge of the optical system being measured, one can simply add 'special' modes that are expected to be present in the wavefront, such as the optical vortex ($\phi = \theta$). The topological charge of the vortex will then be the modal coefficient that is found by the PROD. As has been described, the pixel basis is already sufficient to describe any wavefront. However, in doing so, it will also capture noise, and it does not provide useful physically interpretable modal coefficients; for example the order of a vortex.

Once the over-complete dictionary has been formed by the concatenation of the modal derivatives, a spatial transformer network is also initialized [26]. This is a module which performs a translation of the Zernike modes and a rigid transform of the special modes, with trainable parameters for the rotation angle, θ , and translations, (Δ_x, Δ_y) , enabling the most efficient representation of the wavefront. The Zernike modes and the special modes are each given their own set of rigid transform parameters, to account for the potentially different sources of these distortions. The trainable modal coefficients and rigid transform parameters are then optimized by minimizing the error with the measured wavefront derivatives whilst simultaneously minimizing the L1 norm of the coefficients, according to Section 3.2.

In practice, at initialization we create multiple sets of random coefficients and random (or grid) rigid transform parameters, and optimize them in parallel. By then choosing the one that achieved the minimum cost, the robustness of the technique is increased, without a dramatic increase in computational time due to the properties of GPUs. Due to the rigid transformation, the spatial domain which the over-complete dictionary is defined on is larger than the original derivatives, and depends on the upper bounds of the translation and rotation parameters; during optimization, the predicted wavefront derivatives are simply cropped to the same size as the measured ones for the cost calculation.

3.2. Optimization framework

While L1 minimization problems are traditionally approached using convex optimization algorithms like Iterative Shrinkage-Thresholding Algorithms (ISTA) [27], Interior Point Methods (IPM) [28] or the Alternating Direction Method of Multipliers (ADMM) [29], we opt for the Adaptive Moment Estimation (ADAM) optimizer [30] implemented in PyTorch. The complexity of our optimization problem, particularly with the introduction of center fitting, makes it potentially non-convex. ADAM's ability to handle both convex and non-convex problems effectively makes

it a suitable choice. Moreover, ADAM's efficient GPU-accelerated implementation in frameworks like PyTorch can in practice outperform traditional convex optimization algorithms without similar optimization. We formulate our sparse optimization problem as:

$$\min_{\mathbf{c}, \Delta_x, \Delta_y} \|\mathbf{W} - \Phi(x - \Delta_x, y - \Delta_y)\mathbf{c}\|_2^2 + \lambda \|\mathbf{c}\|_1 \quad (20)$$

where \mathbf{W} is the measured wavefront (or more precisely its slopes), Φ is our over-complete dictionary, \mathbf{c} are the coefficients we're optimizing, δ_x and δ_y are the center coordinates, and λ is a regularization parameter. This approach combines the flexibility to handle potential non-convexity, the speed of GPU acceleration, and the sparsity-inducing properties required for our over-complete wavefront reconstruction problem. The following experiments were performed on a single NVIDIA GeForce RTX 3090 GPU.

3.3. Results

3.3.1. Axicon phase

Firstly, PROD is utilised to perform phase stitching of a wavefront not efficiently expressed in the Zernike basis: a simple linear radial phase, $\phi(r) \propto r$. Such a wavefront can be generated by an Axicon [31], a type of lens formed from a conical piece of glass. Due to its geometry, there is a linear dependence in the radial position and the amount of glass that the pulse has traveled through, resulting in the radial wavefront. When used to focus a Gaussian beam, it creates a Bessel-like beam which does not experience diffraction over a region of interest [32], making it of interest in a wide range of applications.

Here, we create a wavefront by randomly sampling Zernike coefficients from a Normal distribution $\sim \mathcal{N}(0, 0.5)$, up to the 28th mode, before adding the Axicon phase. To test the ability of the PROD technique to handle noise, 3 sampling scenarios are simulated. Normally distributed noise is added to the gradients, with a standard deviation equal to a percentage, $\eta \in [0\%, 2\%, 5\%]$, of the gradients maximum value. The learning rate of the ADAM optimizer was set to 3×10^{-3} , and the regularisation parameter was set to $\lambda = 7 \times 10^{-4}$. The results are shown in Fig. 2.

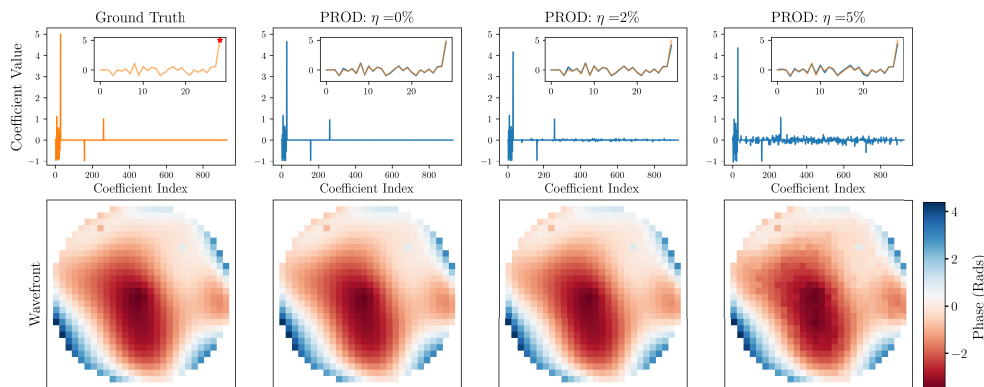


Fig. 2. The wavefront was synthesized by the random initialization of 28 Zernike modes, and the Axicon mode. The top row displays the modal coefficient values, with each subplot showing a zoom of the first 29 coefficients - the Zernikes and Axicon (marked with a red star). The bottom row shows the predicted wavefront, with all plotted on the same color scale. One sees that even in the presence of significant noise (parameterized by a percentage, η , of the maximum value of the derivative), PROD is able to accurately extract the true coefficient values.

It is evident that as the noise level was increased, PROD was able to use the pixel basis in order to fit the noise, meaning that the predictions of the Zernike and Axicon mode remained fairly constant throughout the noise levels, proving the robustness of the technique.

3.3.2. Shifted oblique trefoil

This experiment demonstrates the benefits of the rigid transform aspect of the approach. The phase map is implemented as an off-centered Z_3^3 Zernike mode (oblique trefoil), with 2 pixels of shot noise. Considering the Zernike modes are defined from $[-1,1]$, the new center position is chosen as $[0.2, 0.2]$. Normally distributed random noise ($\eta = 2\%$) was added to the sampled derivatives, before the proposed approach was used to recover the wavefront from the derivatives. For comparison, we use the scikit-learn [33] implementation of LASSO regression [34]. This is a commonly implemented technique to solve Eq. 17, and is described in the appendix. The technique is given the same over-complete dictionary, but possesses no ability to perform the transformation of modes. LASSO is parameterized by a L1 penalization parameter, α , which controls the sparsity of the solution. This was adjusted and the results are displayed in Fig. 3.

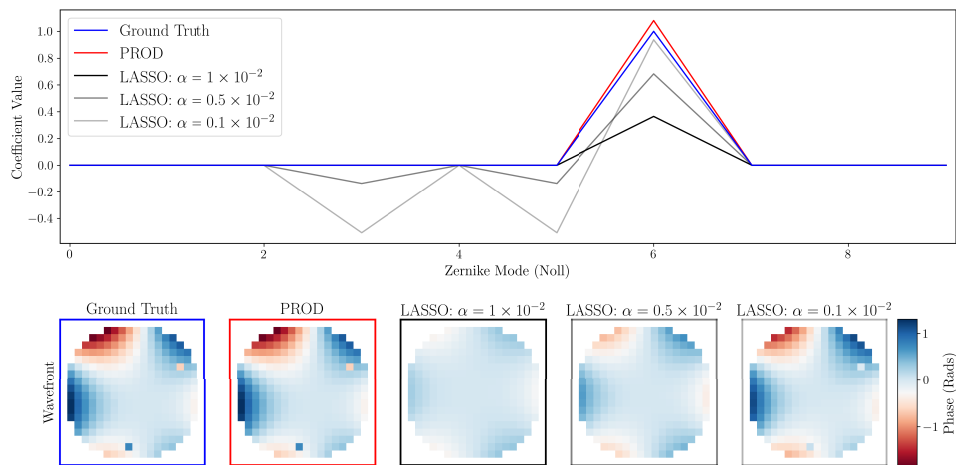


Fig. 3. Demonstrating the utility of PROD's rigid transform module to efficiently represent de-centered modes, in this case the Z_3^3 mode. Here PROD learns the de-center parameter, and identifies the single Zernike mode. Also displayed for comparison is the LASSO technique, which is given the same over-complete dictionary, but possesses no ability to transform the modes. All wavefronts are plotted on the same color scale. If its L1 penalization parameter, α , is large, it cannot represent the wavefront accurately, whereas if α is small, it must utilise lower order Zernike modes to help represent the wavefront.

As alluded to in Section 2.2.2, a shifted Zernike mode introduces other lower order Zernike terms into the expansion. This is seen for the LASSO technique; as α is decreased, more low order Zernike modes are used to help represent the function. As the proposed approach uses a trainable rigid transform, it finds the most efficient representation using just the single Zernike mode.

4. Physical experiments

An experiment was performed to measure an optical vortex. This was performed using a Spectra-Physics Spitfire ACE laser system to generate ultrashort pulses with central wavelength, $\lambda_0 = 798$ nm, and a transform limited duration of 64 fs. To generate the vortex, a setup was used that has been thoroughly discussed in previous work [35,36]. Firstly, a quarter waveplate (QWP)

was used to circularly polarise the incoming pulse, before it passed through a ($l = 2$) structured waveplate (s-waveplate), creating two circularly polarized vortices with opposing handedness. A further QWP was used to transform them into two linearly polarised vortices before a linear polariser was used to select one of them. The optical components used here were achromatic, so the signal was integrated over time when captured. The wavefront measurement was then performed on a home-built Shack Hartmann sensor (focal length, $f = 14.2$ mm, and pitch, $\Lambda = 300 \mu\text{m}$) [37]. The microlens focii for the vortex are seen in Fig. 4.

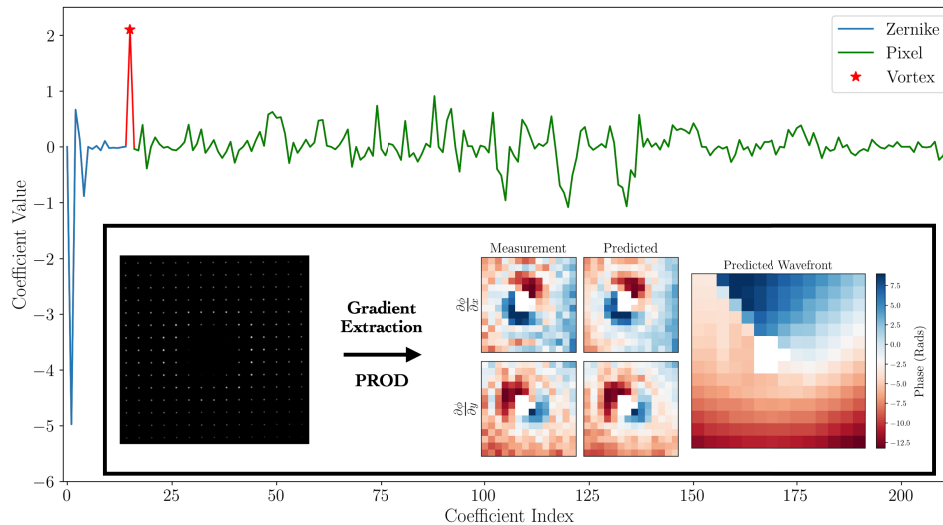


Fig. 4. The use of PROD to characterize an experimental optical vortex. Once the gradients are extracted from the Shack-Hartmann pattern, the technique is used to find the modal coefficients for the vortex, and lower order Zernike modes. The coefficient for the vortex (the topological charge) was 2.09, in good agreement with the theoretical value for this setup of 2.

A reference microlens focii pattern was first captured, before the s-waveplate was introduced to the optical setup. The centroids of each microlens focii were extracted, before the wavefront gradients were found by the calculating the centroid difference between the vortex and the reference measurements. The PROD technique was then used to extract the modal coefficients and the wavefront. The results are displayed in Fig. 4. We see that the PROD technique was able to correctly identify the 2nd order vortex, with a predicted modal coefficient of 2.09. This discrepancy of 5% relative to the theoretical value can be attributed to misalignment in the distance between the sensor and the microlens array, and imperfection in the vortex generation.

5. Conclusions and outlook

Leveraging an over-complete dictionary and trainable rigid transformations, PROD has been introduced as a robust technique to reconstruct a wavefront from measurements of its gradients. Where most existing techniques typically reconstruct in the pixel or the Zernike basis, there exist limitations with both however; the former overfits to noise and doesn't give physically interpretable coefficients, and there are some specialized wavefronts that cannot be efficiently expressed in the Zernike basis, such as the optical vortex. Exploiting one's prior knowledge about a measurement, PROD is able to extract the physically-meaningful modal coefficients for the modes that are within its over-complete dictionary. Furthermore, as its dictionary also incorporates the pixel basis, it can express noise without influencing the other modal coefficients, such as those for Zernikes and special modes.

In its presented form, a PROD reconstruction of a 20×20 wavefront takes on the order of seconds, meaning it is principally for post-analysis. However, if one assumes that all modes are centered and removes the rigid transform module of this technique, LASSO can be used to solve the L1 problem in a far shorter time (ms), allowing a real-time reconstruction for applications in adaptive optics. While this may be suitable in some cases, such as when one is certain that the system is centered, in the most general case this is not a suitable assumption, for example with a decentered vortex, and can lead to an under-sparse representation, and resultingly inaccurate modal coefficients.

This application of overcomplete dictionaries and the principals of sparsity can also be applied to other problems in laser physics. For example, the reconstruction of the spectral phase from measurements can also be performed in different bases, and sometimes involves switching between bases [38,39]. Furthermore, with sufficient computational resources, there is no reason why this technique cannot be extended to recent work in the spatio-spectral regime, where the wavefront is considered in three dimensions [37].

Appendix A

6.1. Representation of vortex in Zernike basis

We aim to expand the optical vortex phase function $\phi(\theta) = \ell\theta$ in terms of Zernike polynomials $Z_n^m(\rho, \theta)$. The Zernike polynomials form a complete set over the unit disk, allowing any function to be represented as:

$$\phi(\rho, \theta) = \sum_{n=0}^{\infty} \sum_{m=-n}^n c_n^m Z_n^m(\rho, \theta)$$

Due to the orthogonality of the Zernike polynomials, the coefficients c_n^m are given by:

$$c_n^m = \frac{\int_0^1 \int_0^{2\pi} \phi(\theta) Z_n^{m*}(\rho, \theta) \rho \, d\theta \, d\rho}{\int_0^1 \int_0^{2\pi} |Z_n^m(\rho, \theta)|^2 \rho \, d\theta \, d\rho}$$

Since $\phi(\theta)$ is independent of ρ , we can separate the radial and angular components. The denominator simplifies to π due to normalization, and we define the radial integral:

$$S_n^{|m|} = \int_0^1 R_n^{|m|}(\rho) \rho \, d\rho$$

where $R_n^{|m|}(\rho)$ are the radial Zernike polynomials. Our main task is then evaluate the angular integral:

$$A_m = \int_0^{2\pi} \phi(\theta) e^{-im\theta} \, d\theta = \ell \int_0^{2\pi} \theta e^{-im\theta} \, d\theta$$

This integral can be evaluated using integration by parts. Let $u = \theta$ and $dv = e^{-im\theta} d\theta$, so $du = d\theta$ and $v = \frac{e^{-im\theta}}{-im}$. Applying integration by parts:

$$A_m = uv \Big|_0^{2\pi} - \int_0^{2\pi} v \, du = \left[\theta \cdot \frac{e^{-im\theta}}{-im} \right]_0^{2\pi} - \int_0^{2\pi} \frac{e^{-im\theta}}{-im} \, d\theta = \frac{2\pi}{-im}$$

Substituting back, the coefficient c_n^m becomes:

$$c_n^m = \frac{\ell S_n^{|m|} A_m}{\pi} = \frac{\ell S_n^{|m|} \left(\frac{2\pi}{-im} \right)}{\pi} = \ell \frac{2i}{m} S_n^{|m|}$$

This shows that the coefficients c_n^m decay proportionally to $1/m$. The slow decay of c_n^m with increasing m indicates that a large number of Zernike modes are required to accurately represent

the vortex phase $\phi(\theta) = \ell\theta$. Therefore, the Zernike basis is inefficient for representing optical vortices, motivating the use of alternative representations such as an over-complete dictionary.

6.2. LASSO regression

The least absolute shrinkage and selection operator (LASSO) is a technique that leverages the L1 norm to find regularised solutions to least squares problems. Consider one has a measurement vector, \vec{y} , and a dictionary, \mathbf{x} . One wishes to find the sparsest possible coefficient vector c , subject to the condition that $\vec{y} = \mathbf{x}\vec{c}$. The minimization problem described by LASSO is,

$$\min_c \{ \|\vec{y} - \mathbf{x}\vec{c}\|_2 + \alpha \|\vec{c}\|_1 \} \quad (21)$$

where $\|a\|_2 = \sum_j a_j^2$, $\|c\|_1 = \sum_j |c_j|$, and α is a constant determining the level of sparsity.

The scikit-learn implementation solves this problem using coordinate descent. Here, the modal coefficients are optimized one at a time [40]. To optimize parameter c_i , we first define \mathbf{x}_{-i} and c_{-i} as the dictionary and coefficients with index i removed, meaning the remaining residual for component i is found by $(y - \mathbf{x}_{-i}c_{-i})$. Solving the LASSO equation for c_i gives,

$$c_i = \frac{\mathbf{x}_i^T (y - \mathbf{x}_{-i}c_{-i})}{\mathbf{x}_i^T \mathbf{x}_i} - \frac{\alpha}{\mathbf{x}_i^T \mathbf{x}_i} \frac{\partial |c_i|}{\partial c_i} \quad (22)$$

where $\frac{\partial |c_i|}{\partial c_i} = \text{sign}(c_i)$. This equation clearly has a term from the least squares, and a term from the L1 regularisation. To enforce sparsity, a soft threshold is used to set the parameter to zero if its value is under a threshold.

$$c_i = \text{sign}(c_i) \cdot \min \left(\left| \frac{\mathbf{x}_i^T (y - \mathbf{x}_{-i}c_{-i})}{\mathbf{x}_i^T \mathbf{x}_i} \right| - \frac{\alpha}{\mathbf{x}_i^T \mathbf{x}_i}, 0 \right) \quad (23)$$

This process is iterated over all coefficients $i \in [1, N_c]$, and is then repeated until the final error falls below a predefined threshold, or a predefined termination number of loops is reached.

Funding. Deutsche Forschungsgemeinschaft (453619281); Science and Technology Facilities Council (ST/V001655/1); Ministerio de Ciencia, Innovación y Universidades (PID2020-119818GB-I00, PID2023-149836NB-I00); Consejería de Educación, Junta de Castilla y León (SA136P20).

Acknowledgments. We would like to acknowledge useful discussions with the groups of Professor Peter Norreys and Dr. Andreas Döpp.

This work was supported by the Independent Junior Research Group "Characterization and control of high-intensity laser pulses for particle acceleration", DFG Project No. 453619281. We would also like to acknowledge UKRI-STFC grant ST/V001655/1, and the following funding sources: European Regional Development Fund and Consejería de Educación, Junta de Castilla y León (SA108P24); Ministerio de Ciencia e Innovación (PID2020-119818GB-I00, PID2023-149836NB-I00).

Disclosures. The authors declare no conflicts of interest.

Data availability. An implementation of the technique is publicly available [41].

References

1. R. W. Gerchberg, "A practical algorithm for the determination of phase from image and diffraction plane pictures," *Optik* **35**, 237–246 (1972).
2. C. Zuo, J. Li, J. Sun, *et al.*, "Transport of intensity equation: a tutorial," *Optics and Lasers in Engineering* **135**, 106187 (2020).
3. J. Primot and L. Sogno, "Achromatic three-wave (or more) lateral shearing interferometer," *J. Opt. Soc. Am. A* **12**(12), 2679–2685 (1995).
4. B. C. Platt and R. Shack, "History and principles of shack-hartmann wavefront sensing," (2001).
5. F. Dai, J. Li, X. Wang, *et al.*, "Exact two-dimensional zonal wavefront reconstruction with high spatial resolution in lateral shearing interferometry," *Opt. Commun.* **367**, 264–273 (2016).
6. B. Pathak and B. R. Boruah, "Improved wavefront reconstruction algorithm for shack-hartmann type wavefront sensors," *J. Opt.* **16**(5), 055403 (2014).

7. W. H. Southwell, "Wave-front estimation from wave-front slope measurements," *J. Opt. Soc. Am.* **70**(8), 998–1006 (1980).
8. J.-C. Chanteloup, "Multiple-wave lateral shearing interferometry for wave-front sensing," *Appl. Opt.* **44**(9), 1559–1571 (2005).
9. X. Tian, M. Itoh, and T. Yatagai, "Simple algorithm for large-grid phase reconstruction of lateral-shearing interferometry," *Appl. Opt.* **34**(31), 7213–7220 (1995).
10. J. Liang, B. Grimm, S. Goelz, *et al.*, "Objective measurement of wave aberrations of the human eye with the use of a hartmann–shack wave-front sensor," *J. Opt. Soc. Am. A* **11**(7), 1949–1957 (1994).
11. Y. He, Z. Liu, Y. Ning, *et al.*, "Deep learning wavefront sensing method for shack-hartmann sensors with sparse sub-apertures," *Opt. Express* **29**(11), 17669–17682 (2021).
12. L. Seifert, H. Tiziani, and W. Osten, "Wavefront reconstruction with the adaptive shack–hartmann sensor," *Opt. Commun.* **245**(1-6), 255–269 (2005).
13. S. Howard, J. Esslinger, R. H. Wang, *et al.*, "Hyperspectral compressive wavefront sensing," *High Power Laser Sci. Eng.* **11**, e32 (2023).
14. G. Harbers, P. Kunst, and G. Leibbrandt, "Analysis of lateral shearing interferograms by use of zernike polynomials," *Appl. Opt.* **35**(31), 6162–6172 (1996).
15. F. Dai, Y. Zheng, Y. Bu, *et al.*, "Modal wavefront reconstruction based on zernike polynomials for lateral shearing interferometry," *Appl. Opt.* **56**(1), 61–68 (2017).
16. Y. Shen, X. Wang, Z. Xie, *et al.*, "Optical vortices 30 years on: Oam manipulation from topological charge to multiple singularities," *Light:Sci. Appl.* **8**(1), 90 (2019).
17. D. McGloin and K. Dholakia, "Bessel beams: diffraction in a new light," *Contemp. Phys.* **46**(1), 15–28 (2005).
18. N. K. Efremidis, Z. Chen, M. Segev, *et al.*, "Airy beams and accelerating waves: an overview of recent advances," *Optica* **6**(5), 686–701 (2019).
19. R. J. Noll, "Zernike polynomials and atmospheric turbulence," *J. Opt. Soc. Am.* **66**(3), 207–211 (1976).
20. R. K. Tyson and B. W. Frazier, *Principles of adaptive optics* (CRC Press, 2022).
21. T. B. Andersen, "Efficient and robust recurrence relations for the zernike circle polynomials and their derivatives in cartesian coordinates," *Opt. Express* **26**(15), 18878–18896 (2018).
22. S. S. Chen, D. L. Donoho, and M. A. Saunders, "Atomic decomposition by basis pursuit," *SIAM Rev.* **43**(1), 129–159 (2001).
23. R. Rubinstein, A. M. Bruckstein, and M. Elad, "Dictionaries for sparse representation modeling," *Proc. IEEE* **98**(6), 1045–1057 (2010).
24. E. J. Candès, J. Romberg, and T. Tao, "Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information," *IEEE Trans. Inf. Theory* **52**(2), 489–509 (2006).
25. E. J. Candès and T. Tao, "Near-optimal signal recovery from random projections: Universal encoding strategies?" *IEEE Trans. Inf. Theory* **52**(12), 5406–5425 (2006).
26. M. Jaderberg, K. Simonyan, A. Zisserman, *et al.*, "Spatial transformer networks," *Advances in neural information processing systems* **28** (2015).
27. A. Beck and M. Teboulle, "A fast iterative shrinkage-thresholding algorithm for linear inverse problems," *SIAM J. Imaging Sci.* **2**(1), 183–202 (2009).
28. S.-J. Kim, K. Koh, M. Lustig, *et al.*, "An interior-point method for large-scale ℓ_1 -regularized least squares," *IEEE J. Sel. Top. Signal Process.* **1**(4), 606–617 (2007).
29. S. Boyd, N. Parikh, E. Chu, *et al.*, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *FNT in Machine Learning* **3**(1), 1–122 (2010).
30. D. P. Kingma, "Adam: A method for stochastic optimization," *arXiv* (2014).
31. J. H. McLeod, "The axicon: a new type of optical element," *J. Opt. Soc. Am.* **44**(8), 592–597 (1954).
32. V. Garcés-Chávez, D. McGloin, H. Melville, *et al.*, "Simultaneous micromanipulation in multiple planes using a self-reconstructing light beam," *Nature* **419**(6903), 145–147 (2002).
33. F. Pedregosa, G. Varoquaux, A. Gramfort, *et al.*, "Scikit-learn: Machine learning in Python," *J. Machine Learning Research* **12**, 2825–2830 (2011).
34. R. Tibshirani, "Regression shrinkage and selection via the lasso," *J. Royal Statistical Society Series B: Statistical Methodology* **58**(1), 267–288 (1996).
35. I. Lopez-Quintas, W. Holgado, R. Drevinskas, *et al.*, "Optical vortex production mediated by azimuthal index of radial polarization," *J. Opt.* **22**(9), 095402 (2020).
36. M. López-Ripa, Í. J. Sola, and B. Alonso, "Bulk lateral shearing interferometry for spatiotemporal study of time-varying ultrashort optical vortices," *Photonics Res.* **10**(4), 922–931 (2022).
37. N. Weisse, J. Esslinger, S. Howard, *et al.*, "Measuring spatio-temporal couplings using modal spatio-spectral wavefront retrieval," *Opt. Express* **31**(12), 19733–19745 (2023).
38. B. Alonso, W. Holgado, and Í. J. Sola, "Compact in-line temporal measurement of laser pulses with amplitude swing," *Opt. Express* **28**(10), 15625–15640 (2020).
39. M. López-Ripa, Ó. Pérez-Benito, B. Alonso, *et al.*, "Few-cycle pulse retrieval using amplitude swing technique," *Opt. Express* **32**(12), 21149–21159 (2024).
40. R. Tibshirani, "Coordinate descent," *Course Convex Optimization* pp. 10–725 (2022).
41. <https://github.com/sunnyhoward/overdictionary>.