# ANALYSIS AND SYNTHESIS OF THREE-DIMENSIONAL ILLUMINATION USING PARTIAL COHERENCE

A DISSERTATION

SUBMITTED TO THE DEPARTMENT OF ELECTRICAL ENGINEERING

AND THE COMMITTEE ON GRADUATE STUDIES

OF STANFORD UNIVERSITY

IN PARTIAL FULFILLMENT OF THE REQUIREMENTS

FOR THE DEGREE OF

DOCTOR OF PHILOSOPHY

Zhengyun Zhang

August 2011

This dissertation is online at: http://purl.stanford.edu/vn229hj9775

I certify that I have read this dissertation and that, in my opinion, it is fully adequate in scope and quality as a dissertation for the degree of Doctor of Philosophy.

**Marc Levoy, Primary Adviser**

I certify that I have read this dissertation and that, in my opinion, it is fully adequate in scope and quality as a dissertation for the degree of Doctor of Philosophy.

**Mark Horowitz**

I certify that I have read this dissertation and that, in my opinion, it is fully adequate in scope and quality as a dissertation for the degree of Doctor of Philosophy.

**George Barbastathis**

Approved for the Stanford University Committee on Graduate Studies.

**Patricia J. Gumport, Vice Provost Graduate Education**

*This signature page was generated electronically upon submission of this dissertation in electronic format. An original signed hard copy of the signature page is on file in University Archives.*

# Abstract

There exists many devices that generate three-dimensional illumination patterns. Analysis of the capabilities of three major device families using measures of partial coherence reveals specific illumination patterns that cannot be generated by each device family. Ray-based devices cannot achieve patterns with high resolution, coherent holographic devices cannot achieve certain intensity patterns, even in the two-dimensional case, and volumetric devices cannot simulate occlusions and suffer from out-of-focus blur. Synthesis of more versatile illumination patterns is proposed by computing the mutual intensity representation of a desired partially coherent beam from application specifications and then generating the beam using time-multiplexing methods based on coherent modes. The mutual intensity can be computed directly from a simple scene description, or it can be computed through a novel algorithm using nonlinear conjugate gradients from a desired three-dimensional intensity volume. Equivalent coherent mode decomposition representations for the same mutual intensity will be considered in terms of optimality in efficiency. For cases when the computed mutual intensity is fairly incoherent, a new "quasi-Schell" mode decomposition is proposed to reduce the number of patterns needed at the SLM by introducing partially coherent sources. Use of arbitrary partially coherent beams for three-dimensional illumination enables the versatile generation of light patterns not possible with current devices and is a promising new field for exploration.

# Acknowledgements

There are many people to whom I owe thanks for helping me arrive at where I am today. First and foremost, I am indebted to the guidance and assistance that my advisor Marc Levoy provided me during this journey. He is one of the most reasonable professors I have ever met and one of the most passionate about diving headlong into research. His qualities and his guidance inspired me and helped me find a topic that I would be passionate about. Even when he started seeing my research topic wander off and leave his area of expertise, he tried to do what he could for me and introduced me to George Barbastathis, an expert professor in the field I was starting to enter.

George's experience and expertise in the field helped introduce me to the field of optics and its culture, and even the occasional heated debate taught me how to think logically and present lines of reasoning that would be needed to convince the optics research community. Furthermore, at a point when funding was scarce, George's generosity allowed me to continue my research in optics without having to scramble to find other funding sources. Without his selfless interest in me and my research, I would not have been able to complete this dissertation at all.

Just as George has helped me with the specifics and details of the optics field, Mark Horowitz has helped me take a step back and look at the big picture and made me think about problems more intuitively. Mark is a professor in the Electrical Engineering department who often attends Marc's group meetings, and his wise comments often imparted upon me different points of view for the same problem. These realizations effected actual understanding of my research, beyond simple knowledge of the results and conclusions.

I would like to also thank many other professors here at Stanford, including Bruce

The photograph of the dog used in this dissertation is a part of the ImageStack software package developed by Andrew Adams and the other students in Marc's group, and its license text is reproduced below according to the terms of the license:

```
* Copyright (c) 1995-2010, Stanford University
* All rights reserved.
*
* Redistribution and use in source and binary forms, with or without
* modification, are permitted provided that the following conditions are met:
*     * Redistributions of source code must retain the above copyright
*        notice, this list of conditions and the following disclaimer.
*     * Redistributions in binary form must reproduce the above copyright
*        notice, this list of conditions and the following disclaimer in the
*        documentation and/or other materials provided with the distribution.
*     * Neither the name of Stanford University nor the
*        names of its contributors may be used to endorse or promote products
*        derived from this software without specific prior written permission.
*
* THIS SOFTWARE IS PROVIDED BY STANFORD UNIVERSITY ''AS IS'' AND ANY
* EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED
* WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE
* DISCLAIMED. IN NO EVENT SHALL STANFORD UNIVERSITY BE LIABLE FOR ANY
* DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES
* (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES;
* LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND
* ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT
* (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS
* SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.
```

# Contents

# List of Figures

# Chapter 1

# Introduction

Generation of three-dimensional illumination is useful in many application areas, including displays for 3D entertainment devices, microscopes for biology, and photo-lithography used in semiconductor manufacturing. Many devices operating on vastly different principles have been proposed, but they each have different advantages and disadvantages. These devices can be roughly divided into three main families – ray devices, holographic devices and volumetric devices.

The illumination produced by each of these devices can be described mathematically by the second order characteristics of the randomly fluctuating optical field that has been generated. This characterization treats the light produced by each device as a partially coherent field, and the second order characteristics measure the state of coherence of this field. A unified analysis of the capabilities of each device family using the ideas of partial coherence will reveal that each device family covers a different region in the space of all possible illumination patterns.

Furthermore, analysis of partial coherence using what is known as a coherence modes decomposition hints at a time-multiplexing method to control the statistics of a generated beam, at least in the ideal case. Given this control, we would be able to generate many more illumination patterns than what these three existing families of devices would allow. We would also be able to reduce the design of three-dimensional illumination to a software problem – for each application, we would simply need to

compute a partially coherent optical field given a desired set of qualities this three-dimensional illumination needs to have, whereas previously we would need to either pick or design a hardware device in tandem.

## 1.1 Applications

There are many applications for the creation of three-dimensional light patterns. The most obvious class of applications would be autostereoscopic displays – devices that impart visual depth information to the viewer without the need for special glasses. One such application area would be medicine, where a three-dimensional display would enable a physician to intuitively access volumetric data obtained from three-dimensional imaging modalities such as computed tomography (CT) or magnetic-resonance imaging (MRI). Another application area would be in entertainment, where three-dimensional displays impart three-dimensional scenes to the viewer. For example, Nintendo recently released a handheld gaming system called the Nintendo 3DS which displays a three-dimensional scene to the player without the need for special glasses.

One might also consider control of three-dimensional light patterns in designing illumination for a scene. For example, full control over illumination would allow for the reduction of unwanted light/glare when imaging a scene [1]. Control of light in a microscope could also enable the microscopist to be able to adjust contrast and/or enable different imaging modalities [2].

Since light carries energy, three-dimensional light patterns can also be used to selectively interact with matter in a volume. For example, a high intensity light beam can function as optical tweezers, which can be used to grab and move small particles [3, 4]. Another application area would be in the optical stimulation of neurons in a thick tissue specimen through the use of channelrhodopsins, a family of compounds which can change the membrane structure of neurons when hit by light, thereby causing a neuron to fire [5–9]. Lastly, control over the three-dimensional structure of light is needed in photo-lithography, where a light pattern is used to burn a pattern onto photo-resist, which is then indirectly used to etch patterns in silicon [10–12]. The

light beam needs to have a very long depth-of-field (and hence needs to be controlled in 3D) in order to be able to overcome surface variations on the wafer as well as being able to etch through thick photo-resist.

## 1.2 Generation

Just as there are many application areas for three-dimensional light patterns, there also currently exist many ways to generate them. These methods can be roughly classified into three classes – ray, holographic and volumetric.

Ray devices attempt to create a three-dimensional light pattern by controlling the radiance along each ray in space, i.e. the light field [13, 14]. Aperture scanning devices [15] scan a pinhole in the Fourier plane of a time-multiplexed two-dimensional display so that each image displayed creates roughly a parallel bundle of rays emanating from the device, thus creating a set of parallel bundle of rays for each ray direction. Integral imaging [16] devices using microlens arrays or lenticular sheets over a two-dimensional display covert spatial images locally into angular distributions of radiance. The light field illuminator [2] is one such device, where a microlens array coupled with a projector is used to create three-dimensional light patterns in a microscopic specimen. Yet another device is the aforementioned Nintendo 3DS, which uses a parallax barrier instead of a lenticular sheet. A parallax barrier is simply a lenticular sheet where pinholes replace the lenses.

A second class of devices, holographic devices, involves the modulation of a coherent (laser) beam spatially in amplitude and/or phase to create a three-dimensional light pattern through the use of interference and diffraction. Classically, holograms have been created by exposing film to an interference pattern between the coherent light of the source and the coherent light after passing through/reflecting off of an object [17]. The hologram is then viewed by hitting it with coherent light. With the advent of more powerful computers, holograms can also be computed digitally and then printed onto some medium to display. These holograms can produce synthetic patterns unrelated to real scenes and are called computer generated holograms (CGH) [18–20]. Lastly, recent availability of spatial light modulators (SLM), which

allow computers to directly modulate the amplitude and/or phase of an incoming co-
herent beam, have enabled real-time control over the display of holograms [21]. With
regards to the applications mentioned before, optical tweezers are usually generated
by controlling the amplitude and/or phase of a coherent laser beam using SLMs.
Holographic displays have also been a topic of research for 3D television.

The last family of devices are volumetric displays. These devices essentially pro-
duce a set of light emitters in space or the appearance thereof through either reflection,
scattering or emission.

In general, there are two types of volumetric displays – static volume and swept
volume displays. In a static volume display, the volume is filled with an unmoving
medium, and the medium is locally excited and emits light. One example would
be the use of two-photon absorption, where intense laser beams are used to cause
a particle in the medium to absorb two photons (of equal or differing wavelength)
at the same time and then relax, causing fluorescent emission [22]. Finite emission
decay times coupled with a temporally multiplexed excitation source that scans the
entire volume quickly allows the display of a light emitting volume.

In swept volume displays, image points are generated either physically or virtually
in space, and these points are multiplexed in time to generate the appearance of a light
emitting volume due to finite integration times in the sink system, e.g. persistence of
vision in the human eye. One example would be the use of a spinning diffuse screen
upon which temporally multiplexed two dimensional light patterns are projected [23].
The screen scatters the incoming light, creating the illusion of a three-dimensional
light emitting plane that sweeps out an entire volume over time. Another example
would be the use of a temporally multiplexed two-dimensional display coupled with
a varifocal mirror that produces images of the display at different depths over a short
time interval [24].

## 1.3   Outline

We will start by developing a formal definition for the illumination generation problem
and review the concept of partial coherence in Chapter 2. In this chapter, we will

show that for most application areas, coherence statistics of the optical field such as the mutual intensity are sufficient to describe the illumination pattern generated by a given device. The tools introduced in this chapter will then be used in Chapter 3 to analyze the three illumination device families under a unified framework. This analysis will show that each device family has different inherent limitations on the types of illumination patterns that can be generated. More specifically:

- ray devices can only produce limited resolution patterns due to the uncertainty principle,

- coherent holographic devices cannot generate incoherent fields and are unable to reproduce specific intensity patterns due to coherence forcing every point on the wave function to interact with each other, and

- volumetric devices suffer from out-of-focus blur and cannot simulate the appearance of occluders and astigmatic effects.

Chapter 4 will start with a discussion on how to generate arbitrary partially coherent fields using temporally multiplexed spatial light modulators and then proceed to explore various algorithms needed to compute a desired mutual intensity and temporal patterns. A straight-forward algorithm will be presented to compute a desired mutual intensity from a simple scene description, and implementation details will be discussed. For a fairly incoherent scene, many modes (i.e. frames in the temporal SLM pattern) are needed for high quality images. A novel algorithm based on nonlinear conjugate gradients will also be presented for the application of controlling the total intensity of light at a three-dimensional lattice of points (i.e. voxels). Results from two different example test cases suggest that partially coherent fields recreate three-dimensional intensity patterns with higher fidelity than fully coherent fields even when both methods are allotted the same number of degrees of freedom of control. An analysis of the iterated projections algorithms currently used in the optics literature indicates that they do not minimize error in the intensity in the least squares sense, and this conclusion is corroborated by comparisons with the nonlinear conjugate gradients algorithm through a test case.

The final section of Chapter 4 discusses two potential areas of research. First, there exists an infinite number of coherence mode representations for any given mutual intensity pattern, and an optimization problem is proposed for improving the energy efficiency of a temporally multiplexed partially coherent light generation system. Second, a large number of coherent modes are needed to recreate fairly incoherent beams, as shown earlier in the chapter. A new decomposition of a partially coherent beam into an incoherent mixture of partially coherent "quasi-Schell" modes can reduce the number of frames needed in a temporal sequence by using partially coherent sources as opposed to fully coherent sources. This decomposition is formulated as a weighted low-rank rank approximation problem and is applied to the same scene simulation example to reduce the number of frames needed by a factor of 16, demonstrating the potential for this approach. Finally, Chapter 5 summarizes the main points of this dissertation.

The main contributions of this thesis are:

- a unified analysis of the capabilities of three illumination device families using partial coherence. The concepts regarding the uncertainty principle in ray-based devices have been explored indirectly via imaging in a talk at the 2009 IEEE International Conference on Computational Photography [25] and directly through microscope illumination in a journal article in the Journal of Microscopy [2]. A simpler form of the proof on impossible two-dimensional patterns has been presented at the 2011 OSA Digital Holography and Three-Dimensional Imaging conference [26] and a full version is currently being considered for submission to a journal.

- a novel algorithm for computing a desired partially coherent beam from a set of intensity constraints in a three-dimensional volume. Work is in progress for submission to JOSA A.

- a method of decomposing a partially coherent beam into multiple "quasi-Schell" modes and framing of the decomposition as a weight low-rank approximation problem.

# Chapter 2

# Illumination

We begin our investigation by establishing a formal definition of the illumination problem, where we wish to design a source system so that it generates a suitable illumination pattern to elicit a desired response out of a sink system. Given this definition, we can then apply several assumptions about the sink system that are valid for most application areas to show that any incoming illumination can be fully characterized by the second-order characteristics of a partially coherent beam. We will then conclude by reviewing the basics of partially coherent light and phase space before proceeding to an analysis of illumination system families in the next chapter.

## 2.1 Sources and sinks

In any illumination setup, there must be component(s) that generate light and a component(s) that receive the light. For this manuscript, we will refer to the light generating component(s) as a *source* system and the light receiving component(s) as a *sink* system. The source system consists of one or more light emitters, whose outgoing light can be modified by one or more occluders (which can be thought of as amplitude masks) and one or more lenses (which can be thought of as phase masks). The sink system will have a similar composition – incoming light is modulated by one or more occluders and/or lenses, and the modulated light is finally absorbed by one or more detectors. Let's define the *illumination* to be the optical radiation after it

Figure 2.1: We can formally define the illumination problem as one where we wish to design a source system such that it generates an illumination pattern that causes the sink system to produce a desired output response. A source system will in general consist of emitters and optical elements such as occluders and lenses. A sink system will in general consist of detectors and its own occluders and lenses.

finishes interacting with the optical elements in the source system and before it starts interacting with the optical elements in the sink system. Given these definitions, we can define the illumination problem as:

> Given a specified sink system and a desired response from
> its detectors, design a source system such that the exitant
> illumination from it drives the sink system in such a way
> that the desired response is obtained.

To solidify this definition, let us explore a few examples of illumination setups. One example setup would be a person viewing a three-dimensional display consisting of a liquid crystal display (LCD) panel with an attached lenticular sheet. In this setup, the source system is the three-dimensional display, the illumination is the light emanating from the display, and the sink system consists of the person's eyes. The backlight in the LCD panel contains the emitters in the source system, the LCD matrix itself is an occluder that modulates the color and intensity of the light, and

the lenticular sheet is a set of lenses that refract the outgoing light. In the sink system, the irises are occluders that limit the incoming angle of light, and the lenses of the eyes focus that light onto the retina, which consists of many detectors in the form of neurons. For this particular setup, the problem being solved is the generation of appropriate three-dimensional patterns of light from the display such that they create a specific response in the retinal neurons of the viewer, causing the viewer to see the desired three-dimensional image.

Another example setup would be the optical stimulation of neurons using a computer generated holographic device. The laser (emitter), spatial light modulator(s) (occluder and/or lens) and microscope optics (occluders and/or lenses) form the source system. The channelrhodopsin labeled neurons (detectors) in the live tissue sample form the sink system. The illumination would be the coherent light emitted from the objective of the microscope. For this particular setup, the problem being solved is the generation of appropriate three-dimensional patterns of light from the CGH device so that the illumination deposits sufficient energy on the labeled neurons, causing them to fire.

Recall that we summarized three different illumination device families in the previous chapter. Each device, in this formal definition, will be considered a source system, and the application of this device to a particular problem will specify the sink system. Thus, to establish how well these different devices perform, we will need to compare the types of illumination patterns each device or device family can achieve.

## 2.2 Quantifying illumination

To enable a rigorous comparison of illumination patterns, we require a rigorous mathematical description of the exitant illumination from a source system. One possibility is to aim for the most general of descriptions, an optical scalar field oscillating over space and time:

$$U(\mathbf{r}, t) \tag{2.1}$$

With the exclusion of polarization effects and evanescent fields, this description suffices to fully quantify the light being emitted by a source system. However, this representation is unwieldy – its size and format makes it hard to extract intuitive and useful descriptions of the light. Furthermore, it contains too much information, most of which would be lost by the detectors in the sink system.

To arrive at a more intuitive and compact representation of the illumination, let us make some simplifying assumptions which hold for the majority of use cases:

1. **The sink system consists of only linear components outside of the detectors themselves.**
   Most optical components in common everyday use are linear. This includes any passive occluders and/or lenses. The Kerr effect [27–29] and two-photon effects [30, 31] would be two effects which violate this assumption. Harmonic-generating nonlinear crystals would also violate this assumption.

2. **Components in the sink system have time-invariant behavior.**
   In general, optical components are static and do not vary in their behavior over time. If components vary slowly over time, then we can consider several separate time-invariant sink systems, one for each time segment.

3. **The detectors in the sink system can not detect the wave itself, but only detect the integral of incoming intensity of the light over a relatively long time compared to the oscillation period and any path-length differences induced by the sink system.**
   Only electromagnetic waves at optical frequencies fall within the scope of this manuscript, and it is rare for a device to react to the temporal oscillations of the EM wave itself at those frequencies. In most cases, optical detectors simply integrate the amount of energy it has received over a period of time. The integration time assumption should only be violated in either very large systems or in systems that have a lot of recursive reflections.

Given these assumptions, we can formulate an expression for the energy $Y_i$ received by the $i^{th}$ detector in the sink system as a function of the optical scalar field $U(\mathbf{r}, t)$

emanated from the source system:

$$Y_i = \int \left| \iint h_i(\mathbf{r}, \tau) U(\mathbf{r}, t - \tau) d\mathbf{r} d\tau \right|^2 dt \tag{2.2}$$

where $h_i(\mathbf{r}, t - \tau)$ is the contribution of the field at $\mathbf{r}$ at time $\tau$ to the $i^{th}$ detector at time $t$. We're modeling the energy as the magnitude squared of the incoming field to the detector and we can use a linear time-invariant expression because of the first two assumptions. The integral over $t$ is due to the final assumption.

We can expand the magnitude squared of a double integral term in Equation (2.2) into a quadruple integral:

$$
\begin{aligned}
Y_i &= \int \left| \iint h_i(\mathbf{r}, \tau) U(\mathbf{r}, t - \tau) d\mathbf{r} d\tau \right|^2 dt \\
&= \int \iint h_i(\mathbf{r}_1, \tau_1) U(\mathbf{r}_1, t - \tau_1) d\mathbf{r}_1 d\tau_1 \iint h_i^*(\mathbf{r}_2, \tau_2) U^*(\mathbf{r}_2, t - \tau_2) d\mathbf{r}_2 d\tau_2 dt \\
&= \iiiint\!\!\int h_i(\mathbf{r}_1, \tau_1) h_i^*(\mathbf{r}_2, \tau_2) U(\mathbf{r}_1, t - \tau_1) U^*(\mathbf{r}_2, t - \tau_2) d\mathbf{r}_1 d\mathbf{r}_2 d\tau_1 d\tau_2 dt
\end{aligned}
\tag{2.3}
$$

Then, by performing the following change of variables:

$$\bar{\tau} = \tau_1 \quad \tau' = \tau_2 - \tau_1 \tag{2.4}$$

we arrive at:

$$
\begin{aligned}
Y_i &= \iiiint\!\!\int h_i(\mathbf{r}_1, \bar{\tau}) h_i^*(\mathbf{r}_2, \bar{\tau} + \tau') U(\mathbf{r}_1, t - \bar{\tau}) U^*(\mathbf{r}_2, t - \bar{\tau} - \tau') d\mathbf{r}_1 d\mathbf{r}_2 d\bar{\tau} d\tau' dt \\
&= \iiiint h_i(\mathbf{r}_1, \bar{\tau}) h_i^*(\mathbf{r}_2, \bar{\tau} + \tau') \int U(\mathbf{r}_1, t - \bar{\tau}) U^*(\mathbf{r}_2, t - \bar{\tau} - \tau') dt d\mathbf{r}_1 d\mathbf{r}_2 d\bar{\tau} d\tau'
\end{aligned}
\tag{2.5}
$$

Applying our third assumption, we assume that the range of $\bar{\tau}$ over which $h(\mathbf{r}, \bar{\tau})$ is nonzero is very small compared to the limits of integration of $t$, and further simplify

the expression to:

$$\begin{aligned} Y_i &= \iiiint h_i(\mathbf{r}_1, \bar{\tau}) h_i^*(\mathbf{r}_2, \bar{\tau}+\tau') \Gamma_U(\mathbf{r}_1, \mathbf{r}_2, \tau') d\mathbf{r}_1 d\mathbf{r}_2 d\bar{\tau} d\tau' \\ &= \iiint \Gamma_h(\mathbf{r}_1, \mathbf{r}_2, -\tau') \Gamma_U(\mathbf{r}_1, \mathbf{r}_2, \tau') d\mathbf{r}_1 d\mathbf{r}_2 d\tau' \end{aligned} \qquad (2.6)$$

where $\Gamma_h$ is the cross correlation of the impulse response $h(\mathbf{r}, \tau)$ and $\Gamma_U$ is known in optics as the mutual coherence function of the stochastic scalar field $U(\mathbf{r}, t)$. That is, the energy received at any detector is only a function of some illumination-invariant impulse response and the mutual coherence function of the illumination. Describing an optical field by a mutual coherence field means we are viewing light as a *partially coherent field*. Let us now review what partial coherence means and use this opportunity to establish some notation for the subsequent chapters.

## 2.3   Partial coherence

The physical reason we need to study coherence is that light is a stochastic process – a single "monochromatic" point emitter does not typically emit a constant wave in a pure frequency. That is, point emitters do not generate photons at a constant rate; they create photons randomly. Therefore, if we consider the oscillation emanating from the point emitter over time, we'll see that it is locally monochromatic over short time intervals, but that there will be random phase shifts of the wave over time. What this does is to reduce the amount of interference resulting from mixing light from two different emitters.

As an example, let us consider the average intensity at a point receiving contributions from two different point emitters with vastly different phase shifts. Let $U_1(t)$ be the field contribution from the first emitter and $U_2(t)$ be the field contribution from

the second emitter. Then, the average intensity $I$ is a time average given by:

$$\begin{aligned}
I &= \left\langle |U_1(t) + U_2(t)|^2 \right\rangle \\
&= \left\langle U_1(t)U_2^*(t) + U_2(t)U_1^*(t) + |U_1(t)|^2 + |U_2(t)|^2 \right\rangle \\
&= 2\mathrm{Re}\left\langle U_1(t)U_2^*(t) \right\rangle + \left\langle |U_1(t)|^2 \right\rangle + \left\langle |U_2(t)|^2 \right\rangle \qquad (2.7)
\end{aligned}$$

However, random phase shifts cause $U_1(t)$ to be decorrelated from $U_2(t)$, and hence the first term is zero, yielding the result that the average intensity measured at this point is simply a sum of the *intensity* contributions from the two point emitters. In this case, we can say that the light from the two emitters are mutually *incoherent*.

Now, consider the case that $U_1(t)$ and $U_2(t)$ are actually correlated. For instance, $U_1(t)$ and $U_2(t)$ are measurements of the field at two different locations of a laser beam. Then, the first term in the expression in equation (2.7) will not be zero and will in fact vary based on the relative phase shift of $U_1(t)$ and $U_2(t)$. This perturbation is what causes interference, since the intensity of the sum of two fields is not the sum of the intensities of the two fields. We would call this situation the fully *coherent* case.

In general, in a "monochromatic" system with multiple point emitters, the light generated will tend to be *partially coherent*, where interference has a partial effect on the intensity. Suppose we wish to analyze what happens when we mix contributions from different parts of a partially coherent field with a possible time delay between components. Mixing the partially coherent field at location $\mathbf{r}_1$ with the partially coherent field at a time $\tau$ later at location $\mathbf{r}_2$ yields an average intensity $I_{sum}$ given by:

$$\begin{aligned}
I_{sum} &= \left\langle |U(\mathbf{r}_1, t) + U(\mathbf{r}_2, t+\tau)|^2 \right\rangle \\
&= \left\langle |U(\mathbf{r}_1, t)|^2 \right\rangle + \left\langle |U(\mathbf{r}_2, t+\tau)|^2 \right\rangle + 2\mathrm{Re}\left\langle U(\mathbf{r}_1, t)U^*(\mathbf{r}_2, t+\tau) \right\rangle \\
&= I(\mathbf{r}_1) + I(\mathbf{r}_2) + 2\mathrm{Re}\left\langle U(\mathbf{r}_1, t)U^*(\mathbf{r}_2, t+\tau) \right\rangle \qquad (2.8)
\end{aligned}$$

That is, the average intensity of the sum is the sum of the average intensities plus twice the real part of the cross correlation. This correlation is referred to in the optics

literature as the mutual coherence function [32]:

$$\Gamma(\mathbf{r}_1, \mathbf{r}_2, \tau) = \langle U(\mathbf{r}_1, t)U^*(\mathbf{r}_2, t + \tau) \rangle \tag{2.9}$$

This function describes the entire coherence state of the partially coherent field and can be used to derive the average intensity at any point. Spatial coherence is described by the value of the mutual coherence function when $\mathbf{r}_1 \neq \mathbf{r}_2$. Temporal coherence is described by the value of the mutual coherence function when $\tau \neq 0$. Furthermore, since the mutual coherence function is the expected value of the product of two fields, linear operations on coherent fields can be applied to partially coherent fields by simply applying it to the two components (corresponding to $\mathbf{r}_1$ and $\mathbf{r}_2$ separately. Therefore, the mutual coherence function is sufficient to describe most optical phenomena outside of quantum and nonlinear optics.

In 1982, Wolf proposed a new way of looking at this function [33]. Taking the Fourier transform with respect to $\tau$ of the mutual coherence function, yields what is known as the *cross spectral density*:

$$W(\mathbf{r}_1, \mathbf{r}_2, \omega) = \frac{1}{2\pi} \int \Gamma(\mathbf{r}_1, \mathbf{r}_2, \tau)e^{-j\omega\tau}d\tau \tag{2.10}$$

The cross spectral density transforms the partially coherent field into a set of different optical frequencies and describes the spatial coherence between points at each frequency independently. Furthermore, Wolf shows that the cross spectral density at any particular frequency $\omega$ can be factored into a sum of coherent fields $\sqrt{\lambda_n}\phi_n(\mathbf{r}, \omega)$, i.e. *coherent modes*:

$$W(\mathbf{r}_1, \mathbf{r}_2, \omega) = \sum_n \lambda_n(\omega)\phi_n^*(\mathbf{r}_1, \omega)\phi_n(\mathbf{r}_2, \omega) \tag{2.11}$$

We will revisit this result in subsequent chapters, as this is a powerful tool for both analyzing and synthesizing partially coherent fields.

For the simpler case of quasi-monochromatic light, when the bandwidth of light is very small compared to the average frequency $\omega_0$, we can take a slice of the cross

spectral density at $\omega = \omega_0$, yielding an expression with only spatial variables. This expression is a more rigorous definition of the *mutual intensity*, a function that describes the spatial coherence of the field:

$$J(\mathbf{r}_1, \mathbf{r}_2) = W(\mathbf{r}_1, \mathbf{r}_2, \omega_0) = \Gamma(\mathbf{r}_1, \mathbf{r}_2, 0) \tag{2.12}$$

The equality in the above equation holds due to degeneracy of the integral in Eq. (2.10) caused by quasi-monochromaticity, i.e. $\Gamma$ is a pure phasor in $\tau$.

Furthermore, for a quasi-monochromatic partially coherent beam in the paraxial regime, we need only the mutual intensity function at a single transverse plane to obtain the mutual intensity function (and thus optical qualities of the beam) at other planes. This is because propagation of light is a linear operation and thus we can apply a propagation kernel such as the Fresnel diffraction integral to the two halves of the cross correlation at one plane to get the cross correlation at a different plane:

$$J_{z_1}(x_1, y_1, x_2, y_2) = \frac{1}{\lambda^2(z_1 - z_0)^2} \times$$
$$\iiiint J_{z_0}(\xi_1, \eta_1, \xi_2, \eta_2) e^{\frac{jk}{2(z_1 - z_0)}\left[(x_1 - \xi_1)^2 + (y_1 - \eta_1)^2 + (x_2 - \xi_2)^2 + (y_2 - \eta_2)^2\right]} d\xi_1 d\eta_1 d\xi_2 d\eta_2 \tag{2.13}$$

where $k = 2\pi/\lambda$ and $\lambda$ is the wavelength. Therefore, the mutual intensity in a three-dimensional volume of free propagation can be fully determined by the mutual intensity at a particular plane.

## 2.3.1 Discretization

Since optical systems are inherently bandlimited, whether by wavelength or by the systems' numerical aperture, the optical field and coherence measures such as the mutual intensity can be fully described by a discretized representation according to the Nyquist-Shannon sampling theorem. Given a sufficiently small sampling interval $\Delta$ and a bandlimited two-dimensional field $U(x, y)$, the *discrete form* $U[m, n]$ of the

field is the following two-dimensional discrete function:

$$U[m, n] = U(m\Delta, n\Delta) \tag{2.14}$$

where $m$ and $n$ are integers. Likewise, the *discrete form* $J[m_1, n_1, m_2, n_2]$ of a bandlimited mutual intensity $J(x_1, y_1, x_2, y_2)$ is the following four-dimensional discrete function:

$$J[m_1, n_1, m_2, n_2] = J(m_1\Delta, n_1\Delta, m_2\Delta, n_2\Delta) \tag{2.15}$$

Since most optical systems are linear, it would be useful to be able to operate on fields using linear algebra. Assuming a field $U[m, n]$ of finite extent where $m, n \in \{1, \ldots, N\}$, let us form a *vector form* $\mathbf{u}$ of the field. First, let us write the discrete form of the field as a matrix $U$:

$$U = \begin{pmatrix} U[1,1] & U[1,2] & \ldots & U[1,N] \\ U[2,1] & U[2,2] & \ldots & U[2,N] \\ \vdots & \vdots & \ddots & \vdots \\ U[N,1] & U[N,2] & \ldots & U[N,N] \end{pmatrix} = \begin{pmatrix} \mathbf{u}_1 & \mathbf{u}_2 & \ldots & \mathbf{u}_N \end{pmatrix} \tag{2.16}$$

where $\mathbf{u}_n$ are the columns of $U$. Then, the vector form of the field is created by the vertical concatenation of these columns:

$$\mathbf{u} = \begin{pmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \\ \vdots \\ \mathbf{u}_N \end{pmatrix} \tag{2.17}$$

Note that each row in $\mathbf{u}$ corresponds to a specific spatial location. For each row index $i$, let us define the corresponding spatial locations by two discrete functions:

$$m[i] = [(i-1) \mod N] + 1$$
$$n[i] = \lfloor (i-1)/N \rfloor + 1$$

The process of converting a discrete representation into vector form is illustrated pictorially in Figure 2.2.



Figure 2.2: A bandlimited continuous field can be sampled and converted into a discrete form $U[m,n]$. A bounded discrete form $U[m,n]$ can be converted into vector form by vertically stacking the columns in $U[m,n]$.

Since we have converted the field from two-dimensional discrete function down to a one-dimensional vector representation, we should be able to convert the mutual intensity down from a four-dimensional discrete function into a two-dimensional matrix representation. We will now define the *matrix form* $J \in \mathbb{C}^{N^2 \times N^2}$ of the mutual intensity to be:

$$
\begin{pmatrix}
J[m[1],n[1],m[1],n[1]] & J[m[1],n[1],m[2],n[2]] & \cdots & J[m[1],n[1],m[N^2],n[N^2]] \\
J[m[2],n[2],m[1],n[1]] & J[m[2],n[2],m[2],n[2]] & \cdots & J[m[2],n[2],m[N^2],n[N^2]] \\
\vdots & \vdots & \ddots & \vdots \\
J[m[N^2],n[N^2],m[1],n[1]] & J[m[N^2],n[N^2],m[2],n[2]] & \cdots & J[m[N^2],n[N^2],m[N^2],n[N^2]]
\end{pmatrix}
\tag{2.18}
$$

Note that $J$ is a Hessian matrix. Given this matrix representation, Wolf's coherence mode decomposition can be described by a singular value decomposition:

$$
J = F\Sigma F^H
\tag{2.19}
$$

where $\Sigma$ is a diagonal matrix consisting of the singular values $\sigma_i$ and $F \in \mathbb{C}^{N^2 \times N^2}$ is an orthonormal matrix with columns $\mathbf{f}_i$:

$$F = \begin{pmatrix} \mathbf{f}_1 & \mathbf{f}_2 & \dots & \mathbf{f}_{N^2} \end{pmatrix} \tag{2.20}$$

The matrix/vector form equivalent of the coherent mode decomposition in Equation (2.11) is given by:

$$J = \sum_{i=1}^{N^2} \sigma_i \mathbf{f}_i \mathbf{f}_i^H \tag{2.21}$$

The matrix form of the mutual intensity thus allows us to easily see how coherent the partially coherent field is. A fully coherent field would only consist of one mode and hence would be a rank-one matrix. A fully incoherent field has no nonzero terms off the diagonal and therefore would have up to $N^2$ modes. The matrix representation of coherence was initially introduced by Gamo [34] and properties of the matrix representation (e.g. rank vs. coherence) are explored in Ozaktas et al.'s work [35].

## 2.4   Phase space

The idea of partial coherence has also been explored in the field of phase space optics, a topic that has acquired recent popular interest. Of importance are two quadratic space-frequency functions [36] computed via a Fourier transform and coordinate system change of the mutual intensity function on a plane [37]. They are the Wigner distribution function,

$$B(x, y, f_\xi, f_\eta) = \iint J\left(x + \tfrac{\xi}{2}, y + \tfrac{\eta}{2}, x - \tfrac{\xi}{2}, y - \tfrac{\eta}{2}\right) e^{-j2\pi\left(f_\xi \xi + f_\eta \eta\right)} d\xi d\eta \tag{2.22}$$

and its Fourier conjugate (including a coordinate reversal), the ambiguity function,

$$A(f_x, f_y, \xi, \eta) = \iint J\left(x + \tfrac{\xi}{2}, y + \tfrac{\eta}{2}, x - \tfrac{\xi}{2}, y - \tfrac{\eta}{2}\right) e^{-j2\pi(f_x x + f_y y)} dx dy \tag{2.23}$$

The Wigner distribution function is a function of space and frequency variables. Since pure spatial frequencies of the optical field on a plane correspond to plane

waves along different directions (i.e. the angular spectrum), we can intuitively treat the frequency variables of the Wigner distribution function as directions. In fact, one way to interpret the Wigner distribution function is that it is a quasi-probability distribution of the position and momentum of photons. Since position (intersection point with the plane in question) and direction (momentum) is all that is needed to specify a ray of light, the Wigner distribution function can be thought of as a more rigorous ray-model representation of light that incorporates interference and diffraction phenomena. The use of the Wigner distribution function as a measure of "generalized radiance" of rays was the reason it was first introduced to the field of optics [38–42]. Unlike real radiance, "generalized" radiance can be negative, which creates interference effects.

The Wigner distribution is an intuitive mathematical representation of light that reduces complicated direct operations on the field to simpler operations – propagation in a paraxial beam from one transverse plane to another can be computed via a linear shear in the coordinate system, and lensing by ideal thin lenses can also be represented by shears in an orthogonal direction. The actual intensity of the field at the plane can be calculated by projecting along the frequency axis, and the angular distribution of light can be calculated by projecting along the spatial axis. A summary of these methods can be found in Bastiaan's paper on the subject [43].

Likewise, the ambiguity function also exhibits useful properties, the most important of which is that the Fourier transform of the intensity at different planes in a paraxial beam are simply slices of the ambiguity function [44,45]. This property forms the backbone of the development of wavefront coding systems for extended depth of field [46]. In the full four-dimensional case of the ambiguity function in 3D space, each slice through the origin corresponds to the 2D Fourier transform of an image formed through an astigmatic lens system; the slopes correspond to the focal lengths along the two transverse axes.

Lastly, since the Wigner distribution and ambiguity function are derived directly and reversibly from the mutual intensity, they are equivalent representations of the coherence state of a partially coherent field.

## 2.5   Summary

The illumination problem can be formulated as the design of a source system that can produce a specific illumination pattern so that it elicits a desired response in the sink system. If we make some assumptions which correspond to the most common situations, we can show that detector responses in the sink system are simply a function of the mutual coherence function (mutual intensity in the quasi-monochromatic case) of the illumination generated by the source system. The mutual intensity can also be viewed in three other equivalent representations:

1. as an incoherent sum of coherent modes – a partially coherent field can be described as multiple fully coherent fields which do not interfere with each other,

2. as a Wigner distribution – a partially coherent field can be modeled as a set of generalized rays with possibly negative radiance,

3. as an ambiguity function – a partially coherent field can be modeled as a set of 2D intensity patterns formed by imaging through an (astigmatic) lens with differing focal lengths along the two transverse axes.

For ease of notation and analysis, we will from this point onward only consider quasi-monochromatic light, as extension to polychromatic light is straightforward from the independence of different frequencies in the cross spectral density. Except where noted, we will also consider only paraxial beams; extension to non-paraxial scalar fields can be done using the Rayleigh-Sommerfeld diffraction integral or angular spectrum methods [47].

# Chapter 3

# Analysis

We will now use the ideas of partial coherence discussed in the previous chapter to analyze the capabilities of the following three illumination device families:

1. **Ray devices**
   The desired illumination is created by generating a set of light rays with varying radiance. This approach models illumination by specifying a radiance for each ray in free space.

2. **Holographic devices**
   The desired illumination is generated by modulating a coherent field (e.g. from a laser) using spatial light modulator(s). This approach models illumination as a fully coherent field.

3. **Volumetric devices**
   The desired illumination is generated by scanning planar images or points through space, either through direct excitation of a medium or by imaging an incoherent light source. This approach models illumination as the output of a set of isotropic (diffuse) point emitters in free space.

We will use coherence representations such as the mutual intensity or the ambiguity function to evaluate what illumination patterns can and cannot be generated by each device family. We will show that there are drawbacks for each device family because the underlying model for that family fails to describe certain types of illumination.

## 3.1 Ray devices

Ray devices rely on the geometric optics model of light travelling along rays. Each ray, if unhindered, contains a constant amount of radiance (i.e. optical energy) travelling along it, and individual rays don't interact with each other. The effect of changes in index of refraction (e.g. in the case of lenses) is modeled by the bending of rays through either refraction or reflection, computed using Snell's Law and Fresnel equations. The intensity (power density) of light at a point can be computed by integrating the radiance of all the rays that intersect that point. This approach forms the basis of plenoptic cameras [48, 49], light field rendering [14] and integral imaging [16]. Commercially available multiscopic displays using parallax barriers or lenticular sheets are also based on this approach.

Under geometric optics, we can characterize illumination reaching a viewer by the set of rays that intersect a virtual plane placed directly in front of the viewer. Doing so results in a four-dimensional distribution of (positive) radiance in "ray space", where two of the dimensions are in position and two of the dimensions are in direction. This concept is known as the "light field" in the computer graphics literature [14]. For the rest of manuscript, when we refer to the light field, we will be referring to this concept, as opposed to the concept of an optical field in general. This light field has the same dimensionality as the partial coherence representations covered in the previous chapter. Even though one of these representations, the Wigner distribution, can be thought of as a quasi-distribution on a photon's position-momentum state space and thus mimicking the layout of ray space, quantum mechanics' uncertainty principle stipulates that we cannot know simultaneously with exact precision the position and momentum of a particle. Therefore, any approach that utilizes ray space will be limited by this issue, as we will see when we analyze a few devices in this device family using partial coherence concepts.

### 3.1.1 Integral imaging display

In an integral imaging display, an incoherent 2D image is placed at the front focal plane of a lens array, as shown in Figure 3.1. This lens array demultiplexes the

spatioangular information on this image into a light field. In the ideal case, each pixel generates light that reaches only one lens, and this light is then converted into a "fat" ray emanating from the other side of the lens. Integral imaging displays usually operate in the macroscopic regime, but they have been adapted for use in microscopy as well to create three-dimensional patterns inside a microscopic specimen [2]. Furthermore, instead of a lens array, a pinhole array can also be used, which results in a parallax barrier display.



Figure 3.1: An integral imaging display uses a lens array (ii) to convert the spatioangular data found on an incoherent image plane (i) into an output light field (iii). The lens array has focal length $F$ and is located at $z = 0$. In this figure, two pixels are turned on at the incoherent image plane, which generates two "rays" from this device.

Since the image plane in front of the lens array is incoherent, each "ray" generated by this device from a pixel on this image plane is also incoherent with respect to any other ray. Therefore, we can analyze this using the Wigner distribution and consider each pixel separately. The total result would simply be the sum of all the pixel contributions. To simplify notation, we will only consider the flat-land representation, where we have one-dimensional images and lens arrays. The extension to two dimensions is straightforward. Recall from the previous chapter that we will also assume a quasi-monochromatic system.

Now let's calculate the "impulse response" in the Wigner distribution induced by light from a single pixel. Suppose the pixel in question is at position $x = s_0 + t_0 F$ with

unit intensity and its corresponding lens in the lens array is centered at $x = s_0$ with aperture function $a(x)$. We can think of this pixel as generating the ray at position $s_0$ and angle (slope) $t_0$. We will calculate the impulse response by calculating the following four Wigner distributions in sequence:

1. $B_1$ at the original pixel plane

2. $B_2$ after propagating to just before the lens array plane

3. $B_3$ after applying the focusing (phase) component of the lens

4. $B_4$ after applying the masking (amplitude) component of the lens

At $z = -F$, the Wigner distribution of this pixel is given by:

$$B_1(x, f_\xi; s_0, t_0) = \delta(x - s_0 - t_0 F) \tag{3.1}$$

A propagation of distance $f$ along the axis from the incoherent image plane to the lens array plane shears the Wigner distribution:

$$
\begin{aligned}
B_2(x, f_\xi; s_0, t_0) &= B_1(x - f_\xi \lambda F, f_\xi; s_0, t_0) \\
&= \delta(x - f_\xi \lambda F - s_0 - t_0 F) \tag{3.2}
\end{aligned}
$$

Focusing by the lens causes a shear in the orthogonal direction:

$$
\begin{aligned}
B_3(x, f_\xi; s_0, t_0) &= B_2(x, f_\xi + (x - s_0)/(\lambda F); s_0, t_0) \\
&= \delta(x - \lambda F(f_\xi + (x - s_0)/(\lambda F)) - s_0 - t_0 F) \\
&= \delta(\lambda F f_\xi - t_0 F) \\
&= \delta(f_\xi - t_0/\lambda)/(\lambda F) \tag{3.3}
\end{aligned}
$$

Masking by the aperture function results in a convolution along the angular coordinate for each spatial coordinate with the Wigner distribution $B_a(x - s_0, f_\xi)$ of the aperture

function $a(x - s_0)$ [36]:

$$
\begin{aligned}
B_4(x, f_\xi; s_0, t_0) &= \int B_3(x, f'_\xi; s_0, t_0) B_a(x - s_0, f_\xi - f'_\xi) df'_\xi \\
&= \int \delta(f'_\xi - t_0/\lambda)/(\lambda F) B_a(x - s_0, f_\xi - f'_\xi) df'_\xi \\
&= B_a(x - s_0, f_\xi - t_0/\lambda)/(\lambda F)
\end{aligned}
\tag{3.4}
$$

Therefore, the resulting system is shift invariant, and the impulse response is simply the Wigner distribution of the aperture function. Given an input idealized light field $L(s, t)$, the resulting output Wigner distribution can be written as:

$$
B_{out}(x, f_\xi) = \frac{1}{\lambda F} \iint B_a(x - s, f_\xi - t/\lambda) L(s, t) ds dt
\tag{3.5}
$$

That is, the Wigner distribution of the aperture *blurs* the ideal light field we would wish to render. In practice, $L(s, t)$ isn't a continuous function with infinite support; it is discrete in $s$ and bounded in both dimensions. However, $L(s, t)$ in that case can be represented as a sum of delta functions, and the above equation can still be applied. Coincidentally, *light field imaging* using a microlens array (essentially the integral imaging display in reverse) also results in a convolution of the Wigner distribution microlens aperture function with the incoming ideal Wigner distribution [25], due to Helmholtz reciprocity.

A more intuitive picture can be obtained by investigating this blur in the Fourier domain. That is, we can look at the ambiguity function $A_{out}(f_x, \xi)$ of the output light as a function of the Fourier transform $\tilde{L}(f_s, f_t)$ of the ideal light field. The convolution theorem implies that we are masking the Fourier transform of the ideal light field by the ambiguity function of the aperture of the lens:

$$
A_{out}(f_x, \xi) = \frac{1}{F} A_a(f_x, \xi) \tilde{L}(f_x, -\xi/\lambda)
\tag{3.6}
$$

That is, areas where $A_a(f_x, \xi)$ has a small value masks out the corresponding areas in $\tilde{L}(f_x, -\xi/\lambda)$. In other words, *high frequency regions of the ideal light field are attenuated and/or removed in an integral imaging display.* This effect is physically caused

by diffraction, as a finite sized lens can only image at limited resolution. The front focal plane of the lens is conjugate with the far field (i.e. angular information) of the light. Therefore, reducing the lens aperture to obtain spatial resolution at the lens array plane results in a reduction in angular resolution due to the diffraction limit. The extension to parallax barriers, where pinholes replace microlenses, is straightforward, since we can describe a pinhole by an aperture function of smaller spatial support compared to that of a lens.

To a first-order approximation, an ambiguity function occupies roughly the same area regardless of what the original function is. We can see this by projecting the ambiguity function along the $f_x$ and $\xi$ coordinates to obtain the magnitude squared of the original function and its Fourier transform, respectively. Since increasing the extent of a function reduces its extent in the Fourier domain and vice versa, attempting to increase the area by increasing the width of the ambiguity function would reduce its height, and so forth. In fact, if the width of the original function is scaled by a factor $\alpha$, the width of the Fourier transform would have to be scaled by a factor $1/\alpha$. The resulting extent of the ambiguity function would still be the same. Thus, masking in the integral imaging display tells us that we will only be able to access a fixed area of the Fourier transform of the ideal light field around the origin.

For the specific case of the light field illuminator in [2], the fixed area limitation in the ambiguity function results in a trade-off between transverse (x-y) resolution of a focused spot and the number of distinctly addressable depths (z). The transverse resolution is directly controlled by the microlens pitch and the objective magnification. However, reducing the microlens pitch to increase transverse resolution causes the angular resolution to drop, which in turn reduces the number of addressable depths, as shown in Figure 3.2. This is because a specific angular extent (i.e. NA) induces a vertical blur in the ambiguity function, causing a perfectly focused spot to become a thick line. In order for focused spots to be distinct, they must not overlap by much in the ambiguity function. However, increasing the transverse resolution reduces the "height" of the ambiguity function mask, reducing the number of these depths that will fit inside the new mask. This also reduces the maximum depth as well. Note that these effects are caused by physical limitations and cannot simply be remedied

Figure 3.2: Pictorial representations of the ambiguity function corresponding to various focused spots with different microlens configurations. A focused spot at a specific depth created by an integral imaging display with limited NA forms a slab in the ambiguity function, whose borders are indicated by the bold lines in (i). For the particular configuration in (i), four distinct depths are achievable, with slabs corresponding to each additional depth outlined in thin lines. Increasing the transverse resolution (increasing the width of the ambiguity function mask) by a factor of two causes the angular resolution to drop (a decrease in the height of the ambiguity function mask), resulting in half as many distinct depths, as shown in (ii) Note that the maximum depth (slope) has been quartered as well.

via decreasing the pixel size on the illumination device.

**A resolution anomaly**

Before we move on, let us take a closer look at the ambiguity function corresponding to a plain rectangular aperture, as shown in Figure 3.3. This would be the "transfer" function for a microlens system in flatland with no apodization. Recall that the ambiguity function consists of slices that correspond to the Fourier transform at different depths. It is interesting to note here that for a particular depth off the primal plane, shown by the dotted line in the figure, this aperture yields higher-than-expected resolution.



Figure 3.3: Use of a plain (unapodized) rectangular aperture yields a mask in the ambiguity function as shown in (i). Positive values are red and negative values are blue. Performing a coordinate transform converts the ambiguity function into a depth varying optical transfer function (OTF), as shown in (ii). The OTF is the Fourier transform of the impulse response in intensity and the vertical axis in (ii) denotes distance from the lens array plane ($z = 0$). Note that there is a slight increase in the width of the OTF past the lens array plane, indicated by the horizontal dotted line. The equivalent cut in the ambiguity function is shown as the slanted dotted line in (i).

This is somewhat counter-intuitive, as we expect diffraction to cause a bundle of

"parallel" rays emanating from a lens to always be gradually diverging. However, the diffraction pattern in the near field of a small aperture does show a "focal" spot shortly after the aperture, as shown in Figure 3.4. Recall that the "transfer" function is simply the ambiguity function of the aperture shape. This "focal" spot is characteristic of small apertures and is related to the concept of focal shift [50–52].



Figure 3.4: The diffraction pattern from a plane wave moving from the left to right hitting a square aperture of width $100\lambda$ at the left side of the figure, computed using the Fresnel diffraction integral. While diffraction from an aperture is sharp at the aperture plane, it also has a moderately sharp spot shortly after the aperture, labeled by the dotted line. This corresponds to the same dotted line in Figure 3.3. The diffraction pattern is not to scale; an equal distance in reality should be ten times longer along the horizontal/longitudinal/z axis than the vertical/transverse/x axis.

However, recall that $L(s,t)$ was also discrete in $s$. This implies that its Fourier transform $\tilde{L}(f_s, f_t)$ must also be periodic in $f_s$. Therefore, we would observe aliased copies at intervals of $1/\Delta_s$ where $\Delta_s$ is the sampling interval along $s$. Thus, for a standard integral imaging display, it is possible that aliasing would negate the extra resolution gain at that depth. One solution to exploit this extra resolution gain would be to have $s$ sampled more finely by temporally multiplexing the incoherent image in the integral imaging display and concurrently shifting the lens array through different positions. For example, to increase sampling in $s$ by two in each direction, we would iterate through four images, each corresponding to a different sub-microlens shifted set of spatial samples in the ideal light field we wish to generate.

Although stepping through different sub-microlens spatial offsets to achieve more resolution sounds like super-resolution, it is not a super-resolution method. This stepping is simply aiming to restore a small amount of resolution at planes near a specific depth, made possible by an oddity in the wave propagation of light through a small aperture. Through Helmholtz reciprocity, this method may be applied to imaging as well, but this differs from typical light field super-resolution techniques [53, 54] since this technique always offers a small net gain in resolution at the cost of temporal resolution and does not rely on deconvolution or priors to achieve this resolution gain. Note that since only a small amount of transverse resolution at a specific depth is restored, intentionally reducing the temporal resolution of the display might be too costly for the benefit gained.

### 3.1.2   Aperture scanning displays

A different approach to ray-based illumination generation scans a pinhole in the aperture (Fourier) plane of a time-varying incoherent image display so that each image creates a roughly parallel bundle of rays with direction determined by the pinhole position, and time-multiplexing the images results in filling ray space with the desired light field, as shown in Figure 3.5 [15]. The device can also be thought of as a fast time-varying display imaged through a telecentric lens system whose telecentric stop is pinhole sized and translated in sync with the display.

We will now apply the same phase space optics principles to this system as we did to the integral imaging display. Let $L(x, u/F)$, a function mapping a position $(x)$ on the image plane and a position $(u)$ on the aperture plane to an intensity value, with $F$ being the focal length of the 4-f system and $a(u)$ be the transparency profile of the pinhole aperture. This function is essentially an ideal light field representation. The resulting output Wigner distribution of the display can be calculated by accumulating across all pinhole positions the Wigner distribution of light produced by a single image-pinhole pair. The output Wigner distribution for a single image-pinhole pair is calculated via these intermediate Wigner distribution functions:

1. $B_1$ at the original pixel plane

Figure 3.5: In an aperture scanning display, an incoherent image plane (i) is imaged through a 4-f system (ii) with focal length $F$ onto an output plane (iii). Many images are shown in rapid succession, and a pinhole mask at the aperture plane (iv) of the 4-f system moves in tandem to select the direction of the parallel ray bundle emanating from the output plane.

2. $B_2$ after propagating through the first 2-f Fourier transforming lens system

3. $B_3$ after applying the pinhole mask

4. $B_4$ after propagating through the second 2-f Fourier transforming lens system

At the original pixel plane, we have an incoherent image with intensity distribution $I(x; u) = L(x, u/F)$, so the corresponding Wigner distribution is this one-dimensional function:

$$B_1(x, f_\xi; u) = L(x, u/F) \tag{3.7}$$

A Fourier transforming lens system performs the following coordinate transform corresponding to a ninety degree rotation:

$$\begin{aligned} B_2(x, f_\xi; u) &= B_1(-\lambda F f_\xi, x/(\lambda F); u) \\ &= L(-\lambda F f_\xi, u/F) \end{aligned} \tag{3.8}$$

Masking by the pinhole aperture induces a convolution along the spatial coordinate:

$$\begin{aligned} B_3(x, f_\xi; u) &= \int B_2(x, f_\xi - f'_\xi, ; u) B_a(x - u, f'_\xi) df'_\xi \\ &= \int L(-\lambda F f'_\xi, u/F) B_a(x - u, f_\xi - f'_\xi) df'_\xi \end{aligned} \tag{3.9}$$

Propagation through the second lens system also performs a Fourier transform and yields a coordinate transformation in the Wigner distribution:

$$\begin{aligned} B_4(x, f_\xi; u) &= B_3(-\lambda F f_\xi, x/(\lambda F); u) \\ &= \int L(-\lambda F f'_\xi, u/F) B_a(-\lambda F f_\xi - u, x/(\lambda F) - f'_\xi) df'_\xi \end{aligned} \tag{3.10}$$

Combining the contributions for each pinhole location, we arrive at the following result for the final output Wigner distribution at the output plane:

$$B_{out}(x, f_\xi) = 1/F \iint L(-\lambda F f'_\xi, u/F) B_a(-\lambda F f_\xi - u, x/(\lambda F) - f'_\xi) df'_\xi du \tag{3.11}$$

The corresponding ambiguity function is:

$$A_{out}(f_x, \xi) = 1/(\lambda F) A_a(-\xi/(\lambda F), -\lambda F f_x) \tilde{L}(-f_x, \xi/\lambda)  \tag{3.12}$$

These results show that this display system also blurs the ideal light field by a Wigner distribution, this time of the aperture function corresponding to the pinhole, and conclusions similar to the integral imaging display can be drawn.

### 3.1.3 Summary

As can be seen from the two derivations, ray-based display devices suffer from an inherent lack of resolution caused by a blur in the Wigner distribution or equivalently an apodization in the ambiguity function. For the specific instance of the light field illuminator for microscopy, this results in a trade-off between transverse resolution and number of individually addressable depths. This effect is precisely due to attempting to characterize outgoing illumination by a set of non-interfering rays, and thus any ray-based display system would suffer from this effect. Although this effect is minimal at the macroscopic scale, where three-dimensional patterns have feature sizes much larger than the diffraction limit of light, the effect is severe at microscopic scales close to the diffraction limit.

## 3.2 Holographic devices

Unlike ray-based devices, holographic devices generate very high resolution illumination patterns. Holography directly takes into account the wave nature of light and models light propagation in space as a wave equation as opposed to a set of rays. Holographic devices modulate an input coherent beam with a spatial pattern in amplitude and/or phase. This spatial pattern, or hologram, causes diffraction and interference patterns which generate the desired three-dimensional light pattern.

The modulation of a coherent beam by a fixed pattern will result in a fully coherent field and this coherence introduces limitations on the types of three-dimensional

patterns that can be generated. This is in addition to the standard limitations applicable to any illumination device, such as the diffraction limit (light patterns cannot be infinitely sharp since the wave function is bandlimited) and conservation of energy (a propagating light beam in free space cannot create a bright transverse plane of light followed by a fully dark transverse plane of light, as this would imply optical energy has been destroyed).

Constraints on possible patterns that can be generated by coherent fields is often described as one of dimensionality, since it is well-known that a coherent field in a volume can be fully described by a two-dimensional manifold in its Fourier transform. This leads naturally to the observation it would be very difficult to create arbitrary higher-dimensional patterns, such as a three-dimensional pattern in a volume. We will now discuss this limitation in more detail as well as show that coherence causes limitations even when the dimensionality matches, such as in the case of attempting to generate a planar intensity pattern.

### 3.2.1   Four-dimensional limitations

We begin our discussion by noting that a fully coherent field cannot by definition represent arbitrary partially coherent fields. Mathematically, there is a dimensionality mismatch between the four dimensions that is needed to describe arbitrary optical fields in any state of coherence and the two dimensions that can represent all fully coherent fields. Full coherence also forces the matrix form of the mutual intensity (as discussed in the previous chapter) to be rank-one. This is a degenerate case of coherent modes, where the mutual intensity only contains a single coherent mode and is thus fully coherent.

We will now reinforce this theory with a concrete example where slightly modifying the Wigner distribution of a well-known fully coherent field causes it to no longer be coherent, even though the resulting partially coherent field is still physically plausible. Let $B$ be the Wigner distribution function corresponding to a Gaussian beam:

$$B(x, y, f_\xi, f_\eta) = 4\sigma^2\pi e^{-4\sigma^2\pi^2 f_\xi^2 - 4\sigma^2\pi^2 f_\eta^2} e^{-x^2/\sigma^2 - y^2/\sigma^2} \tag{3.13}$$

The associated scalar field at $z = 0$ is a Gaussian with standard deviation $\sigma$ and the intensity profile at that plane has standard deviation $\sigma/2$:

$$U(x, y, z = 0) = e^{(-x^2 - y^2)/(2\sigma^2)} \;, \; I(x, y, z = 0) = e^{(-x^2 - y^2)/(\sigma^2)} \tag{3.14}$$

Let $\sigma$ be sufficiently large compared to the wavelength so that non-paraxial effects can be ignored. Now, consider expanding the Wigner distribution of this beam by a factor of $\kappa > 1$ along both frequency axes while maintaining the same intensity profile at $z = 0$:

$$\hat{B}(x, y, f_\xi, f_\eta) = \frac{1}{\kappa^2} B(x, y, f_\xi/\kappa, f_\eta/\kappa) \tag{3.15}$$

One effect of this angular expansion is compression of the intensity profile of the beam longitudinally by a factor of $\kappa$:

$$\hat{I}(x, y, z) = I(x, y, \kappa z) \tag{3.16}$$

We will now show that this modified Wigner distribution cannot be a coherence representation of a fully coherent field by analyzing the corresponding mutual intensity functions. The mutual intensity function corresponding to the original Gaussian beam is:

$$
\begin{aligned}
J(x_1, y_1, x_2, y_2) &= U(x_1, y_1, z = 0)U^*(x_2, y_2, z = 0) \\
&= e^{(-x_1^2 - x_2^2 - y_1^2 - y_2^2)/(2\sigma^2)} \\
&= e^{-\left[(x_1+x_2)^2 + (x_1-x_2)^2 + (y_1+y_2)^2 + (y_1-y_2)^2\right]/(4\sigma^2)} \tag{3.17}
\end{aligned}
$$

The stretch along the frequency axes of the Wigner distribution causes a compression along the spatial distance expressions in the mutual intensity:

$$\hat{J}(x_1, y_1, x_2, y_2) = e^{-\left[(x_1+x_2)^2 + \kappa^2(x_1-x_2)^2 + (y_1+y_2)^2 + \kappa^2(y_1-y_2)^2\right]/(4\sigma^2)} \tag{3.18}$$

If this new mutual intensity function were to represent that of a coherent field, then

it must be a separable function:

$$\hat{J}(x_1, y_1, x_2, y_2) = \hat{U}(x_1, y_1)\hat{U}^*(x_2, y_2) \tag{3.19}$$

This implies that:

$$\hat{J}(\hat{x}_1, \hat{y}_1, \hat{x}_2, \hat{y}_2)\hat{J}(\hat{x}_2, \hat{y}_2, \hat{x}_1, \hat{y}_1) = \hat{J}(\hat{x}_1, \hat{y}_1, \hat{x}_1, \hat{y}_1)\hat{J}(\hat{x}_2, \hat{y}_2, \hat{x}_2, \hat{y}_2) \tag{3.20}$$

for any coordinate variables $\hat{x}_1, \hat{y}_1, \hat{x}_2, \hat{y}_2$ such that they do not refer to the same spatial position:

$$(\hat{x}_1 - \hat{x}_2)^2 + (\hat{y}_1 - \hat{y}_2)^2 > 0 \tag{3.21}$$

Substituting (3.18) into the left hand side of (3.20) yields:

$$\hat{J}(\hat{x}_1, \hat{y}_1, \hat{x}_2, \hat{y}_2)\hat{J}(\hat{x}_2, \hat{y}_2, \hat{x}_1, \hat{y}_1) = e^{-2\left[(\hat{x}_1+\hat{x}_2)^2 + \kappa^2(\hat{x}_1-\hat{x}_2)^2 + (\hat{y}_1+\hat{y}_2)^2 + \kappa^2(\hat{y}_1-\hat{y}_2)^2\right]/(4\sigma^2)}$$
$$\tag{3.22}$$

Substituting (3.18) into the right hand side of (3.20) yields:

$$\hat{J}(\hat{x}_1, \hat{y}_1, \hat{x}_1, \hat{y}_1)\hat{J}(\hat{x}_2, \hat{y}_2, \hat{x}_2, \hat{y}_2) = e^{-\left[4\hat{x}_1^2 + 4\hat{x}_2^2 + 4\hat{y}_1^2 + 4\hat{y}_2^2\right]/(4\sigma^2)} \tag{3.23}$$

In order for the left and right hand sides to be equal:

$$
\begin{aligned}
(\hat{x}_1 + \hat{x}_2)^2 + \kappa^2(\hat{x}_1 - \hat{x}_2)^2 + (\hat{y}_1 + \hat{y}_2)^2 + \kappa^2(\hat{y}_1 - \hat{y}_2)^2 &= 2\hat{x}_1^2 + 2\hat{x}_2^2 + 2\hat{y}_1^2 + 2\hat{y}_2^2 \\
\kappa^2(\hat{x}_1 - \hat{x}_2)^2 + \kappa^2(\hat{y}_1 - \hat{y}_2)^2 &= (\hat{x}_1 - \hat{x}_2)^2 + (\hat{y}_1 - \hat{y}_2)^2 \\
\kappa^2 &= 1 \tag{3.24}
\end{aligned}
$$

Since we've defined $\kappa > 1$, the two sides in Equation (3.24) cannot be made equal and therefore this modified Wigner distribution cannot be the Wigner distribution of a fully coherent field.$\square$

However, if $\kappa$ is small enough such that the expansion in the angular spectrum doesn't exceed the diffraction limit, then the Wigner distribution function in Equation (3.15) does represent a physically valid *partially* coherent beam. Since the convolution of two Gaussians results in a Gaussian with a larger standard deviation, $\hat{B}$ can be

written as a convolution between the original Gaussian signal (with variance $(8\sigma^2\pi^2)^{-1}$ along the frequency axes) and a Gaussian signal in frequency with variance $(\kappa^2 - 1)/(8\sigma^2\pi^2)$ to obtain one with the correct variance of $\kappa^2/(8\sigma^2\pi^2)$:

$$\hat{B}(x, y, f_\xi, f_\eta) = \iint B(x, y, \alpha, \beta)\frac{4\sigma^2\pi}{\kappa^2 - 1}e^{-4\sigma^2\pi^2\left((f_\xi - \alpha)^2 + (f_\eta - \beta)^2\right)/(\kappa^2 - 1)}d\alpha d\beta \quad (3.25)$$

Recall that a partially coherent beam can be analyzed as the incoherent ensemble of multiple coherent beams and that addition of Wigner distribution functions implies incoherent addition. Convolution along the frequency axis in the Wigner distribution creates replicas of the original Wigner distribution at various offsets in frequency. Therefore, the convolution in Equation (3.25) indicates that this "compressed" Gaussian beam can be constructed as an incoherent ensemble of many Gaussian beams having the same beam waist as the original beam. Each of these constituent beams has been modified by an amplitude scale factor as well as a linear phase factor (due to the shift in frequency) that rotates the central axis of the Gaussian beam away from the optical axis. Hence, we have shown that any longitudinal compression of a particular coherent beam causes it to no longer to be a valid coherent beam, but a physically plausible partially coherent field can represent this compressed light pattern. These results are summarized pictorially in Figure 3.6. Do note that since Gaussian convolution can only increase the variance, the same trick cannot be applied in the case of an "expansion" of the Gaussian beam along the optical axis.

In fact, the partially coherent field represented by this stretched Wigner distribution is a member of a family of well-known partially coherent sources called Gaussian Schell-model sources, where the degree of spatial coherence is a spatially invariant Gaussian function and the intensity profile is a Gaussian function as well [55–57]. This can be made obvious by examining the mutual intensity corresponding to the Wigner distribution in Equation (3.25) by performing the inverse Fourier transform along $f_\xi$ and $f_\eta$, which converts the convolution to a multiplication:

$$\hat{J}(x + \tfrac{\xi}{2}, y + \tfrac{\eta}{2}, x - \tfrac{\xi}{2}, y - \tfrac{\eta}{2}) = J(x + \tfrac{\xi}{2}, y + \tfrac{\eta}{2}, x - \tfrac{\xi}{2}, y - \tfrac{\eta}{2})e^{-\frac{\kappa^2 - 1}{4\sigma^2}\left(\xi^2 + \eta^2\right)} \quad (3.26)$$

Figure 3.6: A Gaussian beam, whose Wigner distribution is shown in (i), can be compressed along the optical axis to obtain a new Wigner distribution, (ii). Since only the projection along one axis changed, this cannot be the Wigner distribution of a (coherent) function. However, we can construct this Wigner distribution in (ii) by adding multiple copies of the original Wigner distribution, which are scaled and shifted (iii), due to the fact that convolution of two Gaussians simply adds the variances. This addition in phase space causes the creation of a partially coherent beam, and hence coherence is the main limiting obstacle to creating the desired Wigner distribution in (ii).

Conversion back to the the two coordinate format from the mean and difference format yields the well-known form of the mutual intensity of a Gaussian-Schell model source:

$$\hat{J}(x_1, y_1, x_2, y_2) = e^{-\frac{x_1^2+y_1^2}{2\sigma^2}} e^{-\frac{x_2^2+y_2^2}{2\sigma^2}} e^{-\frac{\kappa^2-1}{4\sigma^2}\left((x_1-x_2)^2+(y_1-y_2)^2\right)} \tag{3.27}$$

$$= I(x_1, y_1)^{1/2} I(x_2, y_2)^{1/2} \mu(x_1 - x_2, y_1 - y_2) \tag{3.28}$$

where $I(x,y) = e^{-\frac{x^2+y^2}{\sigma^2}}$ is the Gaussian intensity profile and

$$\mu(\Delta_x, \Delta_y) = e^{-\frac{\kappa^2-1}{4\sigma^2}\left(\Delta_x^2+\Delta_y^2\right)} \tag{3.29}$$

is the Gaussian degree of spatial coherence.

## 3.2.2 Three-dimensional limitations

Not only are arbitrary mutual intensity functions unrealizable for coherent fields, arbitrary bandlimited three-dimensional wave functions are also not realizable. It is well known that the Fourier transform of a coherent three-dimensional scalar field far away from any evanescent sources yields a three-dimensional function that is zero everywhere except on the surface of a sphere with radius $1/\lambda$, where $\lambda$ is the wavelength of the optical field in question. In other words, a two-dimensional manifold can describe any realizable coherent three-dimensional field. Therefore, there many of three-dimensional fields that cannot be realized [58,59], e.g. a "plane wave" whose wavelength is twice as long as the wavelength of the field:

$$U(x, y, z) = e^{\frac{j2\pi z}{2\lambda}} \tag{3.30}$$

The limitation arises from Huygens's principle, where a coherent field convolved with a spherical wave must result in the same coherent field. The Fourier transform of a spherical wave lies entirely on the surface of a sphere of radius $1/\lambda$. Therefore, through the convolution theorem, any physically valid field must also be on the surface of this sphere. In other words, enforcing that the three-dimensional scalar field is a

proper solution to the Helmholtz equation results in the loss of one dimension.

Intuitively, full coherence implies that light from *any* two points will interfere with each other. Since light propagates, the field at one point interacts with the field at every other point. In order for such a system to be stable, perturbations of the field at one point would necessarily involve perturbations of the field at many other points, much in the same way that compressing a water balloon at one point would cause other parts of the water balloon to bulge.

However, in many application areas, only the intensity and not the phase of the coherent field is important. Since the intensity is simply the field multiplied by its complex conjugate, this means that the possible extent of the Fourier transform of the intensity is equivalent to the autocorrelation of the hollow sphere. This operation does "fill" three-dimensional space, making it more difficult to determine a specific pattern that would be impossible to generate using a fully coherent field. However, recall that this autocorrelation operation is still a function from a two-dimensional manifold to a three-dimensional pattern. Therefore, the set of possible three-dimensional intensity patterns must have size less than or equal to the set of possible three-dimensional field patterns and thus the same limitations still apply, although the freedom to choose the phase may yield potential gains.

With the discussion so far, limitations of coherent fields have been the result of obvious dimension-mismatch issues, but in the following section, we will demonstrate that not all bandlimited *two-dimensional* intensity patterns can be generated by a coherent field, either. This conclusion will make it obvious that not all bandlimited three-dimensional intensity patterns can be generated either, as any three-dimensional intensity pattern containing an impossible two-dimensional intensity pattern will also be impossible to generate using a coherent field.

### 3.2.3 Two-dimensional limitations

We will now investigate limitations on two-dimensional intensity patterns due to full coherence. More specifically, we will consider a planar slice $U(x, y)$ of a fully coherent field propagating in the $+z$ direction, where the maximum horizontal and maximum

vertical components of spatial frequency are strictly less than a specific frequency $f_{max}$, i.e. the two-dimensional Fourier transform of $U(x, y)$ has square spatial support with a width and height of $2f_{max}$. For example, the output plane of an optical Fourier transformer possesses this property if the input plane is fully coherent and contains negligible energy outside a square shaped region. We will refer to such a field by the term *square-bandlimited coherent field*.

Since intensity is the product of the field with its complex conjugate, the support of the Fourier transform of the intensity must be a square with width and height $4f_{max}$ due to the convolution theorem. By Nyquist's sampling theorem, a sampled lattice at intervals of $\Delta = 1/(4f_{max})$ is sufficient to fully describe the intensity pattern in this plane. This sampling pattern is also obviously sufficient to fully specify the field $U(x, y)$, since the field has a smaller bandlimit. It is the goal of this section to demonstrate a family of intensity patterns that satisfy the Fourier domain support requirements but cannot be generated via a square-bandlimited coherent field.

**Sampling and degrees of freedom**

Let us consider a square-bandlimited coherent field with negligible energy outside a $2N\Delta \times 2N\Delta$ square region, as shown in Figure 3.7. Critically sampling the field (i.e. no two samples are dependent) implies a sampling lattice pattern with sampling intervals of $2\Delta$ along each axis. This results in $N^2$ unique samples of the field inside this square region, and the field inside this region can be controlled simply by specifying these $N^2$ complex values. However, the sampling interval required to sample the intensity pattern is $\Delta$ along each dimension, giving rise to $4N^2$ sample points inside this square region. That is, in order to specify all possible intensity patterns in this square region, we would need to specify $4N^2$ independent real values. Since we only have the ability to control $N^2$ complex values via the field, this means we can only fully specify $2N^2$ of the $4N^2$ samples of the intensity. Thus, there must exist some intensity patterns that cannot be realized in this situation, even if the intensity patterns are bandlimited properly. Furthermore, this mismatch is always present, regardless of the size of the region chosen (i.e. what $N$). In other words, for a two-dimensional square-bandlimited coherent field, the increased sampling rate needed to fully specify

Figure 3.7: In a $2N\Delta \times 2N\Delta$ square region of a planar slice of a square-bandlimited coherent field, critical sampling of the field yields $N^2$ complex-valued samples, and thus yielding $2N^2$ degrees of freedom of control. However, the doubled sampling rate along each dimension for the intensity requires $4N^2$ real samples to properly specify the intensity. Since we only have $2N^2$ degrees of freedom of control in this region, we cannot specify all possible intensity patterns through controlling the coherent field.

the intensity causes a mismatch in the degrees of freedom available through control of the coherent field and the degrees of freedom needed to fully specify the intensity.

**An unrealizable pattern**

With this degree of freedom mismatch in mind, we will now formally prove that unrealizable bandlimited intensity patterns do exist. Given a square-bandlimited coherent field $U(x, y)$, let us define $U[m, n]$ to be the unique discretization of this field and $I[m, n]$ be the unique discretization of the intensity of this field:

$$U[m, n] = U(m\Delta, n\Delta) \ , \ I[m, n] = I(m\Delta, n\Delta) = |U(m\Delta, n\Delta)|^2 \qquad (3.31)$$

Let $S[m, n]$ be the discretization of the field $S(x, y)$ of a diffraction limited spot that can be generated on this plane given the bandlimit constraints on the field:

$$S[m, n] = S(m\Delta, n\Delta) = (mn\pi^2/4)^{-1} \sin(m\pi/2) \sin(n\pi/2) \qquad (3.32)$$

where $\Delta = 1/(4f_{max})$ is the sampling interval required for the intensity (and thus is sufficient for the field) as defined previously. Note that $S[m, n]$ is zero when either $m$ or $n$ is even, unless $m = n = 0$, in which case it has value 1.

We will show in the following section that a square bandlimited coherent field cannot produce an intensity pattern formed by summing three diffraction limited spots of various brightness centered at points $P_1,$, $P_2$ and $P_3$ with respective discrete coordinates $(m_1, n_1)$, $(m_2, n_2)$ and $(m_3, n_3)$:

$$I_0[m, n] = \sum_{i=1}^{3} I_i \, |S[m - m_i, n - n_i]|^2 \qquad (3.33)$$

with the restriction that $m_i$ and $n_i$ are even integers, $I_i > 0$, the points $P_1$, $P_2$ and $P_3$ are unique and that the three points are not in the same row (same $n$) or column (same $m$). Symmetries in the system ($m$ and $n$ coordinates can be swapped and the signs of the $m$ and $n$ coordinates can be flipped) allow for a simpler set of assumptions

without loss of generality:

$$\forall i, I_i > 0 \tag{3.34}$$

$$n_1 < n_2 \le n_3 \tag{3.35}$$

$$|m_1 - m_2| + |m_2 - m_3| > 0 \tag{3.36}$$

$$|n_2 - n_3| + |m_2 - m_3| > 0 \tag{3.37}$$

Condition (3.35) ensures that at most two points share the same row, and the ordering can be swapped via reflection along $n$. Condition (3.36) ensures that at most two points share the same column, with no restrictions on ordering. Condition (3.37) ensures that $P_2$ and $P_3$ are not the same point. Examples of intensity patterns which fit these criteria are shown in Figure 3.8.



(i)                                (ii)

Figure 3.8: Examples of 2D bandlimited intensity patterns that cannot be generated by a square-bandlimited coherent field. The intensity pattern is shown using a heat map palette, where the color changes from black to red to yellow to white with increasing intensity.

**Proof**

We will now proceed with a proof by contradiction, showing that the intensity patterns described in the previous section is not realizable via a square-bandlimited coherent field. First, assume that a square-bandlimited coherent field discretization $U_0[m, n]$

exists such that:

$$|U_0[m,n]|^2 = I_0[m,n] \tag{3.38}$$

If such a field exists, then it should have amplitude equal to the square root of the desired intensity $I_0$ and some phase value which we are free to assign at each sample point. We will show that it is impossible to assign phase values to ensure that both the desired intensity is satisfied (3.38) and that the resulting field is correctly bandlimited.

Let us define without loss of generality the phase of the discretized wave function $U_0[m,n]$ at the three diffraction spot centers $P_1$, $P_2$ and $P_3$ to be $\phi_1$, $\phi_2$ and $\phi_3$, respectively. Given these values, we can first determine the phase of the points in the same column as one of the points $P_i$ ($U_0[m_i,n]$ with $n$ varying) and the points in the same row as one of the points $P_i$ ($U_0[m,n_i]$ with $m$ varying).

We start by considering the field at points $[m,n]$ where both $m$ and $n$ are even. Recall that $S[m,n]$ is zero when either $m$ and $n$ are even unless $m = n = 0$, when it is 1. Thus, this means that the desired intensity $I_0$ at a point where $m$ and $n$ are both even contains contributions from at most one of the diffraction limited spots; i.e. the summation in Equation (3.33) has at most one nonzero term. That is, $U_0[m,n]$ is either:

- $I_i^{1/2}e^{j\phi_i}$, if $m = m_i$ and $n = n_i$ for some diffraction spot $P_i$.

- 0 otherwise.

Using this observation, we can conclude that for the column $U_0[m_i,n]$ in which diffraction spot $P_i$ resides, and for $n$ even (and $m_i$ is obviously even because of where the diffraction spot centers are located):

- if $n = n_i$, then $U_0[m_i,n] = I_i e^{j\phi_i}$, since the location of zeros in $S[m',n']$ removes any contributions from any other diffraction spot.

- if $n = n_j$ for some other diffraction spot $P_j$ that happens to share the same column as $P_i$, then $U_0[m_i,n] = I_j^{1/2}e^{j\phi_j}$ by the same reasoning.

- otherwise, $U_0[m_i,n] = 0$

Now let us consider the phase of $U_0[m_i, n]$ when $n$ is odd. Since $U_0[m_i, n]$ is a discretized axis-aligned slice of a square-bandlimited function, these one-dimensional discrete representations must also have maximum frequency $f_{max}$. Therefore, convolving the pulse train $\sum_n U_0[m_i, n]\delta(x, y - n\Delta)$ (conversion of the discrete representation to the continuous domain) by the bandlimited sinc function $S(0, y)$ (which results in ideal reconstruction for any signal with maximum frequency $f_{max}$) and then sampling at intervals of $\Delta$ should result in $2U_0[m_i, n]$. The extra scale factor comes from the $2\times$ oversampling. This convolution and sample process is equivalent to a discrete convolution:

$$U_0[m_i, n] = \sum_{l \neq n} U_0[m_i, l]S[0, n - l] \tag{3.39}$$

Since $S[0, n - l]$ is zero if $n \neq l$ and $n - l$ is even, this means that the summation's only nonzero entries are when $l$ is even. From the previous discussion, $U_0[m_i, l]$ only has nonzero values when it is on top of a diffraction spot if $l$ is even. Thus, we can rewrite the infinite sum in the previous equation to a finite sum, as the majority of the entries are zero:

$$U_0[m_i, n] = \sum_{j=1}^{3} I_j^{1/2} e^{j\phi_j} S[m_j - m_i, n_j - n] \tag{3.40}$$

This includes a summation for all three spots, but the $S[m_j - m_i, n_j - n]$ quantity ensures no contributions if $m_j \neq m_i$, i.e. if the diffraction spot doesn't lie on our current column under consideration. The same analysis can be applied to the rows as well to obtain a similar result:

$$U_0[m, n_i] = \sum_{j=1}^{3} I_j^{1/2} e^{j\phi_j} S[m_j - m, n_j - n_i] \tag{3.41}$$

With these results, we now know the phase of the field $U_0$ at all rows and columns which contain one of the diffraction spot centers.

Now let's consider some point $U_0[\hat{m}, \hat{n}]$ where $\hat{m}$ and $\hat{n}$ are odd integers. By applying the same convolution-invariance argument with regards to the band limit in the two dimensional regime (and hence a factor of four due to oversampling), we

arrive at an expression similar to Equation (3.39):

$$3U_0[\hat{m}, \hat{n}] = \sum_{k,l \neq \hat{m}, \hat{n}} U_0[k, l] S[\hat{m} - k, \hat{n} - l] \tag{3.42}$$

$S[\hat{m} - k, \hat{n} - l]$ is nonzero only:

- when both $k$ and $l$ are even (case i)

  In this case, $U_0[k, l]$ will only be nonzero at the diffraction spot centers $(m_i, n_i)$.

- when $\hat{m} = k$ and $l$ is even (case ii)

  In this case, $U_0[k, l]$ will only be nonzero at the points $(k, l) = (\hat{m}, n_i)$, where $n_i$ is the coordinate of one or more diffraction spot centers

- when $\hat{n} = l$ and $k$ is even (case iii)

  In this case, $U_0[k, l]$ will only be nonzero at the points $(k, l) = (m_i, \hat{n})$, where $m_i$ is the coordinate of one or more diffraction spot centers

Combining the above cases with Equation (3.42) yields:

$$3U_0[\hat{m}, \hat{n}] = \sum_{i=1}^{3} U_0[m_i, n_i] S[\hat{m} - m_i, \hat{n} - n_i] + U_0[\hat{m}, n_i] S[0, \hat{n} - n_i] + U_0[m_i, \hat{n}] S[\hat{m} - m_i, 0]$$
$$\tag{3.43}$$

Substituting equations (3.40) and (3.41) into the above and further simplification yields:

$$U_0[\hat{m}, \hat{n}] = \sum_{i=1}^{3} I_i^{1/2} e^{j\phi_i} S[\hat{m} - m_i, \hat{n} - n_i] \tag{3.44}$$

Applying Eq. (3.33) to the above yields:

$$\sum_{i=1}^{3} I_i \left| S[\hat{m} - m_i, \hat{n} - n_i] \right|^2 = \left| \sum_{i=1}^{3} I_i^{1/2} e^{j\phi_i} S[\hat{m} - m_i, \hat{n} - n_i] \right|^2 \tag{3.45}$$

Equation (3.45) states that at discrete locations where $m$ and $n$ are both odd, the expected output of the incoherent sum of three diffraction spot intensities is equal to the intensity of the coherent sum of three diffraction spots. This result can be further

simplified to the following linear form via an extended version of the law of cosines:

$$\alpha_{1,2}\psi_{1,2} + \alpha_{1,3}\psi_{1,3} + \alpha_{2,3}\psi_{2,3} = 0 \tag{3.46}$$

where $\alpha_{i,j} = S[\hat{m} - m_i, \hat{n} - n_i]S[\hat{m} - m_j, \hat{n} - n_j]$ and $\psi_{i,j} = (I_i I_j)^{1/2} \cos(\phi_i - \phi_j)$.

Thus, a collection of such candidate points $(\hat{m}, \hat{n})$ forms a system of linear equations. We will now show that it is possible to pick four points such that the corresponding linear system of equations has only one solution, i.e. the matrix whose entries are $\alpha_{i,j}$ has rank equal to 3. This single solution would be where all the $\psi_{i,j}$ are zero, leading to a contradiction that will be shown at a later point.

We will now show that by choosing some positive odd integer $\beta$, we can ensure that the four points $Q_1 = (m_1 - \beta, n_1 - 1)$, $Q_2 = (m_2 + \beta, n_1 - 1)$, $Q_3 = (m_1 - \beta, n_2 + 1)$ and $Q_4 = (m_2 + \beta, n_2 + 1)$ creates a rank-3 matrix. This particular choice of points will always yield a matrix whose first column has a single value across all rows, regardless of $\beta$:

$$\begin{pmatrix} abcd & aceh & bdeh \\ abcd & adeg & bceg \\ abcd & bcfh & adfh \\ abcd & bdfg & acfg \end{pmatrix} \begin{pmatrix} \psi_{1,2} \\ \psi_{1,3} \\ \psi_{2,3} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} \tag{3.47}$$

where

$$
\begin{aligned}
a &= S[0,1] & b &= S[0, n_2 - n_1 + 1] & e &= S[0, n_3 - n_1 + 1] & f &= S[0, n_3 - n_2 - 1] \\
c &= S[\beta, 0] & d &= S[m_2 - m_1 + \beta, 0] & g &= S[m_3 - m_1 + \beta, 0] & h &= S[m_3 - m_2 - \beta, 0]
\end{aligned}
\tag{3.48}
$$

All eight values in Eq. (3.48) are nonzero because one coordinate for the $S[m, n]$ function is 0 and the other is odd. Subtracting the even rows from the odd rows yields a linear system involving $\psi_{1,3}$ and $\psi_{2,3}$:

$$\begin{pmatrix} ae(ch - dg) & be(dh - cg) \\ bf(ch - dg) & af(dh - cg) \end{pmatrix} \begin{pmatrix} \psi_{1,3} \\ \psi_{2,3} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \tag{3.49}$$

The determinant $D$ of the above $2 \times 2$ matrix is: $ef(a^2 - b^2)(ch - dg)(dh - cg)$. Since

$n_1 < n_2$ according to the assumption in Equation (3.35), $a^2 - b^2 < 0$. As stated before, neither $e$ nor $f$ can be zero. For $D$ to be zero, at least one of the remaining two terms would have to be zero, and therefore either $c/g = d/h$ or $d/g = c/h$. If at least one of those were true, then at least one of the following equations has to be true (derived by plugging in the definition of $S[m, n]$):

$$|\beta/(\Delta_1 + \beta)| \;=\; |(\Delta_1 + \Delta_2 + \beta)/(\Delta_2 - \beta)| \tag{3.50}$$

$$|\beta/(\Delta_1 + \beta)| \;=\; |(\Delta_2 - \beta)/(\Delta_1 + \Delta_2 + \beta)| \tag{3.51}$$

where $\Delta_1 = m_2 - m_1$ and $\Delta_2 = m_3 - m_2$. For at least one of the above equations to be true, then at least one of the four following equations must be true:

$$2\beta^2 + 2\Delta_1\beta + \Delta_2(\Delta_1 + \Delta_2) = 0$$
$$2\beta(\Delta_1 + \Delta_2) + \Delta_2(\Delta_1 + \Delta_2) = 0$$
$$2\beta^2 + 2\Delta_1\beta - \Delta_1\Delta_2 = 0$$
$$2\Delta_2\beta + \Delta_1\Delta_2 = 0$$

$$\tag{3.52}$$

To proceed with the proof, let's consider two possibilities separately – whether $\Delta_2$ is zero or not.

**When $\Delta_2 \neq 0$...**

If $\Delta_1$ is also nonzero, then there can be at most 6 unique values of $\beta$ that satisfy at least one of the above equations, since none of them are degenerate in this case. If $\Delta_1$ is zero, then the above four equations can be reformulated as:

$$2\beta^2 + \Delta_2^2 = 0 \tag{3.53}$$
$$2\beta\Delta_2 + \Delta_2^2 = 0 \tag{3.54}$$
$$2\beta^2 = 0 \tag{3.55}$$
$$2\Delta_2\beta = 0 \tag{3.56}$$

The latter two are impossible because $\beta$ has to be an odd integer. This leaves us with

at most 3 unique values of $\beta$ that satisfies at least one of the equations in (3.52).

Thus, only a finite number of $\beta$ satisfies at least one of the equations in (3.52) and there are an infinite number of $\beta$ we can pick. Therefore, $D$ can be made nonzero by choosing beta from a sequence of ascending odd integers until all four equations are false in (3.52).

Since we can make $D \neq 0$ by choosing a particular $\beta$, the matrix in Eq. (3.49) can be made full rank and thus the system has a unique solution $\psi_{1,3} = \psi_{2,3} = 0$. This implies that angles $\phi_1$ and $\phi_3$ are orthogonal, and that angles $\phi_2$ and $\phi_3$ are also orthogonal. Therefore, $\phi_1$ and $\phi_2$ cannot be orthogonal as well. However, plugging $\psi_{1,3} = \psi_{2,3} = 0$ into Eq. (3.47) results in arriving at the conflicting conclusion that $\psi_{1,2} = 0$ as well. It's conflicting because all of the $I_i$ are positive and we cannot have three angles that are all orthogonal to each other, hence making the cosine terms necessarily not all zero in the definition of the $\psi$ terms.

**When $\Delta_2 = 0$...**

In this case, $\Delta_1$ has to be nonzero, because otherwise the assumption given by Equation (3.36) would be violated. Now if $\Delta_1 \neq 0$, then we can subtract the even rows from each other and the odd rows from each other in (3.47) to obtain:

$$\begin{pmatrix} ch(ae - bf) & dh(be - af) \\ dg(ae - bf) & cg(be - af) \end{pmatrix} \begin{pmatrix} \psi_{1,3} \\ \psi_{2,3} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \tag{3.57}$$

The determinant in this case is $gh(c^2 - d^2)(ae - bf)(be - af)$, and since $\Delta_1 = m_1 - m_2 \neq 0$, $c^2 - d^2$ is nonzero, the determinant can only be zero if at least one of the following equations is true:

$$|(\Delta_3 + 1)/1| = |(\Delta_3 + \Delta_4 + 1)/(\Delta_4 + 1)| \tag{3.58}$$

$$|(\Delta_3 + 1)/1| = |(\Delta_4 + 1)/(\Delta_4 + \Delta_3 + 1)| \tag{3.59}$$

where $\Delta_4 = n_3 - n_2 \geq 0$ and $\Delta_3 = n_2 - n_1 > 0$. Since we assumed that $\Delta_2 = 0$, then $\Delta_4 \neq 0$ because of assumption (3.37) and $\Delta_4 \geq 0$ because of assumption (3.35). Thus we know that $\Delta_4 > 0$. Thus, (3.59) has to be false, because $|(\Delta_4 + 1)/(\Delta_4 + \Delta_3 + 1)| < 1$ and $|(\Delta_3 + 1)/1| > 1$. For (3.58) to be true, then at

least one of the following equations must be true:

$$\Delta_3 \Delta_4 = 0 \qquad (3.60)$$

$$\Delta_4 = -\frac{2 + 2\Delta_3}{2 + \Delta 3} \qquad (3.61)$$

Neither can be true because the former needs $\Delta_4 = 0$ and the latter needs $\Delta_4 < 0$. Thus, (3.57) has only a unique solution. Following similar reasoning for when $\Delta_2$ was not zero, we obtain that $\psi_{1,3} = \psi_{2,3} = 0$ and the conflicting conclusion from plugging into (3.47) that $\psi_{1,2} = 0$ as well.

Combining the results from the above two situations yields the conclusion that there exists no valid $\phi_1, \phi_2, \phi_3$ for the phases at points $P_1, P_2, P_3$ that allow for the intensity patterns at an additional four points $Q_1, Q_2, Q_3, Q_4$ to match our desired intensity $I_0$ and still maintain $U_0$'s bandlimit. Hence, our desired intensity $I_0$ is not realizable with the described square-bandlimit coherent field.$\square$

## 3.2.4   Discussion

The above proof demonstrates that there are a family of two-dimensional patterns that cannot be realized with a square-bandlimited coherent field , even though the pattern itself does not contain a spatial frequency greater than the maximum possible spatial frequency realizable using a coherent field. This family of patterns is trivially realizable with partially coherent fields, since a partially coherent field can be viewed as as an incoherent sum (i.e. summing in intensity) of fully coherent fields. Thus, a partially coherent field whose modes correspond to individual diffraction spots generates the desired intensity pattern. Hence, this limitation on two-dimensional intensity patterns must be due to the full coherence of the field. The structure of the proof lends to the possibility of proofs for patterns involving more diffraction spots, but the complexity of the proof would increase drastically in that case. Instead, let us consider the question from the opposite viewpoint – what are some patterns that *are* feasible with coherent fields?

**Provably realizable patterns**

In the previous proof, recall that we specifically asked for three distinct diffraction spots. Obviously, a single diffraction spot can be created by a coherent field, but a two diffraction spot pattern is also easily realizable. The following intensity pattern consisting of the intensity sum of two diffraction limited spots,

$$I(x, y) = I_1 S^2(x - x_1, y - y_1) + I_2 S^2(x - x_2, y - y_2) \tag{3.62}$$

can be easily generated by the coherent sum of one "real" diffraction spot and one "imaginary" diffraction spot:

$$U(x, y) = I_1^{1/2} S(s - x_1, y - y_1) + j I_2^{1/2} S(x - x_2, y - y_2) \tag{3.63}$$

owing to the fact that fields corresponding to diffraction spots are real and that total intensity is an "incoherent sum" of the real and imaginary intensities.

Also recall that we asked for three diffraction spots that did not all lie in the same row or column. Therefore, the following intensity pattern consisting of three diffraction spots lying in the same row would violate that assumption:

$$I[m, n] = \sum_{i=1}^{3} I_i S^2[m - m_i, n - n_0] \tag{3.64}$$

Note that this intensity pattern is *separable*. That is:

$$
\begin{aligned}
I[m, n] &= S^2[0, n - n_0] \sum_{i=1}^{3} I_i S^2[m - m_i] \\
&= I_y[n] I_x[m] \tag{3.65}
\end{aligned}
$$

In order to generate a separable intensity pattern, it is sufficient to find a separable coherent field:

$$U[m, n] = U_y[n] U_x[m] \tag{3.66}$$

such that

$$|U_y[n]|^2 = I_y[n] \tag{3.67}$$

$$|U_x[m]|^2 = I_x[m] \tag{3.68}$$

$U_y[n]$ can simply be $S[0, n-n_0]$. Finding $U_x[m]$ is more complicated, but it is relatively well known in the optics community that the one-dimensional phase-retrieval problem from the intensity is theoretically solvable [60]. Furthermore, this should make sense, since for a region in the one-dimensional case containing N samples of the field, there would be only 2N samples of the intensity, yielding a match in the total degrees of freedom. This result yields an interesting corollary, which is that if a discretization $I[m, n]$ of a bandlimited 2D intensity pattern can be written as a separable function:

$$I[m, n] = I_x[m]I_y[n] \tag{3.69}$$

then this intensity pattern is realizable using a separable square-bandlimited fully coherent field.

**Convexity of set of possible intensities**

We've been talking about the set of possible intensity patterns that can be realized using a square bandlimited coherent field as well as ones realizable by the partially coherent analog. Let's denote $\mathcal{I}_\mathcal{C}$ as the set of all possible coherent intensity patterns with total intensity equal to 1 and $\mathcal{I}_\mathcal{P}$ as the set of all possible partially coherent intensity patterns with total intensity equal to 1. According to coherent mode theory, any partially coherent intensity pattern can always be written as the sum of fully coherent ones:

$$\forall I_P \in \mathcal{I}_\mathcal{P}, \exists I_C^{(i)}, \alpha_i, \beta \geq 0 \text{ such that } \beta I_P = \sum_i \alpha_i I_P \text{ and } \sum_i \alpha_i = \beta \tag{3.70}$$

Therefore, the set $\mathcal{I}_\mathcal{P}$ must be the convex hull of the set $\mathcal{I}_\mathcal{C}$ and is thus convex. Furthermore, we know that every coherent intensity pattern can be written as a partially

Figure 3.9: The set of all possible coherent intensity patterns, $\mathcal{I}_\mathcal{C}$ with unit total intensity is illustrated in the shaded region inside a solid border.  The set of all possible partially coherent intensity patterns $\mathcal{I}_\mathcal{P}$ forms the convex hull of the set $\mathcal{I}_\mathcal{C}$ and is shown inside the dotted border. The diagram is a pictorial representation (and not to scale) of a two-dimensional slice of the hyper-dimensional set, where each axis is the intensity of a point in space. For example, the two green points indicate intensity patterns possible with a coherent field, with the intensity patterns shown using a heat map palette. The red point refers to an intensity pattern that is impossible with a coherent field, and it is simply the average intensity of the two possible intensity patterns.

coherent intensity pattern with a single mode, and thus the following statement is true:

$$\mathcal{I}_\mathcal{C} \subseteq \mathcal{I}_\mathcal{P} \tag{3.71}$$

Since we've shown that there exists partially coherent patterns $I_P \in \mathcal{I}_\mathcal{P}$ that cannot be realized with full coherence, the sets $\mathcal{I}_\mathcal{P}$ and $\mathcal{I}_\mathcal{C}$ cannot be equal, and hence the set of possible intensity patterns realizable using full coherence is only a proper subset of the set of possible intensity patterns realizable using partial coherence:

$$\mathcal{I}_\mathcal{C} \subset \mathcal{I}_\mathcal{P} \tag{3.72}$$

This result is summarized graphically in Figure 3.9.

## 3.3   Volumetric devices

The final family of devices we shall consider are volumetric displays. Use of these devices requires modeling illumination as three-dimensional patterns emitted by a volume of light emitting voxels. One benefit of this model is that it is intuitive to specify three-dimensional scenes. However, the following analysis using partial coherence representations in phase space will demonstrate limitations on illumination patterns that can be generated by these devices.

### 3.3.1   Phase space derivation

We will first derive an expression for the ambiguity function of the illumination generated by a volumetric display device. Volumetric display devices effectively generate a set of incoherent point emitters where light from each emitter causes no interference with light from any other emitter. These point emitters can either be real emitters in space (such as an excited fluorophore) or projected images of emitters. Since there is no interference between separate emitters, the phase space representation of the illumination can be obtained via summing phase space representations of each emitter or group of emitters.

Let us consider a set of point emitters located on the $z = \hat{z}$ plane, with intensity given by the two-dimensional function $I(x, y; \hat{z})$. Then, the ambiguity function corresponding to this set of point emitters is:

$$A_0(f_x, f_y, \xi, \eta; \hat{z}) = \tilde{I}(f_x, f_y; \hat{z})\delta(\xi + \hat{z}\lambda f_x, \eta + \hat{z}\lambda f_y) \tag{3.73}$$

where $\tilde{I}(f_x, f_y; \hat{z})$ is the two-dimensional Fourier transform along $x$ and $y$ of $I(x, y; \hat{z})$. This represents a two-dimensional plane embedded in a four-dimensional space.

Integrating over contributions across all planes $z = \hat{z}$ in space yields the following

expression for the ambiguity function:

$$A_1(f_x, f_y, \xi, \eta) = \int A_0(f_x, f_y, \xi, \eta; \hat{z}) d\hat{z} \tag{3.74}$$

$$= \int \tilde{I}(f_x, f_y; \hat{z}) \delta(\xi + \hat{z}\lambda f_x, \eta + \hat{z}\lambda f_y) d\hat{z} \tag{3.75}$$

$$= \iint \tilde{I}(f_x, f_y; \hat{z}) \delta(\xi + \hat{z}\lambda f_x, \eta + z'\lambda f_y) \delta(z - \hat{z}) d\hat{z} dz' \tag{3.76}$$

$$= \frac{1}{\lambda^2 f_x f_y} \iint \tilde{I}(f_x, f_y; \xi'/(\lambda f_x))$$
$$\times \delta(\xi - \xi', \eta - \eta') \delta(\xi'/(\lambda f_x) - \eta'/(\lambda f_y)) d\xi' d\eta' \tag{3.77}$$

$$= \frac{1}{\lambda^2 f_x f_y} \tilde{I}(f_x, f_y; \xi/(\lambda f_x)) \delta(\xi/(\lambda f_x) - \eta/(\lambda f_y)) \tag{3.78}$$

The delta function in the result is worth noting, since it shows that the ambiguity function is zero at locations where:

$$\xi/(\lambda f_x) \neq \eta/(\lambda f_y) \tag{3.79}$$

That is, the ambiguity function is only nonzero in a three-dimensional subset of the four-dimensional space.

In practice, optical systems are bandlimited, and these point emitters in volumetric displays are not usually entirely isotropic. Therefore, let's also apply an exit pupil function $\tilde{a}(f_x, f_y)$. Let $A_a(f_x, f_y, \xi, \eta)$ be the ambiguity function corresponding to the inverse 2D Fourier transform $a(x, y)$ of the exit pupil function. Then, the ambiguity function including these effects is simply:

$$A_2(f_x, f_y, \xi, \eta) = \iint A_a(f_x, f_y, \xi - \xi', \eta - \eta') A_1(f_x, f_y, \xi', \eta') d\xi' d\eta' \tag{3.80}$$

We expect the exit pupil function to be large enough so that most of the details in the images themselves are retained. Therefore, $A_a(f_x, f_y, \xi, \eta)$ should be a function that is wide in $f_x, f_y$ and short in $\xi, \eta$. Thus, what this function does to the ambiguity function of the plane is that it induces a small amount of approximately constant vertical blur along the $\xi, \eta$ dimensions, creating a "fuzzy" three-dimensional shape in

Figure 3.10: The ambiguity function formed by a volumetric display is shown here pictorially. The graph in (i) shows the ambiguity function projected along the $f_y$ and $\eta$ axes onto $f_x - \xi$ space. Note that the ambiguity function is blurred vertically. The green circle in (i) indicates the position $(f_{x_0}, \xi_0)$ and the graph in (ii) is a slice of the ambiguity function at that position, i.e. what had been projected down onto the point circled in (i). The green line is simply a line through the green circle in (i); note that the graph in (ii) is clustered around a line of the same slope through the origin. The dotted curves circle difficult areas in the ambiguity function for volumetric devices to control. Volumetric devices have difficulty creating light patterns with anisotropic/astigmatic/occlusion effects (A) and also have a difficult time controlling out-of-focus blur (B) caused by the vertical blur in (i).

four-dimensional space.

## 3.3.2   Phase space analysis

We will now analyze the result obtained in the previous section and derive insight as to what types of illumination volumetric displays can and cannot generate. Figure 3.10 (i) shows a pictorial representation of the ambiguity function in (3.80) projected along the $f_y$ and $\eta$ axes, and (ii) shows a slice of the ambiguity function along the $f_y$ and $\eta$ axes when $f_x = f_{x_0}$ and $\xi = \xi_0$. Ideally, the set of all possible illumination patterns should allow for nearly arbitrary patterns in the ambiguity function as well. These images show very clearly some issues in attaining that goal.

As we look at (i), recall that there is a vertical blur along the $\xi$ direction due to the limited aperture. Since the planes are more tightly packed near the origin, but the blur size is approximately the same, this means there's more "cross-talk" between the different planes at these lower frequencies. Physically, this means there's out of focus blur in a volumetric display – e.g. it is impossible to generate a one-dimensional intensity pattern along the optical axis that is alternating light and pure darkness.

Another, more serious issue can be seen in (ii). Recall that we had a three-dimensional shape imposed on the ambiguity function before applying the exit pupil. Since the exit pupil only slightly blurs this in the $\eta$ direction, most of the space in this slice is blank. Therefore there must be a lot of different light patterns that cannot be produced using volumetric displays. Physically, when we set $\xi/(\lambda f_x)$ equal to $\eta/(\lambda f_y)$, we are saying the apparent depth of the light is the same whether we take horizontal or vertical slices of the light. That is, the light we see is *stigmatic*. Nonzero values outside of this line-shaped region in (i) result in *astigmatic* light. This should be unsurprising, as the optical systems in volumetric displays simply generate points in space through stigmatic optics.

The following ambiguity function derivation of a scene consisting of a plane of uniform point emitters followed by an occluder will demonstrate that the empty space in (i) precludes occlusions as well. This should make sense, since light from every emitter (real or virtual) in volumetric displays is visible from every angle and hits no occluders.

Now we will derive the ambiguity function corresponding to a scene consisting of a plane of emitters at depth $z = \hat{z} < 0$ with intensity pattern $I(x, y)$ and a occlusion mask at the $z = 0$ plane with transmittance function $m(x, y)$.

Ignoring numerical effects, the light from the plane of emitters has the following ambiguity function:

$$A_0(f_x, f_y, \xi, \eta) = \tilde{I}(f_x, f_y)\delta(\xi + \hat{z}\lambda f_x, \eta + \hat{z}\lambda f_y) \tag{3.81}$$

Figure 3.11: The action of an occluder at the $z = 0$ plane on light from a set of planar emitters at $z = \hat{z} < 0$ is shown. The original ambiguity function is illustrated in (i) and (ii) in a fashion similar to Figure 3.10. The action of the occluder is to cause a convolution along the frequency axes by a kernel (iii) and (iv), and the resulting ambiguity function is illustrated in (v) and (vi). Note that in (vi), the output ambiguity function is no longer constrained to a narrow linear region as in Figure 3.10 (ii).

Incorporating a mask involves a convolution along the frequency axes with the ambiguity function $A_m(f_x, f_y, \xi, \eta)$ of the transmittance function:

$$
\begin{aligned}
A_1(f_x, f_y, \xi, \eta) &= \iint A_m(f_x - f_x', f_y - f_y', \xi, \eta) A_0(f_x', f_y', \xi, \eta) df_x' df_y' \\
&= \iint A_m(f_x - f_x', f_y - f_y', \xi, \eta) \tilde{I}(f_x', f_y') \delta(\xi + \hat{z}\lambda f_x', \eta + \hat{z}\lambda f_y') df_x' df_y'
\end{aligned}
$$

Performing the following substitution:

$$
\xi' = \hat{z}\lambda f_y' \ , \ \eta' = \hat{z}\lambda f_x' \tag{3.83}
$$

will yield:

$$
\begin{aligned}
&A_1(f_x, f_y, \xi, \eta) \\
&= \frac{1}{\hat{z}^2 \lambda^2} \iint A_m(f_x - \xi'/(\hat{z}\lambda), f_y - \eta'/(\hat{z}\lambda), \xi, \eta) \tilde{I}(\xi'/(\hat{z}\lambda), \eta'/(\hat{z}\lambda)) \delta(\xi + \xi', \eta + \eta') d\xi' d\eta' \\
&= \frac{1}{\hat{z}^2 \lambda^2} A_m(f_x + \xi/(\hat{z}\lambda), f_y + \eta/(\hat{z}\lambda), \xi, \eta) \tilde{I}(-\xi/(\hat{z}\lambda), -\eta/(\hat{z}\lambda)) \tag{3.84}
\end{aligned}
$$

Note that now the above expression, compared to Equation (3.78), no longer has a delta function term. Therefore, this ambiguity function is now actually a four dimensional function and hence the empty space in 3.10 must be filled. A summary of the operation can be seen pictorially in Figure 3.11. □

## 3.4 Summary

We can see from the analysis of the three different families of illumination systems that each type has its advantages and disadvantages. The ray-based systems use an intuitive model of light rays passing through space, but they suffer from resolution issues – high resolution areas to the "left" and "right" of the origin are diminished in Figure 3.3. Holographic systems have very good resolution performance, but coherent waves are harder to understand and also have certain patterns they cannot produce as well, illustrated by the region in $\mathcal{I}_\mathcal{P}$ that is not contained in $\mathcal{I}_\mathcal{C}$ in Figure 3.9. Lastly,

volumetric systems have probably the simplest representation of light, but they suffer from defocus blur and inability to represent astigmatic or occlusion effects, as shown by the circled areas in Figure 3.10. Since we've been comparing the types of light patterns that can be generated to the full set of possible mutual intensity patterns (or equivalently, Wigner distribution functions and ambiguity functions), perhaps it would be best to directly aim for the generation of a specific mutual intensity function.This will be the subject of the next chapter.

# Chapter 4

# Synthesis

We will now investigate how direct generation of specific partially coherent fields can be used to solve illumination problems. Conceptually, we can divide the task into two parts:

1. how to physically generate a desired partially coherent field (i.e. specified by its mutual intensity), and

2. how to design a specific partially coherent field (i.e. how to find a desired mutual intensity) to satisfy a given illumination application.

Recall from our review of partial coherence in Chapter 2 that a partially coherent field can be decomposed into an incoherent mixture of fully coherent fields. A straightforward method to generate a desired partially coherent field would be to produce separate coherent fields, each using its own laser and spatial light modulator (SLM), and then combine them using a system of beam splitters. This approach works for a few modes, but it obviously does not scale very well for partially coherent fields containing a large number of modes. As an alternative, De Santis et al. proposed generating a rapid sequence of coherent fields to create a desired partially coherent field [61]. For sink systems with long integration times, illumination generated this way would impact the sink system in the same way as a "true" partially coherent field. Let us use this approach for the rest of the chapter.

Figure 4.1: The optical setup consists of a coherent plane wave propagating from left to right being modulated by a time-varying phase-amplitude SLM, which is then in turn imaged by a 4-f system that selects the zeroth order onto the plane $\Pi_0$, where we seek to control the mutual intensity of the generated partially coherent light beam.

We will first specify the optical setup needed for this approach while reviewing methods to obtain control over amplitude and phase in a computer generated hologram. We will then consider algorithms for computing a desired mutual intensity given two example problems – simulating a "real" scene and computing a partially coherent beam with desired intensity distribution. Lastly, we will discuss intricacies involved with deriving a temporal sequence of modes from a desired mutual intensity as well as the generation of highly incoherent fields.

## 4.1 Generating a partially coherent field

For the generation of the time-multiplexed field, let us use a system akin to one in Figure 4.1. We will start with a plane wave incident on a SLM with control over the phase and amplitude of each pixel over time. There are $N$-by-$N$ pixels on the

Figure 4.2: Different SLM configurations to enable control over both phase and amplitude. Method (i) images a phase SLM onto an amplitude SLM. Method (ii) uses spatial filtering of a binary amplitude SLM to obtain a lower resolution amplitude-phase output. Method (iii) uses a Michelson interferometer to combine two phase functions coherently to produce arbitrary amplitude-phase patterns. Method (iv) uses spatial filtering to combine two phase functions multiplexed on a single SLM to produce lower resolution amplitude-phase output.

SLM plane spaced $\Delta_{SLM}$ wavelengths apart along each axis. The SLM plane is then imaged onto plane $\Pi_0$ by an ideal 4-f optical setup with a square aperture such that only the zeroth order diffraction pattern is retained, removing the effect of SLM pixel shape on the output light. In other words, the 4-f system is a 1-to-1 ideal relay system and the square aperture is designed such that the image of the $N$-by-$N$ points of the SLM critically sample the resulting wave function at any output plane at any instant in time. As a consequence of this sampling, discretization of fields conveniently has the same sampling pattern as the SLM pixel lattice. We seek to control the mutual intensity of the resulting partially coherent field at $\Pi_0$.

In general, spatial light modulators only modulate either the phase or the amplitude of an incoming coherent field, not both. However, one can build an amplitude-phase modulator using several different configurations, as shown in Figure 4.2. The straight-forward method (i) would be to image the output of an amplitude SLM onto a phase SLM. Another method (ii) would be to use a 4-f filter to select a specific diffraction order of a binary phase SLM [62]. Furthermore, since any complex valued function $f(x) = a(x)e^{j\phi(x)}$ can be written as the sum of two phase functions:

$$a(x)e^{j\phi(x)} = Ae^{j\phi_1(x)} + Ae^{j\phi_2(x)} \tag{4.1}$$

where $A = \max_x a(x)$ and

$$\phi_1(x) = \phi(x) + \cos^{-1}(a(x)/A)/2 \tag{4.2}$$

$$\phi_2(x) = \phi(x) - \cos^{-1}(a(x)/A)/2 \tag{4.3}$$

we can combine these two phase functions coherently using dual-phase approaches. In (iii), output from two phase SLMs are combined coherently through the use of interferometric methods [63]. In (iv), the two phase functions are multiplexed onto a single SLM, and they are mixed by spatial filtering [64].

Now that we have established methods for generating custom partially coherent patterns through temporal multiplexing of some coherence mode representation of a desired mutual intensity, let us explore ways to compute a desired mutual intensity given a specific application.

## 4.2 Scene simulation algorithm

We can simulate simple scenes and compute the output mutual intensity at some plane as long as these scenes can be modeled as light propagating in one direction from/through a sequence of transverse planes in a rectangular "tunnel", where each

plane is either a two-dimensional light emitting pattern of finite extent or a two-dimensional complex modulation function of finite support. The former allows incoherent light sources and the latter allows for both occlusions/absorption (amplitude) and lensing (phase). Any light that hits the edge of the tunnel is assumed to be discarded, and no reflections are allowed.

The method to compute the desired mutual intensity from such a scene is straightforward, although it can be very computationally heavy. The theoretical idea has been explored by Gross [65], and the modes representation optimization has been explored by Rydberg and Bentsson [66]. However, it would be useful here to specify the entire algorithm directly and visit some practical implementation details.

Let planes $\Pi_i$ be a series of transverse planes of interest with longitudinal coordinates $z_i, i = 1, 2, \ldots$, where $z = 0$ is the plane at which the mutual intensity is desired. Let us also assume without loss of generality that the $\Pi_i$ are sorted in increasing $z_i$. For each plane $\Pi_i$, let there be a possibly empty set of point emitters $\{P_k^{(i)}\}$ as well as a band-limited transmission function $T_i(x, y)$. That is, $T_i(x, y)$ can be sufficiently sampled at intervals of $\Delta_{SLM}$ into a discrete representation $T_i[m, n]$. Without loss of generality, let us assume light from point emitters at a particular plane is modulated by the transmission function at the same plane.

The overview of the process is as follows (please refer back to section 2.3.1 for details on discretization):

1. Initiate the current value of the (matrix form) mutual intensity $J \in \mathbb{C}^{N^2 \times N^2}$ to be all zeros and let $i = 0$.

2. For each point emitter $P_k^{(i)}$ at plane $\Pi_i$, compute its corresponding (coherent) mutual intensity pattern, $J^{(i,k)}$, and add it to the current value of the mutual intensity $J$.

3. Let $\mathbf{t}_i$ be the vector form of the transmission function $T_i[m, n]$ at plane $\Pi_i$ and element-wise multiply the current mutual intensity $J$ by the matrix $\mathbf{t}_i \mathbf{t}_i^H$.

4. If there are more planes to be considered, let $P_i$ be the matrix representing linear propagation from plane $\Pi_i$ to $\Pi_{i+1}$, set $J = P_i J P_i^H$, increment $i$ and go

to step 2.

5. If there are no more planes to be considered, let $P_i$ be the matrix representing linear propagation from plane $\Pi_i$ to the $z = 0$ plane and then set $J = P_i J P_i^H$.

We'll now examine some of these steps in more detail.

## 4.2.1 Calculation of point emitter mutual intensity

Given a coordinate $x_k^{(i)}, y_k^{(i)}$ and intensity $I_k^{(i)}$ for point emitter $P_k^{(i)}$, we can calculate the discretized field corresponding to this source. Since we are removing all but the zeroth order diffraction pattern from the output light, this means that we are bandlimited such that the pixels on the SLM critically samples the field. Thus, the field resulting from a point emitter would result in a diffraction spot instead, ie. a sinc function along each axis. For our particular point emitter, the resulting field would be:

$$U_k^{(i)}(x,y) = \frac{\sin(\pi(x - x_k^{(i)})) \sin(\pi(y - y_k^{(i)}))}{\pi^2 (x - x_k^{(i)})^2 (y - y_k^{(i)})^2} \tag{4.4}$$

Let the $N^2$ length vector $\mathbf{u}_k^{(i)}$ be the vector form of the section of the above field that lies inside the transverse boundaries of the SLM (i.e. the "tunnel"). The matrix form mutual intensity for this coherent field is an outer product of the vector with itself:

$$J^{(i,k)} = \mathbf{u}_k^{(i)} \mathbf{u}_k^{(i)H} \tag{4.5}$$

## 4.2.2 Calculation of mask mutual intensity

If we consider the interaction of a coherent field with a mask, it is an element-wise multiplication. Since a partially coherent field can be thought of as the incoherent sum of coherent fields, we can apply a masking operation to each mode of the partially coherent field to obtain the desired result.

Let $J$ be the incoming matrix form mutual intensity. Its coherent mode decomposition can be written as the following matrix factorization:

$$J_{in} = UU^H \tag{4.6}$$

Let $\mathbf{t}_i$ be the vector form representation of the mask. To apply a masking operation, we would need to element-wise multiply each column of $U$, i.e. each mode of $J$, by $\mathbf{t}_i$. Let $D_{t_i}$ be a diagonal matrix whose entries are the elements of $\mathbf{t}_i$. Then, application of this mask to a single mode $\mathbf{u}_{in}$ would be:

$$\mathbf{u}_{out} = D_{t_i}\mathbf{u}_{in} \tag{4.7}$$

Hence, applying this to the coherent mode decomposition of incoming mutual intensity results in

$$J_{out} = D_{t_i}UU^H D_{t_i}^H \tag{4.8}$$

This is equivalent to:

$$J_{out} = (\mathbf{t}_i\mathbf{t}_i^H) \odot J_{in} \tag{4.9}$$

where $\odot$ is element-wise multiplication of two matrices. Hence, we can think of $\mathbf{t}_i\mathbf{t}_i^H$ as the mask "mutual intensity".

### 4.2.3 Linear propagation of mutual intensity

Paraxial propagation of a transverse field is a linear shift-invariant operation. The straightforward way to compute propagation would then to be to pad in 2D, perform a 2D Fourier transform, element-wise multiply by a discretized phase function (either an angular spectrum representation of a Fresnel diffraction transfer function [47]), perform an inverse 2D Fourier transform and then crop in 2D. Thus, for a single coherent field, the action of propagation by $\Delta z$ can be written as:

$$\mathbf{u}_{out} = A_{crop}F^{-1}D_{\Delta z}FA_{pad}\mathbf{u}_{in} \tag{4.10}$$

where $A_{pad}$ is the padding operator, $A_{crop}$ is the cropping operator, $F$ is the forward fast Fourier transform operator and $D_{\Delta z}$ is the propagation transfer function. Since the mutual intensity can be represented as an outer product of mode vectors:

$$J = UU^H \tag{4.11}$$

we can apply the operation in Equation (4.10) to each column in the above equation to obtain an expression for the action of propagation on the mutual intensity:

$$
\begin{aligned}
J_{out} &= A_{crop}F^{-1}D_{\Delta z}FA_{pad}U(A_{crop}F^{-1}D_{\Delta z}FA_{pad}U)^H \\
&= A_{crop}F^{-1}D_{\Delta z}FA_{pad}UU^H(A_{crop}F^{-1}D_{\Delta z}FA_{pad})^H \quad (4.12)
\end{aligned}
$$

In practice, we can perform this propagation operation by:

1. Perform propagation by treating all the columns of the input mutual intensity matrix as separate fields.

2. Perform a conjugate transpose.

3. Perform propagation by treating all the columns of the resultant matrix as separate fields.

4. Perform a conjugate transpose.

The padding and cropping operations may require fine-tuning to avoid inefficient use of memory, so let's discuss how to compute the desired padding amount.

Since we are performing a discrete Fourier transform instead of an actual continuous Fourier transform, the operation models the input and output as periodic functions where one period consists of the actual values in the discretized representation. Thus, any expansion of the beam, i.e. diffraction, will cause aliasing. That is, light diffraction off the left edge of the SLM boundary would actually then show up coming in from the right side. Thus, padding is needed to prevent incorrect simulation of propagation.

Given the propagation distance $\Delta z$ and the sampling interval (which determines the bandlimit) $\Delta_{SLM}$, we can calculate the approximate maximum angle of diffraction of the zeroth order, which is the light we would retain. If we sample at intervals of $\Delta_{SLM}$, then the maximum frequency along each axis that can be represented according to sampling theory is:

$$
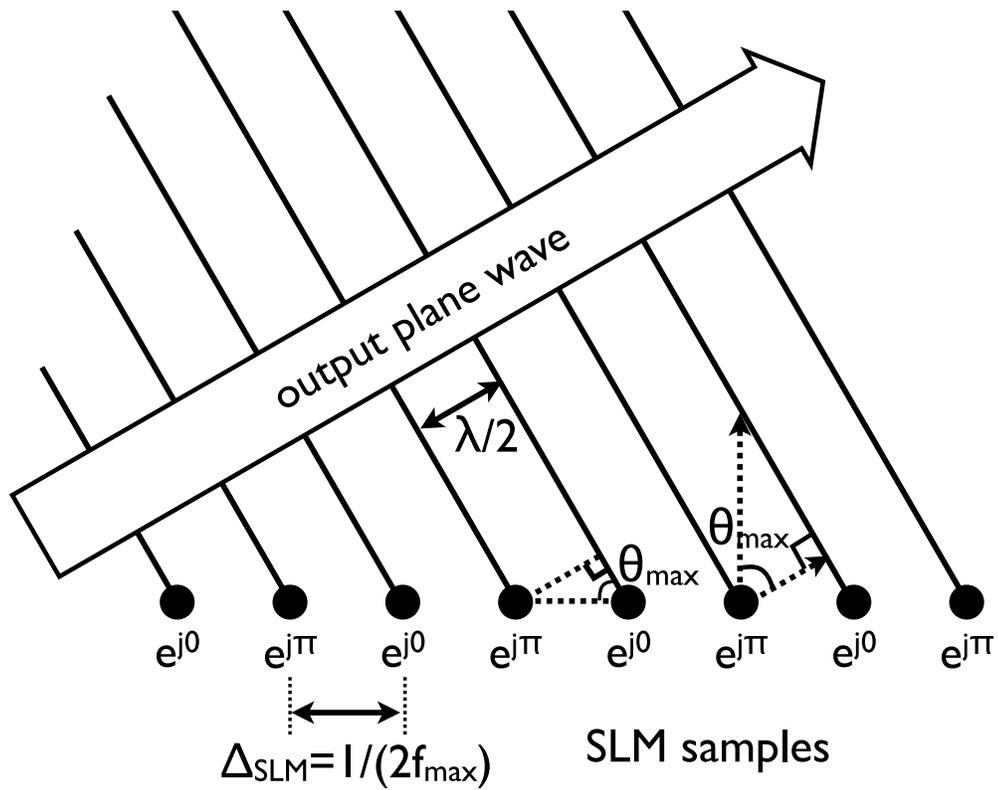f_{max} = \frac{1}{2\Delta_{SLM}} \quad (4.13)
$$

Figure 4.3: The maximum angle plane wave that can be generated from a set of SLM pixels which critically sample the outgoing field is related to the sampling rate by a simple trigonometric relationship. The circles represent SLM pixel centers and the lines in the plane wave are at intervals of $\pi$ phase.

 As shown in Figure 4.3, a plane wave whose transverse phase pattern corresponds to a complex phasor at frequency $f_{max}$ would be at an angle $\theta_{max}$ determined by:

$$\sin(\theta_{max}) = \frac{\lambda}{\Delta_{SLM}} \tag{4.14}$$

Thus, from this maximum angle, we would need to pad transversely along each dimension by:

$$\Delta_{pad} = 2\tan(\theta_{max})\Delta z \tag{4.15}$$

In the discrete domain, we can simply extend with zeros along each dimension $m$ and $n$ for

$$\lceil \tan(\theta_{max})\Delta z/\Delta_{SLM} \rceil \tag{4.16}$$

samples in each direction. In practice, padding by a multiple of the amount generally results in a slightly better results, since this rule of thumb approximates diffraction spread by a ray of the maximum angle plane wave from the edge of the SLM, whereas theoretically an entire plane wave at that angle exists, and not just inside the SLM.

After propagation, the output field/mutual intensity needs to be trimmed down to remove the extra padded samples so that the result fits within the bounds of the SLM again.

### 4.2.4   Scaling and mode representation

The storage and computation for this algorithm can become exceedingly large. Storage of a matrix form mutual intensity pattern corresponding to a $N \times N$ SLM is $O(N^4)$, and the propagation computation is $O(N^4 \log N)$ if we use the fast Fourier transform. While this does not scale very well, some optimizations can be done to partially mitigate the issues.

If there are less than $N^2$ point emitters, then storing the entire mutual intensity matrix is wasting storage, as it cannot be a full-rank matrix in that case. Recall that in the extreme case of a single emitter, we have what is essentially a coherent field and thus a rank-one matrix. An easy alternative would be to retain the modes representation of the mutual intensity instead of the actual matrix. That is, instead of

storing the current mutual intensity $J$, we store its mode representation $U$ where $J = UU^H$. This way, $U$ has only number of columns equal to the number of point emitters. Furthermore, addition of new emitters simply means concatenating $\mathbf{u}_k^{(i)}$ to the current matrix $U$, masking involves multiplying each row in $U$ by the corresponding row in $\mathbf{t}_i$, and propagation requires just a single linear operator instead of two in the case of directly storing and using the mutual intensity. When the number of columns of $U$ exceeds the number of rows, we can perform a singular value decomposition and find an orthogonal representation $\hat{U}$ such that $\hat{U}\hat{U}^H = UU^H$. Furthermore, we can choose to remove modes (columns of $\hat{U}$) that contribute very little energy, as an approximating optimization.

## 4.2.5 Results

As an example, we will now use the described algorithm to compute the mutual intensity corresponding to light that is emitted from a two-dimensional transverse pattern which then hits a two-dimensional occluder. The SLM will be an array of $32 \times 32$ pixels with $20\lambda$ pitch. The emitter pattern will be specified by a $63 \times 63$ pixel lattice with $10\lambda$ pitch. At $z = -3200\lambda$, there will be a plane of emitters with pattern shown by (i) in Figure 4.4. Light then propagates in the positive $z$ direction until $z = 0$, where there will be an occluder pattern shown by (ii) in Figure 4.4.



(i)                    (ii)

Figure 4.4: As shown in (i), a group of emitters in the shape of four copies of the letter "B" is located at $z = -3200\lambda$. White pixels indicate maximum brightness. There were a total of 1268 non-black pixels in this image and 728 of them were at full intensity. At the $z = 0$ plane, there is an occluder in the shape of the letter "A" that blocks any light that falls on a line of the character, as shown in (ii). Black pixels indicate maximum absorption.

Running the described algorithm resulted in a $1024 \times 1024$ mutual intensity matrix. To verify the correct generation of the mutual intensity matrix, the computed output mutual intensity of the occluded light was refocused to obtain the intensity at various transverse planes between the $z = -3200\lambda$ plane and the $z = 0$ plane, and this "focal stack" is shown in Figure 4.5 (i). Furthermore, the intensity image at the $z = 0$ plane was computed after selecting only a square region of the Fourier plane and shifting the square region through 5 positions. This would be as if we passed the light through a 4-f system with a square aperture a quarter the size of the full aperture and shifted the aperture through 5 positions. The goal is to produce "oblique" views of the scene from different directions and this is essentially the aperture scan device discussed in section 3.1.2 run in reverse. The resulting images are shown in Figure 4.5 (ii).



(i)



(ii)

Figure 4.5: A focal stack computed from the output mutual intensity matrix is shown in (i). A tilt-view image sequence generated using a square "pinhole" in the Fourier plane is shown in (ii).

A plot of the square of the singular values of the mutual intensity matrix (i.e. the amount of energy present in each coherence mode) is shown in Figure 4.6, along with a plot of cumulative energy. From the graphs, it is apparent that very negligible energy is present beyond mode number 500. The first 256 modes will capture 89.2% of the energy present in the original mutual intensity matrix. Note that a fully lit

emitter plane without occlusions would be fully incoherent and yield the maximum number of 1024 modes, and a single emitter would yield a coherent field (i.e. one mode). The emitter plane in this case only has 1268 lit pixels out of a maximum possible $63 \times 63 = 3969$ pixels, giving a rough estimate of $1268/3969 * 1024 \approx 327$ modes, agreeing with observations.



Figure 4.6: A plot of energy contained within each coherence mode of the final computed mutual intensity matrix is shown on the left graph. The cumulative energy is shown on the right graph. Energy corresponding to a mode is the square of the corresponding singular value in a singular value decomposition of the matrix.

These observations lead to the feasibility of low rank approximations of the mutual intensity matrix, which uses less modes and thus less patterns to display in rapid succession on an SLM. Since the modes are calculated from a singular value decomposition, optimally reducing the number of modes simply means choosing a subset of

the modes corresponding to the highest energies. In Figure 4.7, we show the effect of reducing modes on the focal stack and in Figure 4.8, we show the effect of reducing modes on the tilt-view images.



Figure 4.7: The resulting focal stack as a function of reducing the number of modes in the mutual intensity via the singular value decomposition. $K$ for each row denotes the number of modes kept, where $K = 1024$ corresponds to keeping all the modes.

From these images, it is apparent that 256 modes is sufficient to achieve nearly indistinguishable results from the full modes case. However, in practice, 256 modes is still a lot of patterns to cycle through on a spatial light modulator, so in section 4.4.2, we will investigate a different approach to generating fairly incoherent fields such as this one.
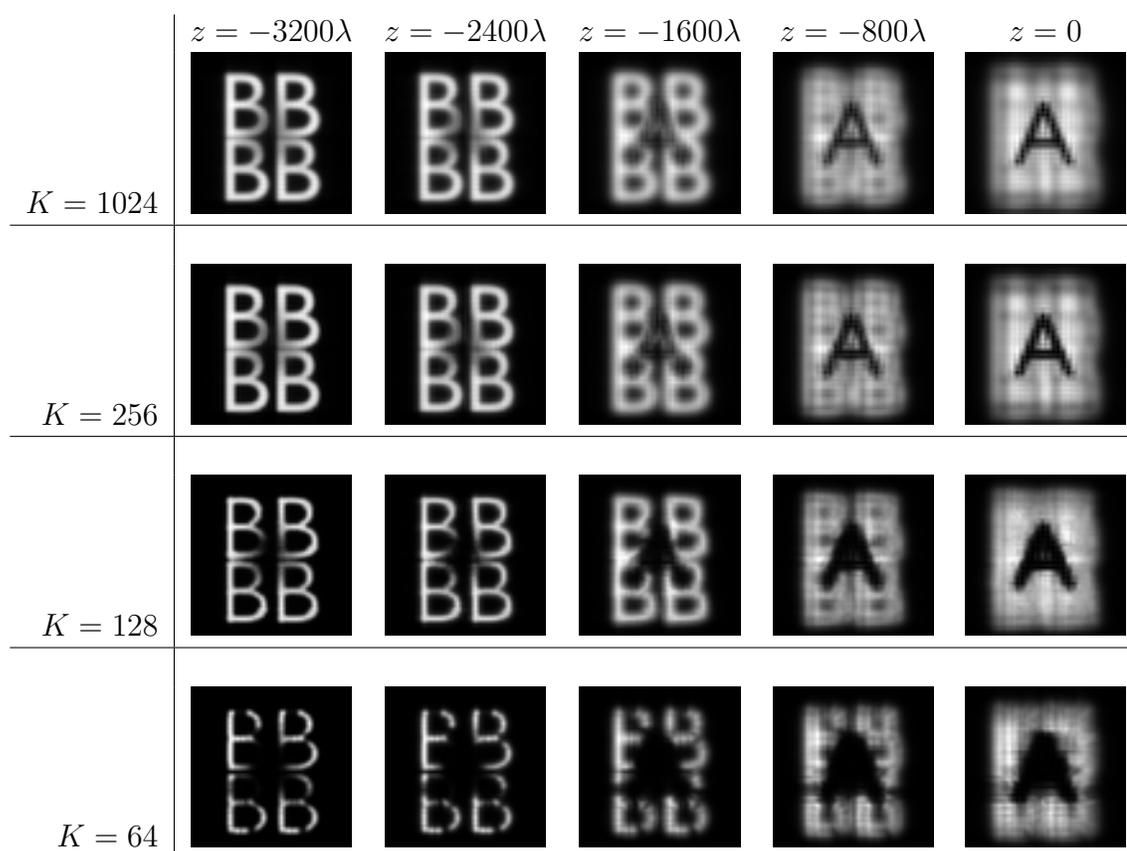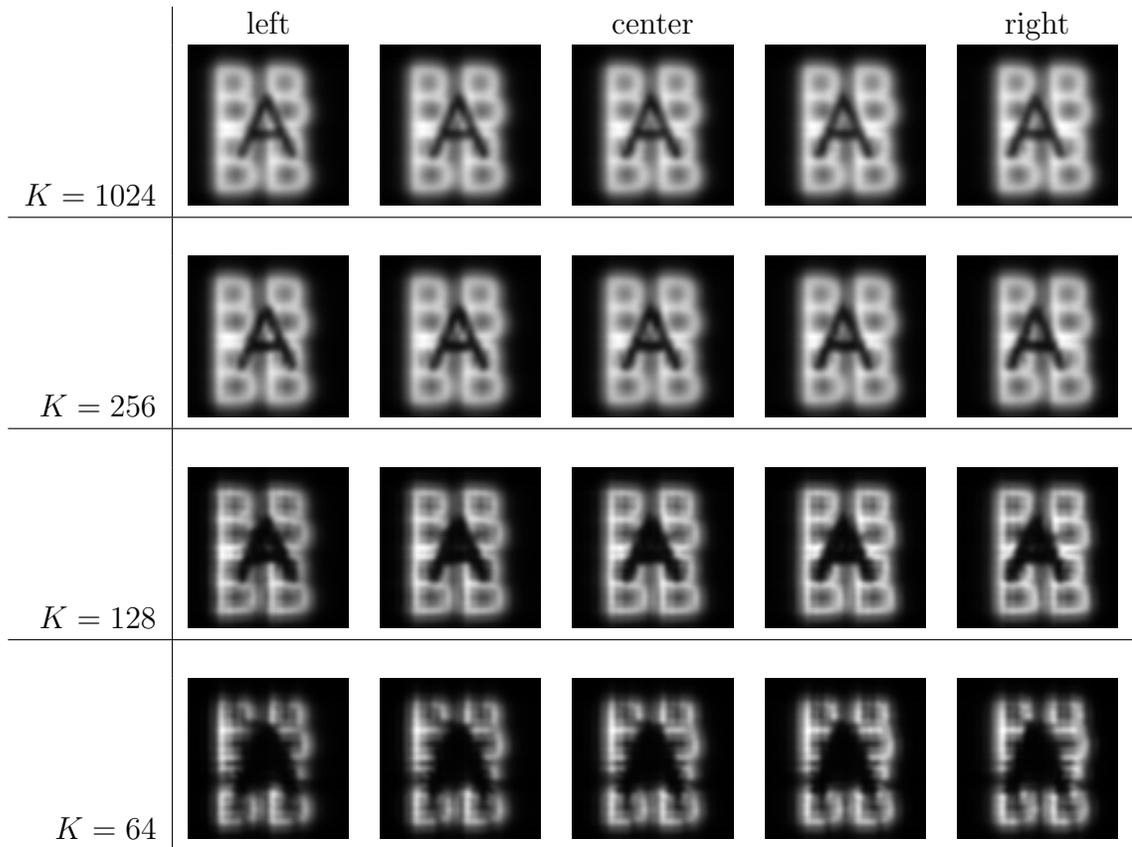
Figure 4.8: The resulting set of tilt-view images as a function of reducing the number of modes in the mutual intensity via the singular value decomposition. $K$ for each row denotes the number of modes kept, where $K = 1024$ corresponds to keeping all the modes.

## 4.3   Voxel intensity algorithm

An alternate way to specify a desired mutual intensity is to require one which has a specific three-dimensional intensity distribution. That is, we would like to specify the light intensity at every voxel in space. Keep in mind that although this is similar to how volumetric displays specify their illumination, it is different in that we are specifying the final intensity that includes effects such as out-of-focus blur, instead of just the emission intensity. Thus, using this technique, it is possible for us to specify intensity patterns that have alternating bright and dark patterns along linear path(s) parallel to the optical axis, whereas attempts to reproduce this using volumetric displays will result in stray light in "black" areas due to out-of-focus blur.

Now let us examine the problem setup in more detail, as shown in Figure 4.9. Compared to Figure 4.1, we have added $M$ transverse planes $\Pi_m$ located at $z = z_m$ at which we would like to control the two-dimensional intensity pattern of light. Since the light is bandlimited, we can sample the intensity patterns at intervals of $\Delta_{SLM}/2$. Together, these sampled planes form a volume composed of voxels, one for each sample. In addition to the specification of the voxels, we will also specify a limit on the number of modes allowed in the creation of a mutual intensity to satisfy the desired voxel intensities. That is, we would like to pick a number $K$ and then compute a mutual intensity $J \in \mathbb{C}^{N^2 \times N^2}$ such that $J$ can be written as $J = UU^H$ where $U \in \mathbb{C}^{N^2 \times K}$ and that the resulting intensities at the desired planes $\Pi_m$ matches the desired voxel intensities. For this computation, it is more efficient to operate on the mode decomposition form of the mutual intensity rather than the full matrix, due to the matrix being never full rank. Therefore, the problem will be formulated as an optimization to derive a set of coherence modes instead of the mutual intensity directly.

Let $\mathbf{u}_k \in \mathbb{C}^{N^2}, k = 1, \ldots, K$ be vectors that form a $K$-mode representation of the mutual intensity of the partially coherent beam. That is, each $\mathbf{u}_k$ is a vectorized form of a discretized field. Let $\mathbf{u} \in \mathbb{C}^{N^2 K}$ be a column vector consisting of the $\mathbf{u}_k$ stacked vertically. Let us also define $\mathbf{y}_m \in \mathbb{R}^{(2N-1)(2N-1)}, m = 1, \ldots, M$ to be the vector form of the intensity at each of the $M$ planes. Along each dimension, samples between

Figure 4.9: The optical setup consists of a coherent plane wave from the left being modulated by a time-varying phase-amplitude SLM, which is then in turn imaged by a 4-f system that removes all but the zeroth diffraction order onto the plane $\Pi_0$. Instead of seeking to just control the mutual intensity at plane $\Pi_0$ like in Figure 4.1, we seek to control the time-averaged intensity at planes $\Pi_1, \Pi_2, \ldots, \Pi_M$ by designing patterns for the SLM.

the original $N$ samples have been included due to the higher sampling rate of the intensity, and thus we have a total of $2N - 1$ samples along each dimension. Denote the $n^{th}$ element of $\mathbf{y}_m$ as $y_{m,n}$.

Furthermore, let us define $P_m \in \mathbb{C}^{(2N-1)(2N-1) \times N^2}$ to be a linear operator (matrix) that propagates light from $\Pi_0$ to $\Pi_m$. That is, each row $\mathbf{p}_{m,i}^H$ in $P_m$ governs how light from all $N^2$ points in a sampled $\Pi_0$ propagates to one of the $(2N-1)(2N-1)$ sample points in $\Pi_m$ that correspond to a sample of the intensity. The operation $P_m$ can be modeled as the composition of the following sequence of operations:

1. padding the input to prevent aliasing after propagation (amount of padding is calculated from how much light spreads during propagation given a NA; alternatively, how many samples are needed in Fourier space to correctly sample the phase delay function)

2. Fourier transform

3. element-wise phase delay according to the propagation distance $z_m$ (this is calculated using the angular spectrum propagation method [47])

4. padding by a factor of two to arrive at the correct sampling rate for the intensity

5. inverse Fourier transform

6. cropping to extract the desired intensity samples

This is summarized in Fig. 4.10.

With the illumination problem reformulated using this discretization, we ideally want the sum of the intensities across time at each sample point to be equal to our desired intensity:

$$y_{m,n} = \sum_{k=1}^{K} \left| \mathbf{p}_{m,n}^H \mathbf{u}_k \right|^2 \tag{4.17}$$

where $m$ indicates which plane this particular intensity sample is located and $n$ indicates the position within that plane.

However, not all intensity patterns are possible, due to limits of physics and also due to the constraint on the number of modes we've placed. Therefore, let us instead

Figure 4.10: Propagation of the sampled SLM amplitude-phase values to samples at an output plane $\Pi_m$ can be calculated easily through a series of linear operations.

try to minimize the squared error (with optional weighting) between the intensity we obtain and the intensity we desire:

$$A(\mathbf{u}) = \sum_{m=1}^{M} \sum_{n=1}^{(2N-1)(2N-1)} w_{m,n}^2 \left( y_{m,n} - \sum_{k=1}^{K} \left| \mathbf{p}_{m,n}^H \mathbf{u}_k \right|^2 \right)^2 \qquad (4.18)$$

where $A(\mathbf{u})$ is a function from $\mathbb{C}^{N^2 K}$ to $\mathbb{R}$ that gives the total weighted square error and $w_{m,n}$ is the weighting factor for the $n^{th}$ pixel on plane $\Pi_m$. Recall that $\mathbf{u}$ is the vertical concatenation of the $K$ vectors $\mathbf{u}_k$. This function is a multivariate non-convex quartic function, and thus it has no closed form solution. An iterative optimization approach will now be derived to minimize it.

## 4.3.1 Optimization method

There are many methods for optimizing general nonlinear non-convex functions, e.g. Newton's method, conjugate gradients or steepest descent. Steepest descent requires

only computation of the gradient, but does not converge as nicely as nonlinear conjugate gradient methods, which only require slightly more computation. Newton's method usually converges the best once inside a positive-definite quadratic area, but having to solve a system of linear equations involving a large Hessian and the fact that the Hessian is not always positive definite makes this method not as desirable due to excessive computational cost. We will take a closer look at steepest descent and Newton's methods later on in the section.

For this optimization, we will use a variant [67] of the Polak-Ribière method [68] of nonlinear conjugate gradients coupled with global line search. As summarized in Fig. 4.11, each iteration of this optimization algorithm will consist of the following steps, which will be explained in the next three subsections:

1. Compute the direction of steepest descent $\mathbf{\Delta} u$ for the merit function $A(\mathbf{u})$.

2. Compute the conjugate gradient step direction $\mathbf{\Lambda} u$ using the modified Polak-Ribiére formula.

3. Compute the step length by finding the global minimum along the line dictated by the conjugate gradient step direction.

**Steepest descent direction**

In order to use conjugate gradients, we need to first compute the direction of steepest descent, i.e. the gradient of our merit function, $A(\mathbf{u})$. Since differentiation is a linear operation and our merit function is a sum across all pixels, the gradient can be computed by finding derivatives corresponding to each pixel first and then summing them afterwards.

For pixel $n$ on plane $\Pi_m$, the portion of the merit function in question is:

$$A_{m,n}(\mathbf{u}) = w_{m,n}^2 \left( y_{m,n} - \sum_{k=1}^{K} \left| \mathbf{p}_{m,n}^H \mathbf{u}_k \right|^2 \right)^2 \tag{4.19}$$

For small changes $\mathbf{\Delta} u$ in $\mathbf{u}$, we can approximate the above equation with a linear

$\mathbf{u}_{(0)}$ (KxNxN random noise pattern in the complex domain)

first iteration

previous iteration

$\Delta\mathbf{u}_{(i-1)}$ $\quad$ $\Lambda\mathbf{u}_{(i-1)}$

$\mathbf{u}_{(i)}$

| compute direction of steepest descent | $\Delta\mathbf{u}_{(i)}$ | compute conjugate gradient direction | $\Lambda\mathbf{u}_{(i)}$ | perform line search | $\alpha_{(i)}$ |

$+$

$\times$

$\mathbf{u}_{(i+1)}$

next iteration

Figure 4.11: The optimization is initialized, usually with a random set of SLM patterns. Then, for each iteration, we compute the direction of steepest descent, and then compute the conjugate gradient step direction using the current and previous steepest descent directions as well as the previous conjugate gradient step direction. Lastly, a line search is performed to find the global minimum along the conjugate gradient step direction and the current iterate is updated.

expansion about $\mathbf{u}$:

$$
\begin{aligned}
\hat{A}_{m,n}(\boldsymbol{\Delta}u) &= w_{m,n}^2 \left( y_{m,n} - \sum_{k=1}^{K} \left| \mathbf{p}_{m,n}^H (\mathbf{u}_k + \boldsymbol{\Delta}u_k) \right|^2 \right)^2 \\
&= w_{m,n}^2 \left( \Delta y_{m,n} - \left( \sum_{k=1}^{K} \left| \mathbf{p}_{m,n}^H \boldsymbol{\Delta}u_k \right|^2 + 2\mathrm{Re}\left\{ \mathbf{u}_k^H \mathbf{p}_{m,n} \mathbf{p}_{m,n}^H \boldsymbol{\Delta}u_k \right\} \right) \right)^2 \\
&\approx w_{m,n}^2 \left( \Delta y_{m,n} - 2 \sum_{k=1}^{K} \mathrm{Re}\left\{ \mathbf{u}_k^H \mathbf{p}_{m,n} \mathbf{p}_{m,n}^H \boldsymbol{\Delta}u_k \right\} \right)^2 &\text{(4.20)} \\
&\approx w_{m,n}^2 \left( \Delta y_{m,n}^2 - 4\Delta y_{m,n} \sum_{k=1}^{K} \mathrm{Re}\left\{ \mathbf{u}_k^H \mathbf{p}_{m,n} \mathbf{p}_{m,n}^H \boldsymbol{\Delta}u_k \right\} \right) &\text{(4.21)}
\end{aligned}
$$

where $\Delta y_{m,n}$ denotes the current error in intensity and $\boldsymbol{\Delta}u_k$ corresponds to a small change in the SLM pattern for mode $k$. In order to understand the operator which takes the real part of the product in the equation, let us decompose the terms in the product into real and imaginary parts by setting:

$$
\boldsymbol{\Delta}u_k = \mathbf{u}_k^{(R)} + j\mathbf{u}_k^{(I)} \tag{4.22}
$$

and

$$
\mathbf{p}_{m,n}\mathbf{p}_{m,n}^H \mathbf{u}_k = \mathbf{a}_{m,n,k} + j\mathbf{b}_{m,n,k} \tag{4.23}
$$

We can then obtain that:

$$
\hat{A}_{m,n}(\boldsymbol{\Delta}u) \approx w_{m,n}^2 \left( \Delta y_{m,n}^2 - 4\Delta y_{m,n} \sum_{k=1}^{K} \mathbf{a}_{m,n,k}^T \mathbf{u}_k^{(R)} + \mathbf{b}_{m,n,k}^T \mathbf{u}_k^{(I)} \right) \tag{4.24}
$$

Hence, the gradient with respect to each $\mathbf{u}_k^{(R)}$ and $\mathbf{u}_k^{(I)}$ can be written as:

$$
\nabla_{\mathbf{u}_k^{(R)}} A = -4 \sum_{m=1}^{M} \sum_{n=1}^{(2N-1)(2N-1)} \Delta y_{m,n} w_{m,n}^2 \mathbf{a}_{m,n,k} \tag{4.25}
$$

$$
\nabla_{\mathbf{u}_k^{(I)}} A = -4 \sum_{m=1}^{M} \sum_{n=1}^{(2N-1)(2N-1)} \Delta y_{m,n} w_{m,n}^2 \mathbf{b}_{m,n,k} \tag{4.26}
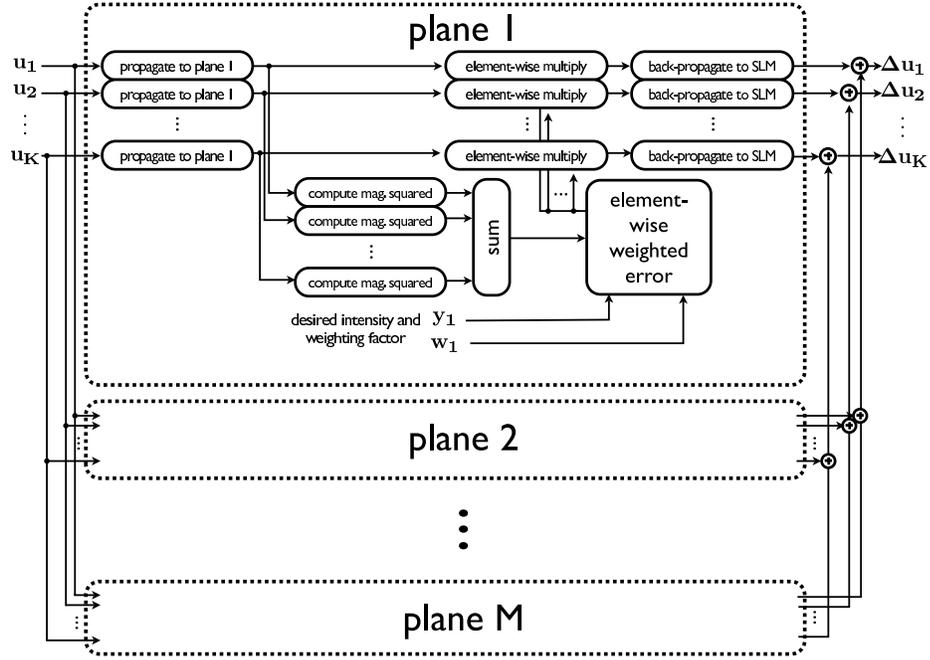$$

Figure 4.12: To find the steepest descent direction, the current SLM patterns $\mathbf{u}_k$ are propagated to one of the $M$ planes, and element-wise multiplied by the weighted error of the obtained intensity. The results are summed together for each of the $K$ patterns to produce the output direction of steepest descent $\Delta u_k$.

If we combined these constituents, we obtain that the direction of steepest descent in complex $\mathbf{u}$ is:

$$\Delta\mathbf{u}_k = 4\sum_{m=1}^{M}\sum_{n=1}^{(2N-1)(2N-1)} \Delta y_{m,n} w_{m,n}^2 \mathbf{p}_{m,n}\mathbf{p}_{m,n}^H \mathbf{u}_k \tag{4.27}$$

This can be computed easily in parallel by observing that for a particular mode $k$, the descent direction is formed by propagating the current SLM pattern to a target pixel $(\mathbf{p}_{m,n}^H\mathbf{u}_k)$, multiplying by the current difference in intensity at that pixel and its weight squared $(\Delta y_{m,n}w_{m,n}^2)$ and then back-propagating back to the the SLM pattern $\mathbf{p}_{m,n}$ and summing this result over all the pixels. Effectively, this can be cheaply computed by a forward propagation, element-wise multiplication by the weighted error, and then a back-propagation, as shown in Fig. 4.12.

## Conjugate gradients

Steepest descent without any line search usually falls very easily into local minima and can be inefficient in that the algorithm can "back-step" since alternating steps are not necessarily orthogonal. These issues are addressed in part by an adaptation of the conjugate gradients method of solving unconstrained quadratic problems to the general nonlinear optimization regime, resulting in a method called the nonlinear conjugate gradients method.

In classic nonlinear conjugate gradients, the new conjugate gradient direction $\mathbf{\Lambda}x$ at step $i$ is formed by adding a scale $\beta$ of the previous conjugate gradient direction to the current steepest descent direction $\mathbf{\Delta}x$:

$$\mathbf{\Lambda}x_{(i)} = \mathbf{\Delta}x_{(i)} + \beta_{(i)}\mathbf{\Lambda}x_{(i-1)} \tag{4.28}$$

where $\beta$ can be calculated using one of many formulas. If we choose the modified Polak-Ribiére formula, we obtain:

$$\beta_{(i)} = \max\left(0, \frac{\mathbf{\Delta}x_{(i)}^T(\mathbf{\Delta}x_{(i)} - \mathbf{\Delta}x_{(i-1)})}{\mathbf{\Delta}x_{(i-1)}^T\mathbf{\Delta}x_{(i-1)}}\right) \tag{4.29}$$

In order to apply this formula to our problem, we need to first convert our problem search space to a space of real numbers. This can be done by simply separating the real and imaginary components of $\mathbf{u}$ again like we did in the previous section. Hence, for the calculation:

$$\mathbf{\Lambda}u_{(i)} = \mathbf{\Delta}u_{(i)} + \beta_{(i)}\mathbf{\Lambda}u_{(i-1)} \tag{4.30}$$

where

$$\beta_{(i)} = \max\left(0, \frac{\mathbf{\Delta}\hat{u}_{(i)}^T(\mathbf{\Delta}\hat{u}_{(i)} - \mathbf{\Delta}\hat{u}_{(i-1)})}{\mathbf{\Delta}\hat{u}_{(i-1)}^T\mathbf{\Delta}\hat{u}_{(i-1)}}\right) \tag{4.31}$$

and $\mathbf{\Delta}\hat{u}_{(i)}$ consists of the real components of $\mathbf{\Delta}u_{(i)}$ concatenated vertically with the imaginary components of $\mathbf{\Delta}u_{(i)}$.

**Line search**

Now that we have computed a direction $\mathbf{\Lambda}u$, we can perform a search for the global minimum (corresponding to a step size of $\alpha_{(i)}$) along that line:

$$\alpha_{(i)} = \arg\min_{\alpha} A\left(\mathbf{u}_{(i)} + \alpha\mathbf{\Lambda}u_{(i)}\right) \tag{4.32}$$

Since $A(\mathbf{u})$ is a quartic in $\mathbf{u}$, making both $\mathbf{u}_{(i)}$ and $\mathbf{\Lambda}u_{(i)}$ constant in the above equation turns this minimization into a single variable ($\alpha$) quartic minimization problem:

$$F(\alpha) = A\left(\mathbf{u}_{(i)} + \alpha\mathbf{\Lambda}u_{(i)}\right) \tag{4.33}$$

The parameter $\alpha$ that corresponds to the global minimum of $A(\mathbf{u})$ along the line defined by the point $\mathbf{u}_{(i)}$ and direction $\mathbf{\Lambda}u_{(i)}$ can be easily determined via the following steps:

1. Determine the five explicit polynomial coefficients to the quartic $F(\alpha)$ explicitly.

2. Solve $\partial F(\alpha)/\partial\alpha = 0$ to find points where local minima may reside. This is a cubic and can be solved either through closed form or through any number of polynomial solvers.

3. Compute $F(\alpha)$ for these candidate points to find the point with the lowest value, and choose the $\alpha$ corresponding to this point as our minimizer.

Being able to derive a closed form expression for $F(\alpha)$ allows the optimization process to jump past local minima if there is a better minimum "visible" elsewhere along the same line. In that sense, we can avoid some local minima with this optimization scheme. However, if the optimization scheme does not ever cause this search line to cross the small region surrounding the global minimum, then this optimization scheme will never reach the global minimum.

**The Hessian and Newton's method**

One might wonder if computing the Hessian and then inverting it in a Newton's method optimization scheme would work better. In addition to the Hessian being

large and thus must be inverted through some sort of inner loop consisting of an iterative algorithm, the problem is that the Hessian is not necessarily positive definite and a straightforward constraint to at least make it positive semi-definite causes Newton's method style iterations to diverge.

If we look at the quadratic terms of our expansion of the merit function before we took the linear approximation in Eq. (4.20), we should have a quadratic term corresponding to $2\mathrm{Re}\left\{\mathbf{u}_k^H\mathbf{p}_{m,n}\mathbf{p}_{m,n}^H\boldsymbol{\Delta}u_k\right\}$ being squared added to a quadratic term corresponding to the $\left|\mathbf{p}_{m,n}^H\boldsymbol{\Delta}u_k\right|^2$ term being multiplied by the $\Delta y_{m,n}$ term. The former obviously yields a positive semi-definite quadratic form. However, the latter yields a negative semi-definite quadratic form due to the negative sign in front of the sum. Therefore, it is very possible for the Hessian to have negative eigenvalues, especially if the difference in intensity $\Delta y_{m,n}$ is very large (and positive-valued).

If we try to constrain the Hessian to be positive semi-definite by dropping the negative semi-definite form from the quadratic approximation, then we end up exactly with a minimization of the sum across all pixels of Eq. (4.20). This may look feasible to minimize, because it has now changed the quadratic programming problem for one iteration of Newton's method into a weighted linear least squares problem. However, this causes too much "interplay" between modes. To see this, we can rewrite Eq. (4.20) as:

$$\hat{A}_{m,n}(\boldsymbol{\Delta}\mathbf{u}_k) \approx w_{m,n}^2\left(y_{m,n} - \sum_{k=1}^{K}\mathrm{Re}\left\{\mathbf{u}_k^H\mathbf{p}_{m,n}\mathbf{p}_{m,n}^H(\mathbf{u}_k + 2\boldsymbol{\Delta}u_k)\right\}\right)^2 \qquad (4.34)$$

What this amounts to is that we are trying to construct the target intensity $y_{m,n}$ from a sum of contributions from the different modes. However, these contributions are different from intensities in that we are looking at the product of the complex conjugate of the propagated old value $\mathbf{u}_k^H\mathbf{p}_{m,n}$ with a propagated new value $\mathbf{p}_{m,n}^H(\mathbf{u}_k + 2\boldsymbol{\Delta}u_k)$. Since these values are generally not complex conjugates of each other, a minimum could be achieved when each contribution contains values that are either negative or have imaginary components, especially if $k > 1$. This is because the negative components from one mode could be cancelled out by overly positive elements

from another mode, and the imaginary components which do contribute to the non-approximated merit function are simply thrown out. This possibility of negative values and imaginary components causes the magnitude of $\mathbf{u}_k$ in the iterations to unpredictably and unnecessarily increase, resulting in either the iteration of linear least squares to diverge if no subsequent line search is performed, or getting the optimization stuck when the gradient is positive, due to the resultant search direction being not a descent direction.

One might also try taking the entire matrix and directly dropping the negative eigenvalues. This may be unfeasible due to the size of the Hessian, since it has rows and columns equal to the total number of pixels on the SLM multiplied by the number of modes.

These issues and the fact that the Hessian is a very large matrix to invert makes nonlinear conjugate gradients very appealing, as it is an improvement on steepest descent with little additional computational cost.

**Projection-based methods**

At this point, it is appropriate to take another look at iterated projection on constraint set methods employed in the field of holography for determining the coherent wave function (or the partially coherent mutual intensity) from a series of intensity measurements. Piestun and Shamir summarize these projection-based coherent approaches and present a parallelized block projection approach for the synthesis of a coherent beam from intensity constraints [69].

In essence, the given solution has to satisfy a set of constraints, each set being a set of intensity values at some transverse plane along the beam. Hence, to satisfy each constraint, the solution wave is propagated to that plane and its amplitude is set equal to the square root of the desired intensity (i.e. projected onto the set of possible wave functions that satisfy the desired intensity). Usually, there is more than one constraint, and thus the projections can either be done in serial or in parallel. When done in parallel, the next step of the iteration is a weighted average of the individual projections of the current solution wave.

For the partially coherent case, Rydberg and Bentsson present a serial algorithm

[70] wherein modes of the partially coherent beam are propagated to one transverse plane at a time wherein all the modes are multiplied by the same amplitude weighting function in order to satisfy the intensity constraint. The action at each plane is a projection, since the minimal change required of all the modes to achieve a desired intensity through amplitude weighting is to keep their phases. For each point in the sampled wave function, the current value of the all the modes can be seen as a point in a higher dimensional space spanned by the real and imaginary components of the wave for each mode. A hyper-cylinder gives the desired intensity. The closest point on the hypersphere to any point will be contained on a line from the origin to that point. Hence, this operation is a projection operation in that space, and these projections are performed serially.

We would like to propose a modified algorithm that combines ideas from both [69] and [70]. This results in two changes to the algorithm presented in [70] – the first of which being that we would like to convert this algorithm into a parallel algorithm. It has been shown at least for iterated projections onto convex sets that the parallel projection method provides better results compared to the serial projections method in the face of inconsistent constraints [71], and for the case of beam synthesis, a desired beam is very likely to be nonphysical and hence would result in inconsistent constraints. A serial approach would also mean that the constraint that received the last projection would have the least error, resulting in an uneven distribution of error across the constraints.

The second change we would like to propose is based on the fact that there is a mismatch in the required sampling rate for the wave function and the intensity. Hence, if we only apply a projection at the sample locations of the SLM mode patterns when we've propagated these modes to a plane with a desired intensity, we would be neglecting three-quarters of the samples of the intensity. In fact, there is no straight-forward way of projecting the modes onto a desired intensity – there are four possibly mutually inconsistent constraints at each plane, formed by four subsets of the intensity samples, i.e. one for samples with even $x$ and $y$ integer lattice coordinates, one for samples with even $x$ coordinates and odd $y$ coordinates, etc. However, we can treat these four constraints separately, and in fact Piestun and Shamir's block

projection method allows for the separation of parts of the transverse plane into different constraints. For the actual change to the algorithm, we can specify that the weighting is the same for the four subsets of the plane. If we do this, then the projection operation is simply propagation to the desired plane, including an upsampling factor as we have done in our nonlinear conjugate gradients algorithm, scaling the modes by an amplitude mask, and then back propagating.

With these changes, each step in the projection-based approach updates each mode with the following expression:

$$\mathbf{u}_{k,(i)} = \frac{1}{M} \sum_{m=1}^{M} \sum_{n=1}^{(2N-1)(2N-1)} \mathbf{p}_{m,n} \sqrt{\frac{y_{m,n}}{\hat{y}_{m,n,(i-1)}}} \mathbf{p}_{m,n}^{H} \mathbf{u}_{k,(i-1)} \qquad (4.35)$$

where

$$\hat{y}_{m,n,(i-1)} = \left| \mathbf{p}_{m,n}^{H} \mathbf{u}_{k,(i-1)} \right|^2 = y_{m,n} - \Delta y_{m,n,(i-1)} \qquad (4.36)$$

is the current propagated intensity for a particular point and $\Delta y_{m,n,(i-1)}$ is the error.

If we make the approximation that forward and back-propagation is an identity operation (this is not necessarily the case unless the intensity samples cover the entire field at that transverse plane, but is a good approximation if the intensity samples cover a majority of the field):

$$I = \frac{1}{M} \sum_{m=1}^{M} \sum_{n=1}^{(2N-1)(2N-1)} \mathbf{p}_{m,n} \mathbf{p}_{m,n}^{H} \qquad (4.37)$$

If we substitute Eq. (4.37) and Eq. (4.36) back into Eq. (4.35), then we obtain an expression for a step in the parallel projections algorithm:

$$\mathbf{\Delta} u_{k,(i)} \approx \frac{1}{M} \sum_{m=1}^{M} \sum_{n=1}^{(2N-1)(2N-1)} \mathbf{p}_{m,n} \frac{\sqrt{y_{m,n}} - \sqrt{y_{m,n} - \Delta y_{m,n,(i-1)}}}{\sqrt{y_{m,n} - \Delta y_{m,n,(i-1)}}} \mathbf{p}_{m,n}^{H} \mathbf{u}_{k,(i-1)} \qquad (4.38)$$

Let us now look at three different cases:

1. When $\Delta y_{m,n,(i-1)}$ is almost equal to $y_{m,n}$. This happens when the current intensity is almost zero.

2. When $\Delta y_{m,n,(i-1)}$ is very negative. This happens when the current intensity is way too high.

3. When $\Delta y_{m,n,(i-1)}$ has a small magnitude compared to $y_{m,n}$. This happens when the current intensity approaches the desired intensity.

For case 1, the very small denominator causes the fraction to blow up. This results qualitatively in the algorithm working very hard to fix under-illuminated regions in the constraint planes.

For case 2, the the fraction approaches $-1$. Therefore, the highest intensities receive the same correction factor as slightly lower intensities. Hence, if the intensity is too high, the algorithm tries to lower them, but doesn't force the highest ones down more strongly.

For case 3, when $\Delta y_{m,n}$ is small compared to $y_{m,n}$, i.e. assume that we are near convergence with this projections algorithm, we can perform a linear Taylor expansion about $y_{m,n}$ for the fraction and obtain:

$$\mathbf{\Delta} u_{k,(i)} \approx \frac{1}{M} \sum_{m=1}^{M} \sum_{n=1}^{(2N-1)(2N-1)} \mathbf{p}_{m,n} \frac{\Delta y_{m,n,(i-1)}}{2y_{m,n}} \mathbf{p}_{m,n}^{H} \mathbf{u}_{k,(i-1)} \qquad (4.39)$$

If we compare this expression for the iteration step to the expression for the steepest descent direction in Eq. (4.27), we find that when the parallel projection algorithm approaches the desired intensity, the step direction is approximately the same as the steepest descent algorithm with the weighting factor $w_{m,n}$ being one over the point-wise square root of the desired intensity. Hence, near convergence, this algorithm should yield a solution where we allow for more error in brighter regions and less error in darker regions. This can become problematic for intensity patterns with large areas of very low intensity, making those regions artificially weighted heavily. This will be investigated later in the results section.

Lastly, without a global line search coupled to the iteration step, it is possible (though rare) for this iterative projections algorithm to have increasing least squares error at the next iteration before convergence.

## 4.3.2 Results

Now that we have described an algorithm for computing a set of modes for approximating a three-dimensional intensity pattern and thus in turn achieving a desired mutual intensity, let us now look at two example problems. For both cases, the algorithm was implemented in MATLAB and run on commodity hardware.

First, in order to understand a bit more about the behavior of the algorithm, we shall look at a simple example – generating a Gaussian beam. Now, obviously a Gaussian beam should have a fully coherent solution, but if we start distorting the beam by compressing or expanding the intensity pattern along the optical axis, we no longer have a physically plausible coherent beam, as discussed in the previous chapter. However, we will see that we can still achieve decent results when the beam is compressed up to a point through the use of partially coherent beams, as the existence of Gaussian Schell-model sources would suggest.

In the second example problem, we shall look at replicating three very different images at different depths. Through this, we will be able to see how well the optimization deals with a fully synthetic pattern, how many modes are required, how our optimization algorithm compares to a iterative projection style algorithm and how closely we can space these three images.

Optimized solutions to both example problems will be verified through simulated propagation of the SLM patterns to see whether they achieve the desired three-dimensional intensity pattern.

**Distorted Gaussian beam**

For the first example problem, we will use a $20 \times 20$ SLM with $\Delta_{SLM} = 20\lambda$ pixel pitch. A total of $M = 17$ planes of $39 \times 39$ voxels will be specified, with the planes evenly spaced between $-z_{max}$ and $z_{max}$, where $z_{max}$ was calculated to be the distance between two transverse planes at which a "ray" travelling at the maximum angle allowed by the $4 - f$ system from the center of one plane will reach the edge of the other:

$$z_{max} = \frac{1}{2}\Delta_{SLM}N\sqrt{4\Delta_{SLM}^2 - 1} \tag{4.40}$$

where $N$ is the number of elements along one dimension (20 in this case). In this particular instance, $z_{max}$ turns to be to approximately $8000\lambda$.

A total of 90 different runs of the optimization algorithm (150 iterations per run) were performed, with variations along two axes – number of modes ($K$) and compression ratio ($\kappa$). The compression ratio determined how "compressed" or "expanded" the Gaussian beam was. The intensity at each voxel $\hat{x}, \hat{y}$ at each plane $m = 1, \ldots, M$ was computed by the classic Gaussian beam intensity formula with additional scaling by $\kappa$ in $z$:

$$I_m[\hat{x}, \hat{y}] = \frac{1}{\sqrt{1 + \lambda^2 \kappa^2 z_m^2/(\pi^2 w_0^4)}} e^{(-2\Delta_{SLM}^2(\hat{x}^2 + \hat{y}^2))/(w_0^2(1 + \lambda^2 \kappa^2 z^2/(\pi^2 w_0^4)))} \tag{4.41}$$

where $z_m$ is the $z$ value for that particular plane and a beam waist $w_0$ of $50\lambda$ was used. The compression ratio $\kappa$ was varied through 15 different values spaced equally in logarithmic space between $1/5$ (stretched out) and 5 (compressed). A subset of the planar intensity patterns are shown in Figure 4.13. The iterations were initialized with a random complex noise pattern.

The results of the optimization are summarized in Figure 4.14. The resulting mode pattern was propagated through the original constraint planes using simulation and the root mean square error across all the planes (the entire volume) was computed for each combination of stretch ($\kappa$) and number of modes ($K$).

The first observation is that for the fully coherent case, only a true Gaussian beam can be generated with minimal error. This agrees with our analysis in the previous chapter using phase space. This also verifies that the algorithm can generate a well-known intensity pattern. We can see the intensity patterns generated by the fully coherent case in Figure 4.15. Note that for the case of longitudinal expansion $\kappa = 0.2$, the generated intensity pattern has many rings, reminiscent of a Bessel beam, a known propagation invariant beam. As the beam is compressed longitudinally beyond a standard Gaussian beam, the fully coherent case falls apart, creating excess high frequency patterns and nulls in the intensity.

The observation to be gleamed from Figure 4.14 is that partial coherence indeed
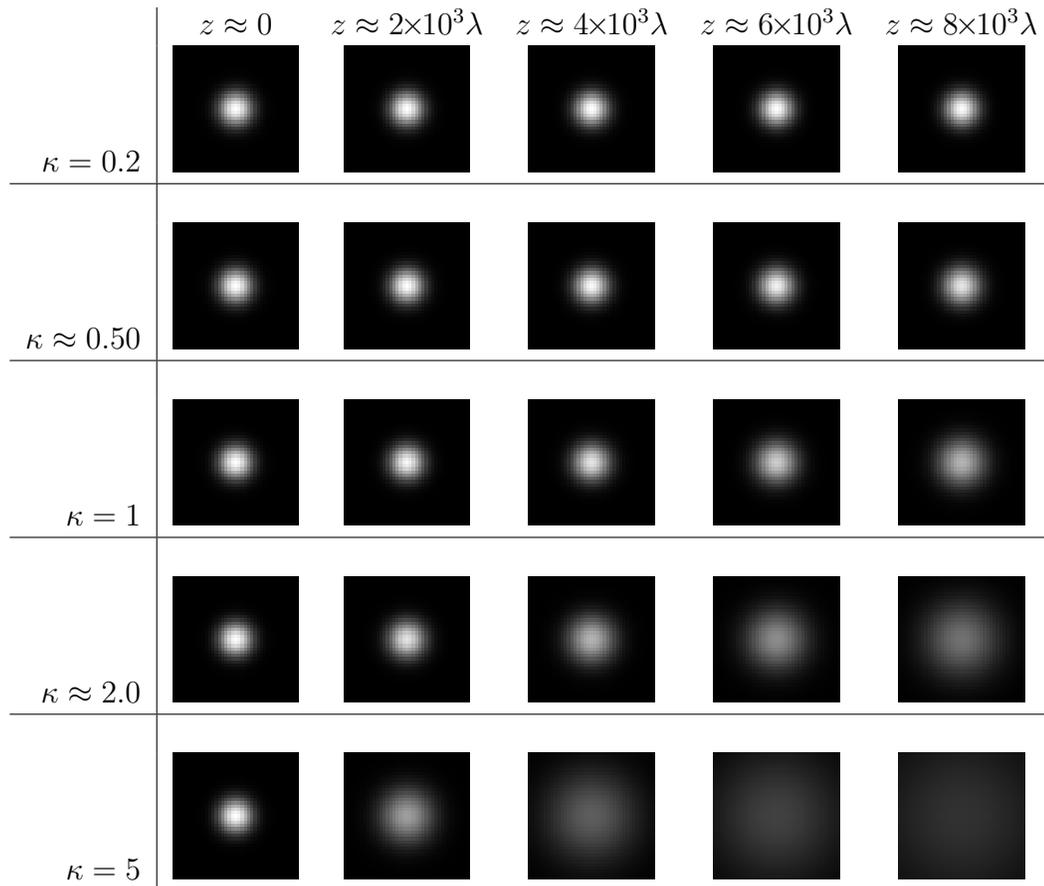
Figure 4.13: Desired intensity patterns at planes $m = 9, 11, 13, 15, 17$ (columns) for various compression ratios $\kappa$ (rows). Whiter indicates higher intensity and all images are on the same scale. Gamma correction of 0.5 has been applied to each image to improve visibility of darker details. Each row contains a set of intensity patterns for a specific compression ratio $\kappa$. Only positive $z$-valued planes are shown, since the voxel intensity pattern is symmetric about $z = 0$. Note that at small compression ratios (first row), the beam becomes almost propagation invariant. At large compression ratios (last row), the beam expands much faster than an actual Gaussian beam (center row).

Figure 4.14: Root mean square error in intensity over the entire volume for an opti-mized (partially) coherent beam attempting to reconstruct an intensity pattern equal to that of a Gaussian beam compressed($\kappa > 1$) or expanded($\kappa < 1$) longitudinally. The horizontal $\kappa$ axis is on a logarithmic scale. The maximum intensity at any point in the desired beam is 1.0. The different curves represent differing numbers of modes used in the optimization. Note that partial coherence $K > 1$ helps up to a point in compressed beams, after which the compression requires angular propagation of light that exceeds the NA of the system. Partial coherence is of relatively little use in the case of expansion, i.e. when the Gaussian beam is stretched towards being almost propagation invariant.
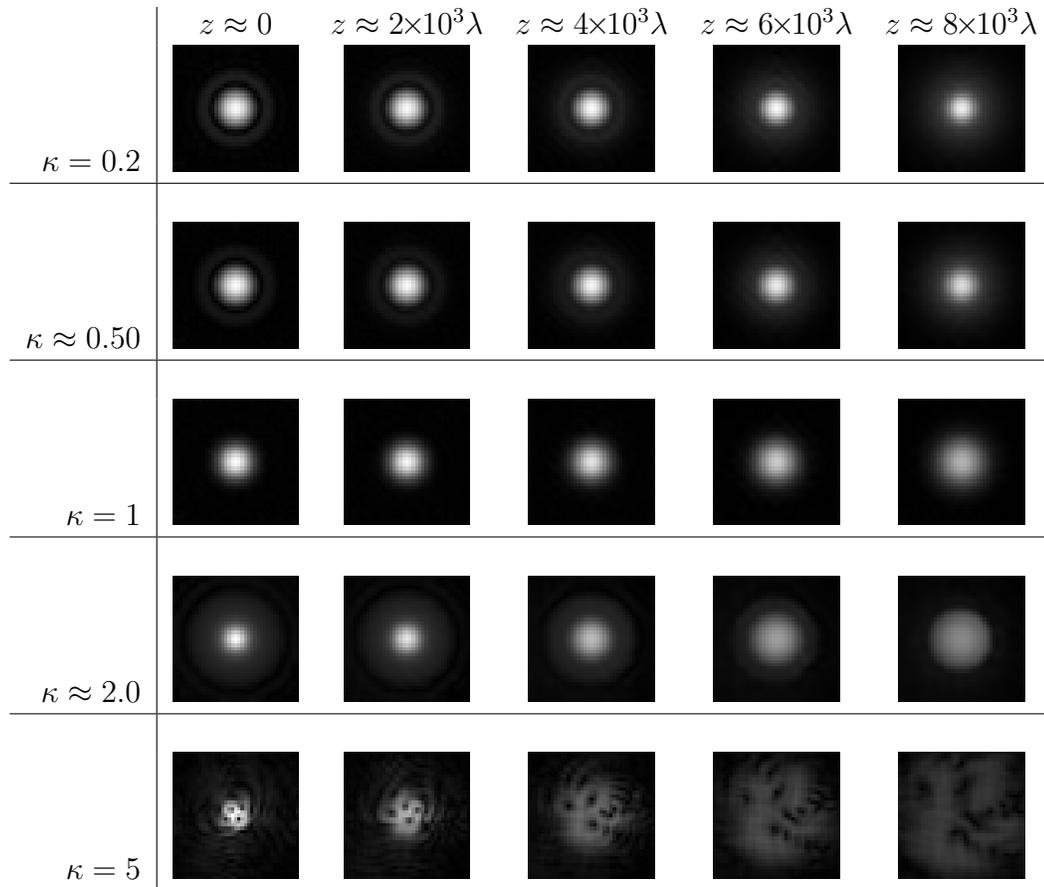
Figure 4.15: Resulting intensity patterns at planes $m = 9, 11, 13, 15, 17$ (columns) from optimization results using $K = 1$ modes for various compression ratios $\kappa$ (rows). Gamma correction of 0.5 has been applied to each image to improve visibility of darker details. Each row contains a set of intensity patterns for a specific compression ratio $\kappa$.

allows for the generation of an intensity pattern when the Gaussian beam is compressed longitudinally – a case where the fully coherent case fails. This agrees with the analysis in the previous chapter where we realized we can decompose a longitudinally compressed Gaussian beam intensity pattern into multiple Gaussian beams oriented at different angles.

Let us now examine detailed results of two different partially coherent setups – a 4-mode setup in Figure 4.16 and a 32-mode setup in Figure 4.17. The $K = 4$ case performs slightly better than the fully coherent case for $\kappa \approx 2.0$ and the $K = 32$ case performs very well. From Figure 4.14, it appears that increased longitudinal compression requires more modes, although anything beyond $K = 8$ seems to yield diminishing returns.

In fact, it appears that no amount of modes can prevent an increase in error at $\kappa \approx 3.5$. As can be seen in the last rows of Figure 4.16 and Figure 4.17, the pattern falls apart, although the $K = 32$ case does approach the desired pattern somewhat. We can still see what appears to be rectangular high frequency components in the $K = 32$ case. This may seem surprising, but there is actually a pretty good explanation for this. If we longitudinally compress the intensity pattern of a Gaussian beam too much, the maximum frequency in the Wigner distribution starts exceeding the maximum allowed by the NA. Thus, the aperture of the system should start factoring in. Hence, the rectangular patterns should not be too surprising, since we are using a square aperture.

Lastly, for the case of expansion $\kappa < 1$, we notice in Figure 4.14 that while the partially coherent case produces slightly less error, the gains are very minimal. In fact, there is minimal difference between the actual generated intensity patterns in Figures 4.15, 4.16 and 4.17 for cases when $\kappa < 1$. This agrees with the fact that the convolution trick explained in the previous chapter cannot be used to create a smaller pattern in the Wigner distribution, as opposed to a bigger pattern.

Observations from the results of this optimization suggest that the algorithm works well enough to verify theoretical predictions. Next, we will investigate the performance of this algorithm with intensity patterns that are not based on a well-known physical light pattern.
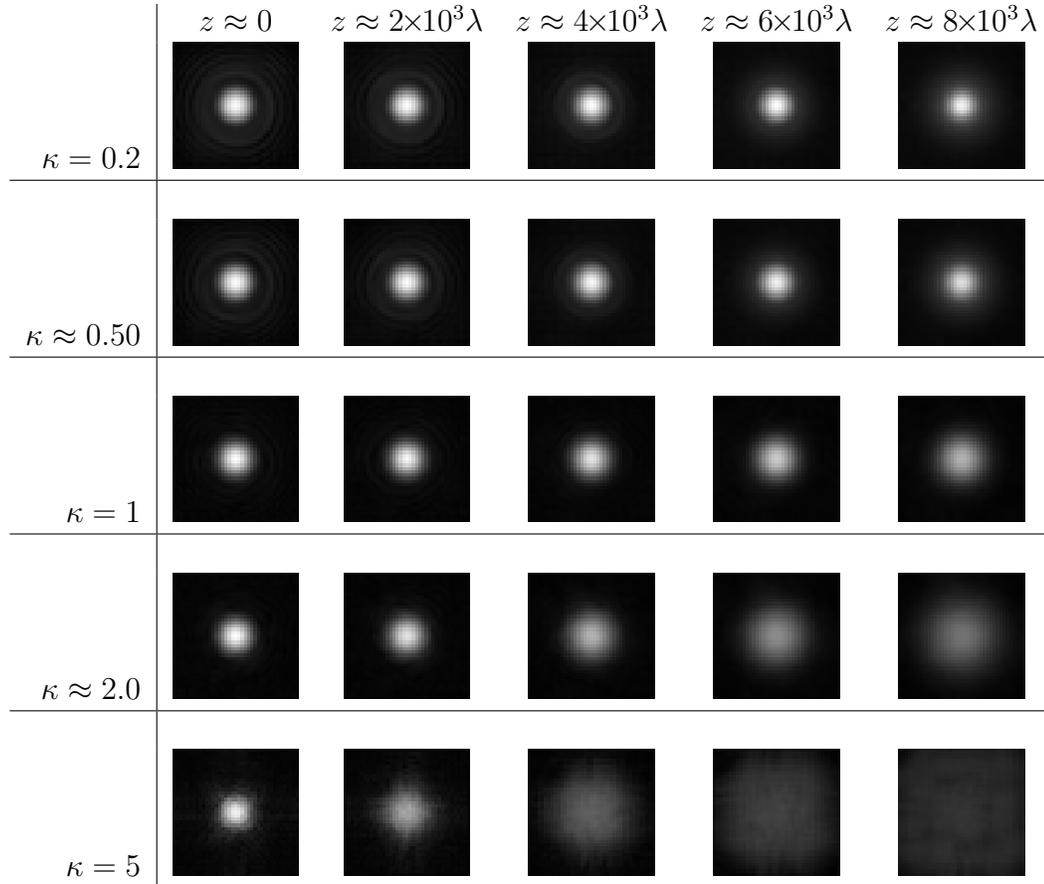
Figure 4.16: Resulting intensity patterns at planes $m = 9, 11, 13, 15, 17$ (columns) from optimization results using $K = 4$ modes for various compression ratios $\kappa$ (rows). Gamma correction of 0.5 has been applied to each image to improve visibility of darker details. Each row contains a set of intensity patterns for a specific compression ratio $\kappa$.
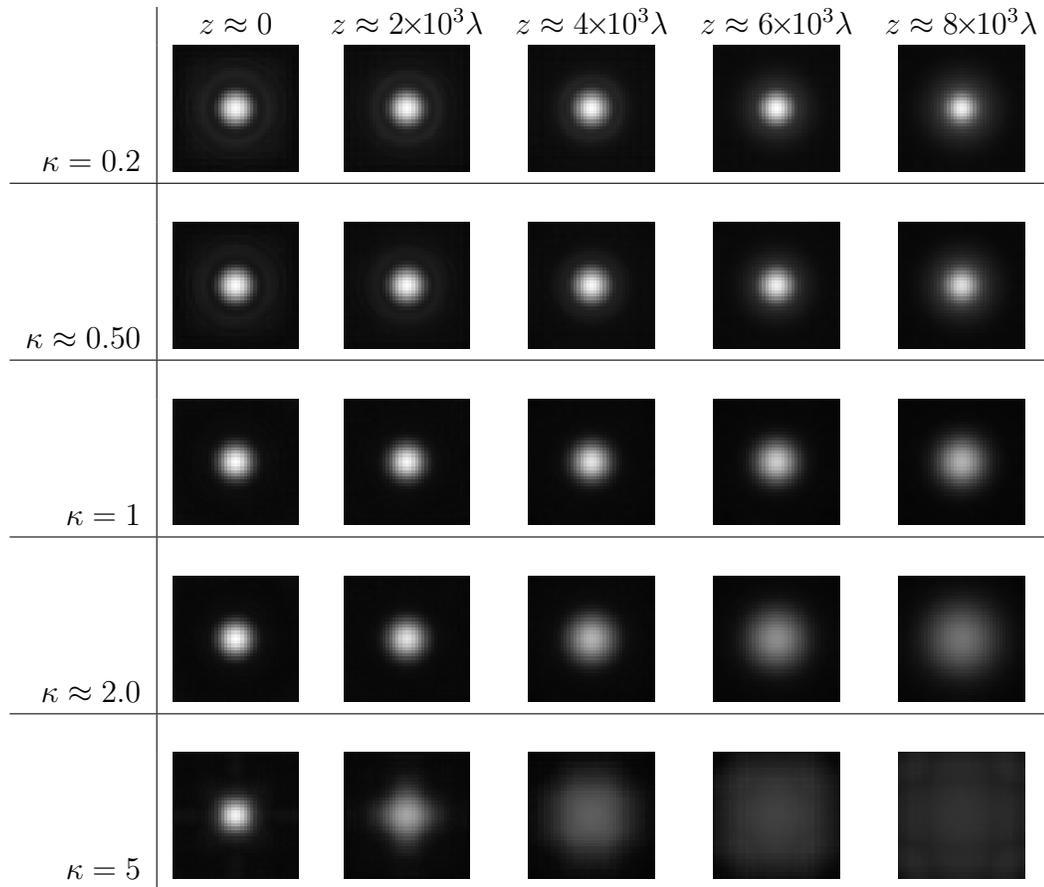
Figure 4.17: Resulting intensity patterns at planes $m = 9, 11, 13, 15, 17$ (columns) from optimization results using $K = 32$ modes for various compression ratios $\kappa$ (rows). Gamma correction of 0.5 has been applied to each image to improve visibility of darker details. Each row contains a set of intensity patterns for a specific compression ratio $\kappa$.

**Three image planes**

For this example problem, we will attempt to produce three different intensity patterns, as shown in Figure 4.18, each at a different transverse plane. We will be using a $64 \times 64$ SLM with pixel pitch $\Delta_{SLM} = 20\lambda$ in order to reproduce the three $127 \times 127$ pixel intensity patterns at locations $z = -\Delta_z$, $z = 0$ and $z = \Delta_z$. The algorithm was run for all combinations of variations along three axes:

1. The number of modes used was either 1, 2, 6 or 12. When the number of modes was 1, a larger $157 \times 157$ SLM as well as a $222 \times 222$ SLM with the same pixel pitch was also implemented.

2. The algorithm used was either the non-linear conjugate gradients (NLCG) algorithm we described, a modified gradient descent algorithm with global line search (GRAD), or the described iterated projection algorithm (PROJ). All algorithms were run for 400 iterations. When algorithm is not specified, NLCG is assumed.

3. The spacing between planes, $\Delta_z$, could be set to any one of $z_{max}$, $z_{max}/2$, $z_{max}/4$ and $z_{max}/8$, where $z_{max}$ has the same definition as in the previous example problem:
$$z_{max} = \frac{1}{2}\Delta_{SLM}N\sqrt{4\Delta_{SLM}^2 - 1} \tag{4.42}$$
where $N$ here is 64, even in the case of the larger SLMs. $z_{max}$ comes to be approximately $25600\lambda$ in this case. When $\Delta_z$ is not specified, assume $z_{max}$.

Quality of the optimization was measured by the image quality at each of the three transverse planes through the peak signal-to-noise ratio (PSNR):

$$PSNR = 10\log_{10}\left[I_{max}^2 / \left(\sum_{i=1}^{127}\sum_{j=1}^{127}(I_d[i,j] - I_s[i,j])^2\right)\right] \tag{4.43}$$

where $I_d[i,j]$ is the desired discrete intensity pattern, $I_s[i,j]$ is the simulated intensity pattern of the output of the optimization and $I_{max}$ is the maximum value of $I_d[i,j]$. Thus, the image plane containing the modified Shepp phantom received less weight
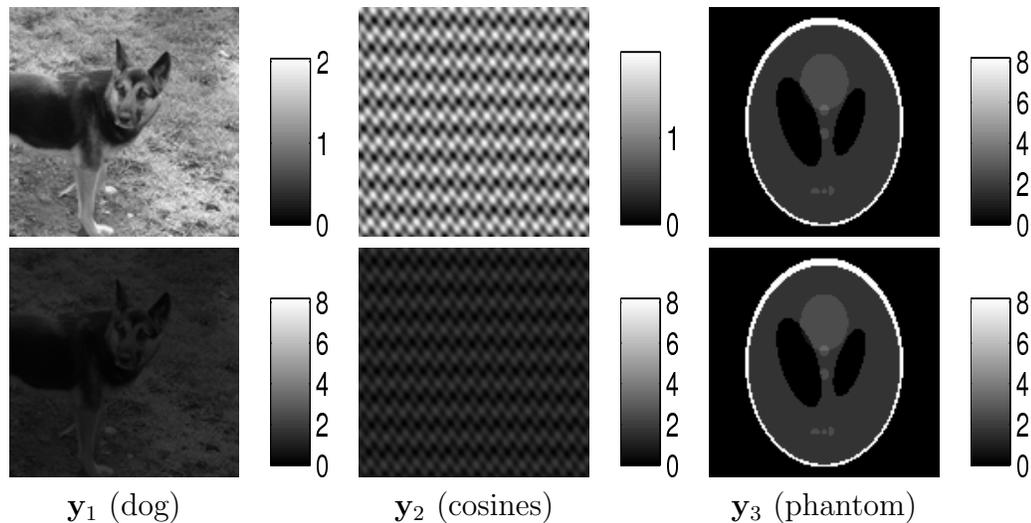
$\mathbf{y}_1$ (dog)              $\mathbf{y}_2$ (cosines)              $\mathbf{y}_3$ (phantom)

Figure 4.18: We would like to show three intensity patterns, one for each transverse plane. Each intensity pattern is a $127 \times 127$ grayscale image representing the desired sampled representation of a continuous and bandlimited intensity function at that plane. The first pattern is a photograph of a dog, which was included with the ImageStack software package distribution. The second is a sum of three raised cosine waves. The third is the modified Shepp phantom produced by MATLAB's `phantom` command with negative values set to zero. Note that the data for the phantom is actually not a sampled representation of an all-positive function due to Gibbs phenomena. All intensity patterns used in the computation were scaled so that the total intensity at each plane was equal. That is why the Shepp phantom image has a higher maximum value, since the majority of its pixels are black. For illustration purposes, the images in the top row have been normalized individually where white represents the brightest pixel. The bottom row shows the actual relative brightness of the three patterns, where white represents the brightest pixel among all the three patterns.

in terms of absolute error across all the algorithms in order to ensure equal weighting of the PSNR.

First, we will examine the effect of the number of modes on the image quality. Since the goal is to reconstruct three distinct planes of $127 \times 127$ pixels, we would need $3 \times 127 \times 127$ or roughly $48,000$ degrees of freedom of control. Each SLM pattern offers $64 \times 64 \times 2$ degrees of freedom of control. Hence, roughly 6 SLM patterns should be needed to create the desired three-plane intensity pattern.

The results of the optimization algorithm as the number of modes is increased can be seen in Figure 4.19. Increasing the number of modes improves the image quality. However, increasing the number of modes beyond 6 does not result in appreciably significant gains.

The best results appear to be fairly good for the photograph and the sinusoidal pattern, although the Shepp phantom exhibits noticeable ringing. The photograph does exhibit slight blurring in the corners, if we carefully compare the image to the desired image from Figure 4.18. This is due to the fact that the corners see less of the SLM and thus have a smaller effective aperture, leading to resolution loss. The ringing from the phantom could also be attributed to this, but there is also something else fundamentally hard with the phantom image – the desired phantom image represented a sampled representation of a real continuous function that was not all positive due to Gibbs phenomena.

A simpler one-dimensional example would be a sequence of samples of value 1 followed by a sequence of samples of value 0. None of these samples are negative, but the bandlimited continuous function represented by these samples would have negative values near the boundary between the 1s and the 0s, which is apparent from the negative side lobes in the sinc function used in bandlimited reconstruction. Hence, the algorithm simply attempted to find the best match in terms of error, even when the input was a desired intensity pattern that violated non-negativity constraints in the continuous domain.

The PSNR score for each plane as a function of iteration number can be seen in Figure 4.20. By iteration 400, gains in PSNR has tapered off in all cases, and therefore PSNR comparisons after 400 iterations should be acceptable. These graphs
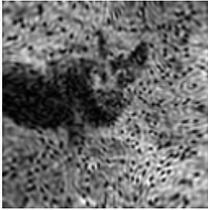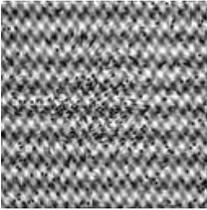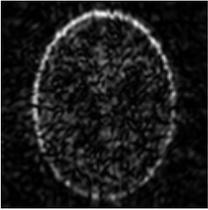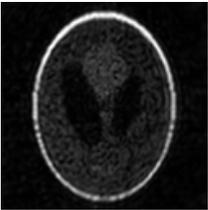
|  | dog(plane 1) | cosines(plane 2) | phantom(plane 3) |
|---|---|---|---|
| $K = 1$ | PSNR=16.4dB | PSNR=18.9dB | PSNR=16.5dB |
| $K = 2$ | PSNR=24.6dB | PSNR=25.8dB | PSNR=19.8dB |
| $K = 6$ | PSNR=27.6dB | PSNR=28.4dB | PSNR=21.3dB |
| $K = 12$ | PSNR=28.1dB | PSNR=28.9dB | PSNR=21.5dB |

Figure 4.19: The resulting intensity patterns at the three target planes from the output modes calculated by the optimization algorithm. Each row signifies the number of modes used in the optimization and each column is a different plane. All images have been scaled such that white corresponds to the brightest pixel in the *desired* intensity pattern.
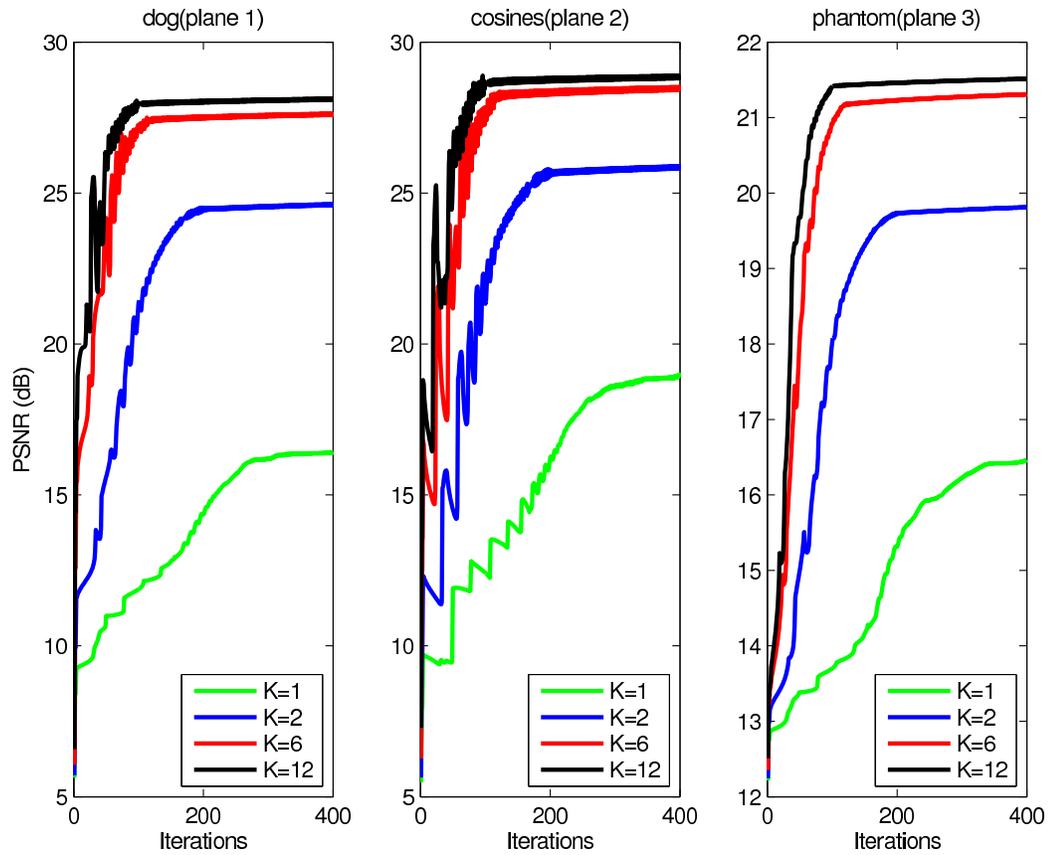
Figure 4.20: Progression of image quality (PSNR) at each plane for each optimization. Each graph is for a separate transverse plane and each line represents a different number of modes used for the optimization. Note that the vertical axis scaling for the third plane is different due to overall poorer performance.

also illustrate numerically that image quality does not improve significantly when we increase the number of modes beyond 6, giving support to our hypothesis that we will need approximately 6 SLM patterns. These results support the obvious conclusion that we need enough degrees of freedom of control (which is lacking in the $K = 1$ and $K = 2$ cases) to create these three separate image patterns

We will now explore another method of increasing the number of degrees of freedom – we will increase the number of pixels on the SLM and retain using a fully coherent setup instead of increasing the number of modes. In order to accomplish this without changing the resolution, the SLM itself should be expanded while retaining the same pitch. The $64 \times 64$ SLM can be expanded to one with $157 \times 157$ pixels to obtain at least a 6-fold increase in pixels or to one with $222 \times 222$ pixels for a 12-fold increase in pixels. If we consider this setup from a different angle, what we are doing is specifying a much smaller region-of-interest than the size of the beam. This has been suggested in the literature before as a way to combat the degrees-of-freedom issue [69].

The resulting images from the algorithm due to this change in configuration can be seen in Figure 4.21. Note that while increasing the degrees of freedom by a factor of 6 does improve the image quality in the fully coherent case, the nulls (black spots) in the resulting intensity images can not be removed even with additional degrees of freedom (i.e. a factor of 12), while the partially coherent case with the same number of degrees of freedom performs quite well. The PSNR graph as a function of iteration count is shown in Figure 4.22

Yet another axis we can explore is the optimization procedure. In addition to the proposed NLCG algorithm, we can also look at removing the conjugate gradients portion to obtain a gradient descent algorithm with global line search (GRAD). Lastly, we should also compare results to the iterated projections algorithm (PROJ). These algorithms were run for the 6-mode case and the resulting images can be seen in Figure 4.23 and the PSNR graphs can be seen in Figure 4.24.

From the figures, it is not apparent that there is much difference between the algorithms. In fact, the PROJ algorithm seems to generate a higher quality image for the Shepp phantom case. The PSNR graphs show a very small advantage to using
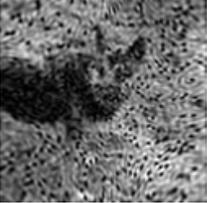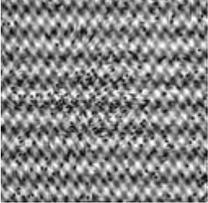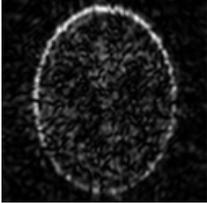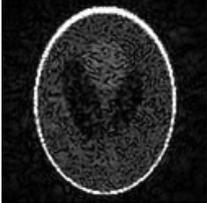
|  | dog(plane 1) | cosines(plane 2) | phantom(plane 3) |
|---|---|---|---|
| $K = 1 \; D = 1/6$ | PSNR=16.4dB | PSNR=18.9dB | PSNR=16.5dB |
| $K = 1 \; D = 1$ | PSNR=19.0dB | PSNR=21.2dB | PSNR=21.1dB |
| $K = 1 \; D = 2$ | PSNR=19.2dB | PSNR=21.2dB | PSNR=21.0dB |
| $K = 6 \; D = 1$ | PSNR=27.6dB | PSNR=28.4dB | PSNR=21.3dB |

Figure 4.21: The resulting intensity patterns at the three target planes from the output modes calculated by the optimization algorithm. Each row signifies the number of modes ($K$) and the relative number of degrees of freedom ($D$) compared to the output used in the optimization and each column is a different plane. All images have been scaled such that white corresponds to the brightest pixel in the *desired* intensity pattern.

Figure 4.22: Progression of image quality (PSNR) at each plane for each optimization. Each graph is for a separate transverse plane and each line represents a different configuration (SLM size, number of modes) used for the optimization. Numbers in parentheses are the relative number of degrees of freedom compared to the output number of pixels. A "coh" indicates fully coherent operation and "pc" indicates partially coherent operation. Note that the vertical axis scaling for the third plane is different due to overall poorer performance.

|  | dog(plane 1) | cosines(plane 2) | phantom(plane 3) |
|---|---|---|---|
| NLCG | PSNR=27.6dB | PSNR=28.4dB | PSNR=21.3dB |
| GRAD | PSNR=27.5dB | PSNR=28.3dB | PSNR=21.2dB |
| PROJ | PSNR=24.7dB | PSNR=23.3dB | PSNR=19.9dB |

Figure 4.23: The resulting intensity patterns at the three target planes from the 6 output modes calculated by various optimization algorithms. Each row signifies the algorithm used in the optimization and each column is a different plane. All images have been scaled such that white corresponds to the brightest pixel in the *desired* intensity pattern.
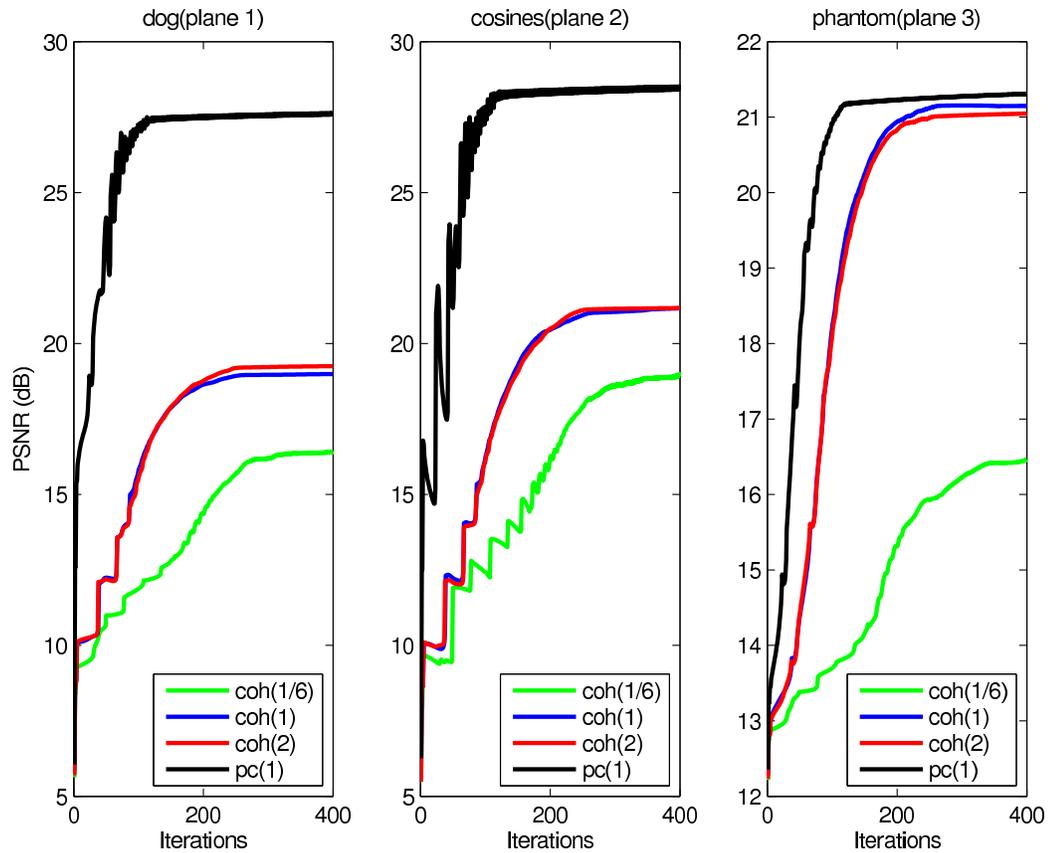
Figure 4.24: Progression of image quality (PSNR) at each plane for each optimization algorithm. Each graph is for a separate transverse plane and each line represents a different optimization algorithm. Note that the vertical axis scaling for the third plane is different due to overall poorer performance.

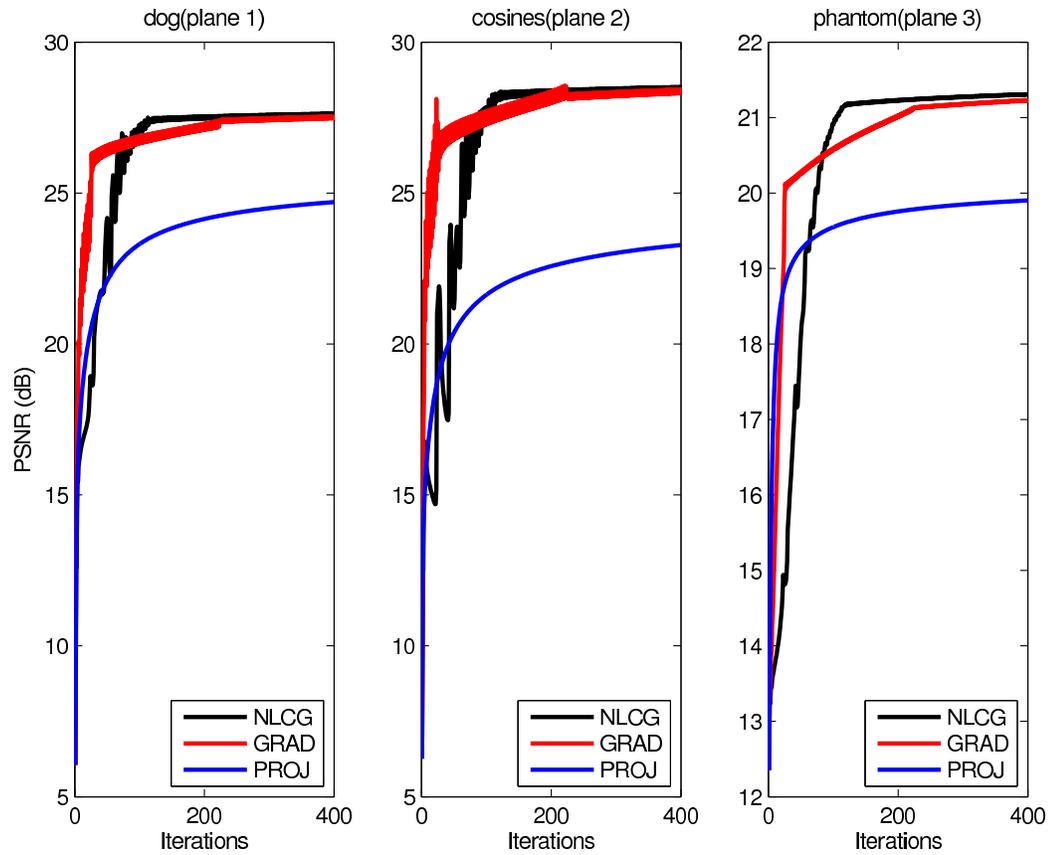the NLCG algorithm over the GRAD algorithm, but a sizable advantage to using the NLCG algorithm over the PROJ algorithm.

Compared to the GRAD algorithm, the NLCG algorithm seems to converge faster, but it also exhibits significantly more oscillations in the PSNR scores over iterations. One possibility is that each of these bumps is when the optimization finds a new valley to go down, sacrificing PSNR at one plane for improved PSNR at other planes. The attempt at mutual orthogonality of subsequent steps in NLCG may cause the algorithm to shoot search directions across a wider region and thus take advantage of the global line search more; recall that the global line search takes the optimization to the global minimum along a line and thus the wider a region the search line reaches, the more likely it will find a much better location.

Numerically, compared to the NLCG algorithm, the PROJ algorithm seems to achieve a lower PSNR score, even though the images look very much the same, with even the Shepp phantom image appearing to have less error. The differences, however, can be seen more clearly by observing the intensity pattern along a horizontal strip across the center of each image plane, as shown in Figures 4.25, 4.26 and 4.27. The reason for PROJ's poor scores in the PSNR is especially clear in Figure 4.27. Here, the PROJ algorithm does a very good job for approximating the medium-low brightness "middle" portion of the phantom with intensity 1.64. However, the NLCG algorithm is closer for the bright edges with intensity 8.20. The error magnitude in the NLCG algorithm much more uniform than the PROJ case, where low intensity regions have much less error than higher intensity regions. In fact, at regions where the intensity should be zero, the PROJ case is very close, whereas the NLCG case shows noticeable ringing. This dependence of error magnitude on target intensity for the PROJ algorithm is suboptimal for a least squares merit function, since high error regions are accentuated. However, this dependence agrees with the analysis of the iterative projections algorithm conducted previously for the case where the resultant intensity is close to the target intensity.

At this point, it should very obvious that the coherent field methods illustrated in [69] are not adequate to generate the images for this particular example. This is illustrated explicitly in Figure 4.28, where the top row of images corresponds to

Figure 4.25: Horizontal cross sections (the center row in each image) of the resulting intensity at the first image plane (dog) for the nonlinear conjugate gradients (NLCG) algorithm and the iterative projections algorithm (PROJ) compared to the ideal desired intensity (ORIG).

Figure 4.26: Horizontal cross sections (the center row in each image) of the resulting intensity at the second image plane (cosines) for the nonlinear conjugate gradients (NLCG) algorithm and the iterative projections algorithm (PROJ) compared to the ideal desired intensity (ORIG).
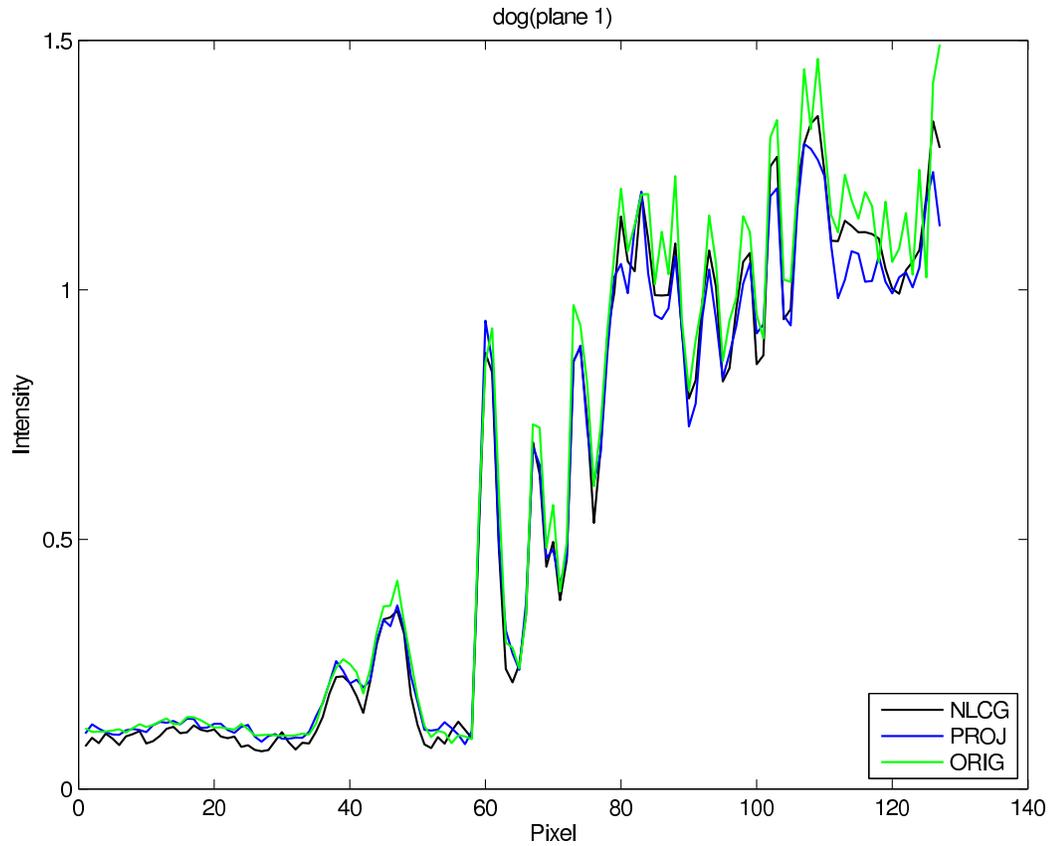
Figure 4.27: Horizontal cross sections (the center row in each image) of the resulting intensity at the third image plane (phantom) for the nonlinear conjugate gradients (NLCG) algorithm and the iterative projections algorithm (PROJ) compared to the ideal desired intensity (ORIG).
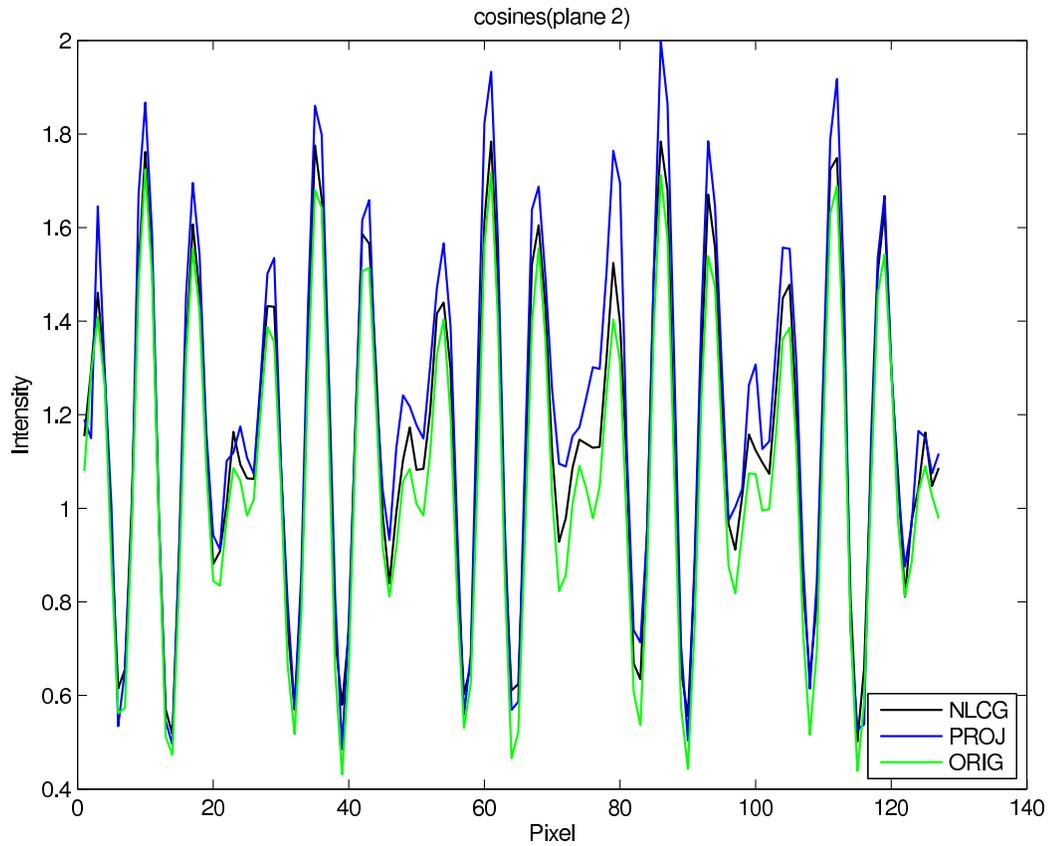
| dog(plane 1) | cosines(plane 2) | phantom(plane 3) |
|:---:|:---:|:---:|
| PSNR=19.0dB | PSNR=21.2dB | PSNR=21.1dB |
| Ideal | Ideal | Ideal |
| PSNR=27.6dB | PSNR=28.4dB | PSNR=21.3dB |

Figure 4.28: The first row of intensity patterns is the result of running an iterative projections algorithm combined with the use of increasing the size of the SLM to support more degrees of freedom in the purely coherent case. The second row of intensity patterns is the ideal desired intensity patterns. The third row of intensity patterns is the result of running the nonlinear conjugate gradients algorithm for a 6-mode partially coherent beam. All images have been scaled such that white corresponds to the brightest pixel in the *desired* intensity pattern.

the use of the PROJ algorithm with the fully coherent degree-of-freedom matching method alluded to in [69]. The partially coherent method from [70] with appropriate sampling and parallel projections performs much better, but is still not as numerically optimal as the NLCG algorithm developed for this manuscript.

Let us also examine what effect decreasing the axial spacing between these three image planes has on image quality. We will use the standard NLCG algorithm with 6 modes while varying $\Delta_z$. Figure 4.29 contains the resulting images as $\Delta_z$ is shrunk and Figure 4.30 shows the corresponding PSNR graphs.

As we squeeze the image planes closer together, we should expect to see the results get worse and worse. This is because each pixel on a plane will see less and less of the other planes due to the limited NA and thus have less degrees of freedom in affecting the pixels on a different plane. For example, it is harder for a highly patterned small patch of pixels to distribute its energy evenly over a second image plane if the light from these pixels cannot reach all the pixels on the second image plane. Therefore, as the distance between image planes decreases, image planes start affecting each other – we see ghosts of the dog at the other two planes, the dog photograph has lost contrast, and there also seems to be a phantom image burned onto the sinusoidal pattern. Interestingly enough, while image quality is adversely affected at a spacing of $z_{max}/8$, image quality seems to have an improvement at $z_{max}/2$! In fact, the images get sharper moving from $z_{max}$ to $z_{max}/2$. It turns out that at a transverse plane $z = z_{max}$, only the center pixel would have its angular window filled by the SLM due to the limited NA; the edge pixels would see a much smaller effective NA and thus suffer from resolution loss. Hence, this resolution loss shouldn't be unexpected.

Lastly, before we move on, let's look at the resulting PSNR for all the different combinations run in the simulation. Notable patterns include that the iterated projections algorithm tends to fare better at the phantom and that the best image quality for partially coherent beams seems to not be at $z_{max}$, although there is an overall positive trend across all the different setups for increasing quality with increasing $\Delta_z$. It's interesting that the fully coherent extended degrees of freedom setups simply fare better when the image planes are spaced farther apart, but that is because the SLMs are larger in area and thus don't introduce the resolution loss issue discussed earlier.

|  | dog(plane 1) | cosines(plane 2) | phantom(plane 3) |
|---|---|---|---|
| $\Delta_z = z_{max}$ | PSNR=27.6dB | PSNR=28.4dB | PSNR=21.3dB |
| $\Delta_z = z_{max}/2$ | PSNR=29.5dB | PSNR=28.9dB | PSNR=23.0dB |
| $\Delta_z = z_{max}/4$ | PSNR=30.0dB | PSNR=27.6dB | PSNR=22.3dB |
| $\Delta_z = z_{max}/8$ | PSNR=26.1dB | PSNR=24.8dB | PSNR=19.4dB |

Figure 4.29: The resulting intensity patterns at the three target planes from the output modes calculated by the optimization algorithm. Each row signifies the spacing between planes and each column is a different plane. All images have been scaled such that white corresponds to the brightest pixel in the *desired* intensity pattern.

Figure 4.30: Progression of image quality (PSNR) at each plane for the optimization algorithm. Each graph is for a separate transverse plane and each line represents a different spacing between desired intensity planes, with the legend in units of $z_{max}$. Note that the vertical axis scaling for the third plane is different due to overall poorer performance.

Figure 4.31: A summary of resulting PSNR for every combination of algorithm (column), number of modes (row) and inter-planar spacing (horizontal axis). The three planes appear as three different lines in each graph.

### 4.3.3 Summary

In this section, we've developed an iterative optimization algorithm for the calculation of mode patterns in order to generate a desired voxel intensity pattern in space. The desired mutual intensity can be calculated by performing an product of the matrix containing the modes arranged as column vectors with the matrix's conjugate transpose. We've also applied this optimization algorithm to two different test cases, both yielding results that confirm our expectations.

## 4.4 Closer look at modes

Now that we have looked at some algorithms to compute the mutual intensity of desired partially coherent illumination for certain situations, let's discuss in more detail mode sequences themselves. That is, let's look closer at temporal multiplexing and the sequence of amplitude-phase SLM patterns used to create partially coherent illumination. For one, given a specific desired mutual intensity, there actually exists an infinite number of temporal patterns that can recreate it. Furthermore, we'll also look at the case of partially coherent beams that have small coherence area, in which case we can actually reduce the number of modes required.

### 4.4.1 Choice of mode pattern
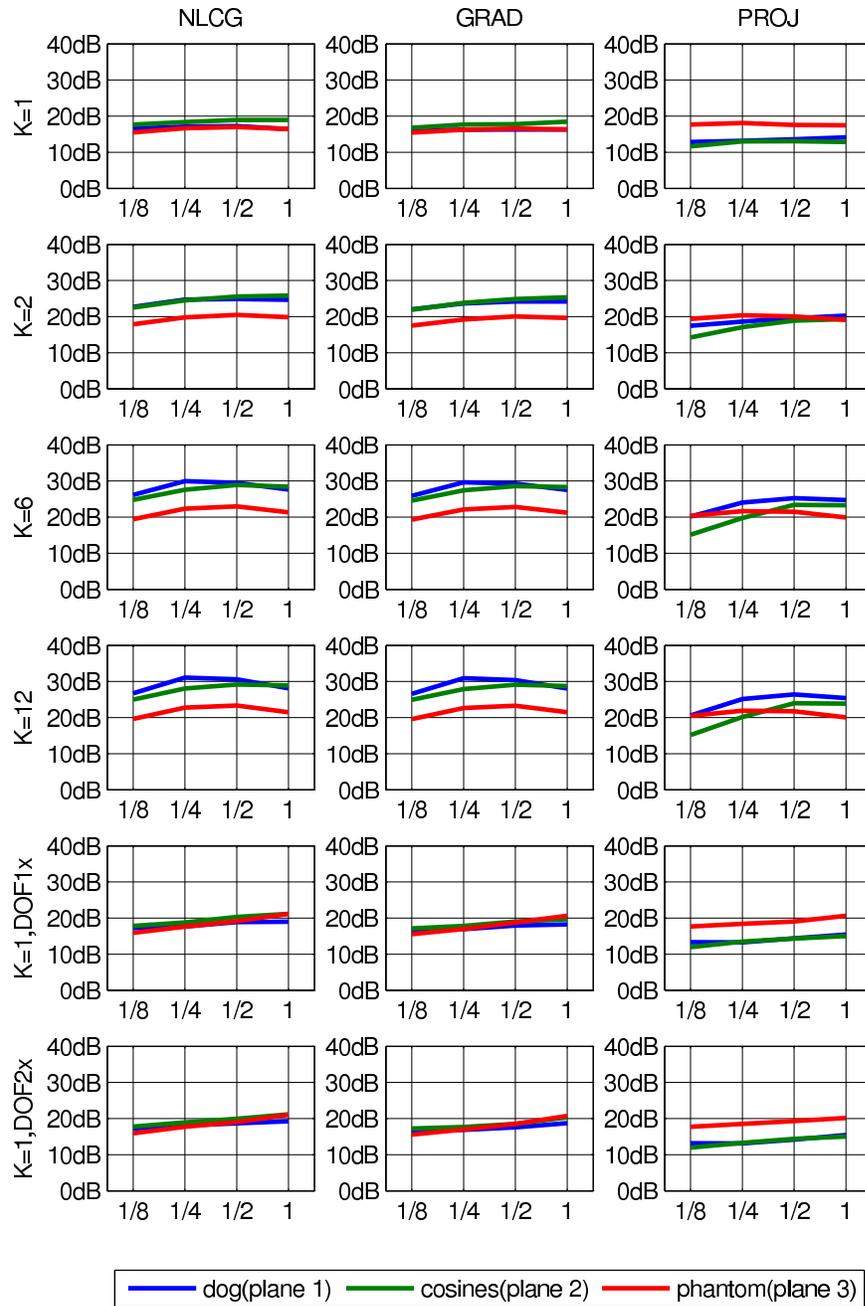
Let $J \in \mathbb{C}^{N^2 \times N^2}$, a Hermitian positive definite matrix, be the discretization of the desired mutual intensity, where row and column indexes each map to specific spatial positions. As long as we have some matrix $U \in \mathbb{C}^{N^2 \times K}$ such that $UU^H = J$, then we can use the columns of matrix $U$ as sequential SLM patterns to generate the desired mutual intensity $J$. For example, the straight-forward way to compute one such $U$ would be to use a singular value decomposition [33, 35]:

$$J = \hat{U}S\hat{U}^H = \hat{U}S^{1/2}(\hat{U}S^{1/2})^H = UU^H \tag{4.44}$$

In other situations, we might have a particular $U$ simply from the optimization technique, e.g. the voxel-based technique described in the previous section.

However, if $U$ is one such matrix, it should be easy to see that any matrix of the form $UQ$ will also generate the desired mutual intensity provided that $Q \in \mathbb{C}^{K \times L}$ has the following property:

$$QQ^H = I \tag{4.45}$$

What this means is that given an input set of modes, we can design with some freedom an output set of modes that recreates the exact same mutual intensity, but with possibly other nice properties.

For instance, if we started with orthogonal modes in decreasing intensity, such as from performing the singular value decomposition (SVD) on the mutual intensity matrix, then we will end up cycling the overall intensity over time and we might be wasting energy unless we can change the laser intensity in sync over time as well. However, if we can use a $Q$ matrix that "smears" and de-orthogonalizes the modes, then we might have a sequence of modes that may have a more even overall intensity distribution, requiring a simple constant adjustment on the input laser power and less attenuation (energy loss) at the SLM.

Suppose we have a matrix $U$ whose columns represent sequential SLM patterns that produce a desired mutual intensity. Let the entry $u_{n,k}$ represent the $n^{th}$ pixel of the $k^{th}$ pattern (mode). Since a SLM can only attenuate light, the incoming light to the SLM must be bright enough to reproduce the pixel with the highest intensity across all the patterns. That is, the incoming intensity at each pixel of the SLM should be:

$$I_{max}(U) = \max_{n,k} |u_{n,k}|^2 \tag{4.46}$$

The average intensity output from each pixel, which is a measure of the total amount of energy passed through the SLM ignoring effects such as inherent efficiency of the SLM and fixed optical setup, can be calculated as:

$$I_{out}(U) = \frac{1}{N^2 K} \sum_{n,k} |u_{n,k}|^2 = \frac{1}{N^2 K} ||U||_F^2 \tag{4.47}$$

Hence, the efficiency of a particular mode pattern $U$ can be calculated as:

$$\eta(U) = \frac{\frac{1}{N^2 K} \sum_{n,k} |u_{n,k}|^2}{\max_{n,k} |u_{n,k}|^2} \tag{4.48}$$

That is, this number gives the amount of energy that is delivered into the beam out of the total incoming energy. A value of 1 is the ideal case, and it means that all the modes are phase-only patterns. Note that if $QQ^H = I$, then $I_{out}(UQ) = I_{out}(U)$, because orthonormal matrices preserve the L2-norm of vectors they operate on (in this case, $Q$ operates on row vectors in $U$). Hence, to increase efficiency while preserving the same output mutual intensity, we must seek to decrease the brightness of the brightest pixel in the set of patterns. Thus, the problem of generating the most efficient set of modes can be formulated as:

> Given a matrix $U \in \mathbb{C}^{N^2 \times K}$, find a matrix $Q \in \mathbb{C}^{K \times K}$ that minimizes $I_{max}(UQ)$ subject to $QQ^H = I$.

This problem can also be thought of as finding the smallest hypercube centered on the origin that fully contains a point cloud. The point cloud has $N^2$ points, one for each row of $U$, and each point has $2K$ coordinates, two for each complex-valued mode.

Unfortunately, a provably optimal algorithm does not seem apparent. However, for a specific mutual intensity $J$, we can establish some bounds on the efficiency of any mode pattern $UU^H = J$ that can generate the specified mutual intensity. First, let's look at the worst case scenario. This can be done by concentrating the energy in the SLM pixel with the highest average intensity into one single mode and zero the other modes for that pixel. In that case,

$$I_{max}^{(worst)} = \max_n \sum_k |u_{n,k}|^2 = \max_n j_{n,n} \tag{4.49}$$

where $j_{n,n}$ is the $(n, n)$ diagonal entry in the desired mutual intensity. The efficiency would then be:

$$\eta^{(worst)} = \frac{\sum_{n,k} |u_{n,k}|^2}{N^2 K \max_n j_{n,n}} = \frac{\sum_n j_{n,n}}{N^2 K \max_n j_{n,n}} \tag{4.50}$$

The best possible allocation would equally allocate the energy of the brightest pixel

evenly among the SLM patterns and would have all the other pixels spread such that they are not brighter than the brightest pixel for any pattern in the SLM sequence. In that case,

$$I_{max}^{(best)} = \frac{1}{K} \max_n j_{n,n} \tag{4.51}$$

Thus,

$$\eta^{(best)} = \frac{\sum_n j_{n,n}}{N^2 \max_n j_{n,n}} \tag{4.52}$$

Note that $j_{n,n}$ also happens to be the average intensity of the $n^{th}$ pixel at the SLM plane. Hence, both the best and worst case scenario efficiency depends on the ratio of the mean temporally-averaged intensity over all the pixels to the temporally-averaged intensity of the "brightest" pixel at the SLM plane. Thus, in some sense, contrast plays a role in how efficient we can generate the desired partially coherent beam. We can get at most a $K$-fold gain in efficiency by designing a clever mode pattern for a specific mutual intensity pattern that can be fully represented by a $K$-mode coherence mode representation.

Now, we did assume that $Q$ was a square matrix in this formulation. What if $Q$ was wide? That is, what if $Q \in \mathbb{C}^{K \times L}$ and that $L > K$? Note that the best case efficiency does not change, since it does not depend on the number of modes. The worst case efficiency could only increase, but that should be pretty obvious, since the worst case for $L > K$ is taking the worst case for $K$ and simply adding mode patterns that are all black. Thus, the best and worst case performance don't increase and we've increased the number of modes needed, which affects our temporal performance, since real devices have limitations on refresh rate.

However, in other application areas, increasing the number of modes might be useful and/or necessary. e.g. to create output modes that are more tailored for specific SLMs, due to the quantization pattern in the complex domain for each pixel. One might imagine many other ways of measuring "niceness" of mode patterns. This can become a fruitful area of research and allow for software and algorithms to compensate for deficiencies in hardware.

## 4.4.2 Small coherence area beams

One other issue with current hardware is the limit on the number of modes that can be generated within a small time window. Thus, it is important to consider how to reduce the number of modes needed. We might perhaps wish to minimize some error measure directly, such as the voxel intensity squared error in the voxel-intensity algorithm or the squared error in the mutual intensity if we perform an SVD on a mutual intensity matrix and only take the first few modes corresponding to the largest singular values. However, let's think about degree of coherence in general and what it says about modes.

If we wanted to generate a fully coherent field, we would simply need one mode, by definition. We'd have control over roughly $2N^2$ degrees of freedom using $N^2$ complex pixel values. Let's suppose, however, that we wanted to generate a fully incoherent field. Its mutual intensity would be a diagonal matrix [35] and it can have up to $N^2$ degrees of freedom. Unfortunately, with the mode-multiplexing method, we would need up to $N^2$ modes or $N^4$ complex pixel values to instantiate this type of illumination. There is obviously some inefficiency here. What we are trying to do is to pick a specific pattern of modes to cancel the off-diagonal entries on the mutual intensity matrix.

However, there's obviously an easy way to generate the desired fully incoherent beam. We can simply illuminate the SLM with an input beam that is fully incoherent and use the amplitude modulation properties of the SLM to create the desired output beam. This should only require $N^2$ pixel values.

Let's take the case of another familiar partially coherent field, the Gaussian Schell-model beam. Recall that this is a beam with a Gaussian intensity profile and Gaussian fall-off in spatial coherence, with mutual intensity as shown in Equation (3.28). With the standard temporal multiplexing method, we would need some number of modes to recreate this beam and in the worst case of fairly incoherent beams, this scales roughly linearly with the ratio of the standard deviation of the intensity profile Gaussian with the standard deviation of the spatial coherence Gaussian [72]. However, as pointed out by De Santis et al., this beam can easily be generated by placing two Gaussian amplitude masks, one at the front focal plane of a lens and one at the back focal

plane [61] with no temporal multiplexing required.

In both the case of the fully incoherent beam and the Gaussian Schell-model beam, we've replaced coherent input illumination with partially coherent illumination created through the use of extended sources. That is, we've used the extended source to create the incoherent aspects of the beam and the SLM to create the coherent aspects of the beam. These observations naturally lead to the possibility of a modified temporal multiplexing system shown in Figure 4.32.



Figure 4.32: An alternative setup for generating partially coherent beams with limited coherence area. Instead of a plane wave source as in Figure 4.1, we use an extended source imaged through a Fourier transforming lens to create a uniform intensity distribution with limited coherence area. The extended source could be, for example, fully incoherent light falling on an amplitude SLM.

The use of this system, however, requires a modified notion of mode decomposition, because the "modes" we generate temporally now are partially coherent in and of themselves. That is, the resulting partially coherent beam we generate has mutual

intensity of the following form:

$$J(\mathbf{r}_1, \mathbf{r}_2) = \sum_n \lambda_n \phi_n^*(\mathbf{r}_1)\phi_n(\mathbf{r}_2)\mu_n(\mathbf{r}_1 - \mathbf{r}_2) \tag{4.53}$$

where $\mu_n(\mathbf{r}_1 - \mathbf{r}_2)$ is the reduction in degree of coherence calculated using the van Cittert-Zernike theorem from the extended source, i.e. the Fourier transform of the extended source intensity profile. Each mode is the coupling of a fully coherent pattern $\phi_n$ with a spatial coherence falloff function $\mu_n$. Due to the similarity of this expression to that of Schell-model sources, let's call each mode a *quasi-Schell mode*, since instead of positive valued separable function that is symmetric, we have a complex valued separable function that is Hermitian. Such a decomposition, which we'll call a *quasi-Schell mode decomposition*, may not always be possible, because $J(\mathbf{r}_1, \mathbf{r}_2)/\mu(\mathbf{r}_1 - \mathbf{r}_2)$ may not be positive definite; in the discrete case, we are scaling down the diagonal values while boosting the off-diagonal values, which has a tendency to make eigenvalues go negative. However, if such a decomposition is possible, then intuitively, we should be able to reduce the number of modes needed to approximate the desired partially coherent beam within some error tolerance. For example, incoherent compositions of Gaussian-Schell model beams can be trivially generated by iterating between the different beams.

In practice, there is another way we can utilize this decomposition. If we consider the case where we have the desired mutual intensity $J$, then the default way to reduce modes from $K$ (the rank of the matrix $J$) modes to say $K' < K$ modes is to compute the factorization $J = \hat{U}\hat{U}^H$ from the SVD and then take the columns in $\hat{U}$ corresponding to the $K'$ largest singular values. This does result in the least squared error in $J$, but that might not be what we desire. For example, it might introduce unwanted high frequency details. This would simply be a case of trying to specify the entire mutual intensity matrix with too few values.

What we can do here is to accept a compromise. Let us intentionally reduce the spatial coherence area of the desired partially coherent beam at the SLM plane. That

is, instead of trying to generate the desired matrix $J$, we generate a matrix $\hat{J}$:

$$\hat{J} = J \odot M \tag{4.54}$$

where $M$ is a matrix where each entry is given by:

$$m_{i,j} = f\left(\sqrt{(x_i - x_j)^2 + (y_i - y_j)^2}\right) \tag{4.55}$$

where $f(r)$ is a positive-valued function with maximum value 1 at $r = 0$ and drops off for increasing $r$, and $(x_i, y_i)$, $(x_j, y_j)$ are the spatial coordinates corresponding to indexes $i$ and $j$ respectively. For example, if $f(r)$ is a delta function, then $M$ would be a diagonal matrix. In general, $M$ should approximate a banded matrix, and thus we are removing "outer" areas of the matrix $J$. That is, $M$ here takes the place of our $\mu(\mathbf{r}_1, \mathbf{r}_2)$ value in Equation (4.53). Obviously, the issue presented earlier about loss of positive definiteness doesn't apply, because we applied $M$ to our desired $J$ and dividing element-wise by $M$ would result in our original $J$, which was positive definite to begin with. Note that reducing the spatial coherence area at the SLM plane reduces the effective resolution off the SLM plane.

We can then seek to find an optimal $U$ such that the following quantity is minimized:

$$M \odot (J - UU^H) = \hat{J} - M \odot (UU^H) \tag{4.56}$$

Here, what we are doing is that we are masking out unimportant elements of $J$ using the $M$ matrix, reducing the number of degrees of freedom of control required. Therefore, we should need less columns in $U$ to retain the same amount of error compared to the case when we are not masking by $M$. By finding a low rank approximation this way instead of a standard SVD, we are trading off resolution away from the primal plane ($\Pi_0$) for a smaller set of modes, making the error introduced in the low rank approximation "nicer".

The mathematical form of this problem has actually been investigated in the filter design and data-mining literature. Specifically, the following optimization problem we are trying to solve is called a *weighted low-rank approximation* problem:

Given a matrix $\hat{J} \in \mathbb{C}^{N^2 \times N^2}$, a matrix $M \in \mathbb{C}^{N^2 \times N^2}$ with non-negative elements and $K < N^2$, find a $U \in \mathbb{C}^{N^2 \times K}$ to minimize:

$$\left| \hat{J} - M \odot (UU^H) \right|^F \qquad (4.57)$$

and there exists numerous papers on solving this problem [73–77].

As a simple proof of concept, we will apply a global line search gradient descent algorithm to the problem and attempt to express the mutual intensity obtained previously in section 4.2 for the simulation of a scene with a plane of emitters and a planar occluder. Employing the idea of gradually changing the problem from the unweighted low rank approximation problem to the weighted low rank approximation problem [74], we will start our initial guess to be the solution to the unweighted low rank approximation and gradually modify the weighting during each iteration such that the net difference from the final weighting follows an exponential decay pattern. For a particular iteration, the direction of steepest descent for the current $U$ can be calculated by:

$$\Delta_U = 4(M \odot M \odot (J - UU^H))^H U \qquad (4.58)$$

For this particular example, $M$ was set to a mutual intensity that represented uniform illumination with spatial coherence drop-off equivalent to a Gaussian with standard deviation $60\lambda$. This algorithm was implemented in MATLAB and run for 1000 iterations on the mutual intensity from the scene simulation algorithm.

The resulting focal stacks and tilt-view images are shown in Figures 4.33 and 4.34 respectively. A set of images corresponding to the unaltered mutual intensity has been included in each figure as well for comparison.

As can be seen from the figures, only 16 quasi-Schell modes are required to recreate the scene with fidelity matching that of the 256-mode coherent mode decomposition. One way of viewing this result is that if a partially coherent field is fairly incoherent, then there is actually less information content present in this field. If carried to its logical conclusion, a fully incoherent field should have very few information, and this agrees with the fact that a fully incoherent field can be fully characterized by a

Figure 4.33: The resulting focal stack as a function of reducing the number of modes in the mutual intensity via quasi-Schell modes. $K$ for each row after the first row denotes the number of modes kept, while the first row corresponds to the unaltered mutual intensity function and is shown for comparison. The asterisk after the number indicates that a quasi-Schell mode decomposition had been approximated through the weighted low-rank approximation algorithm.
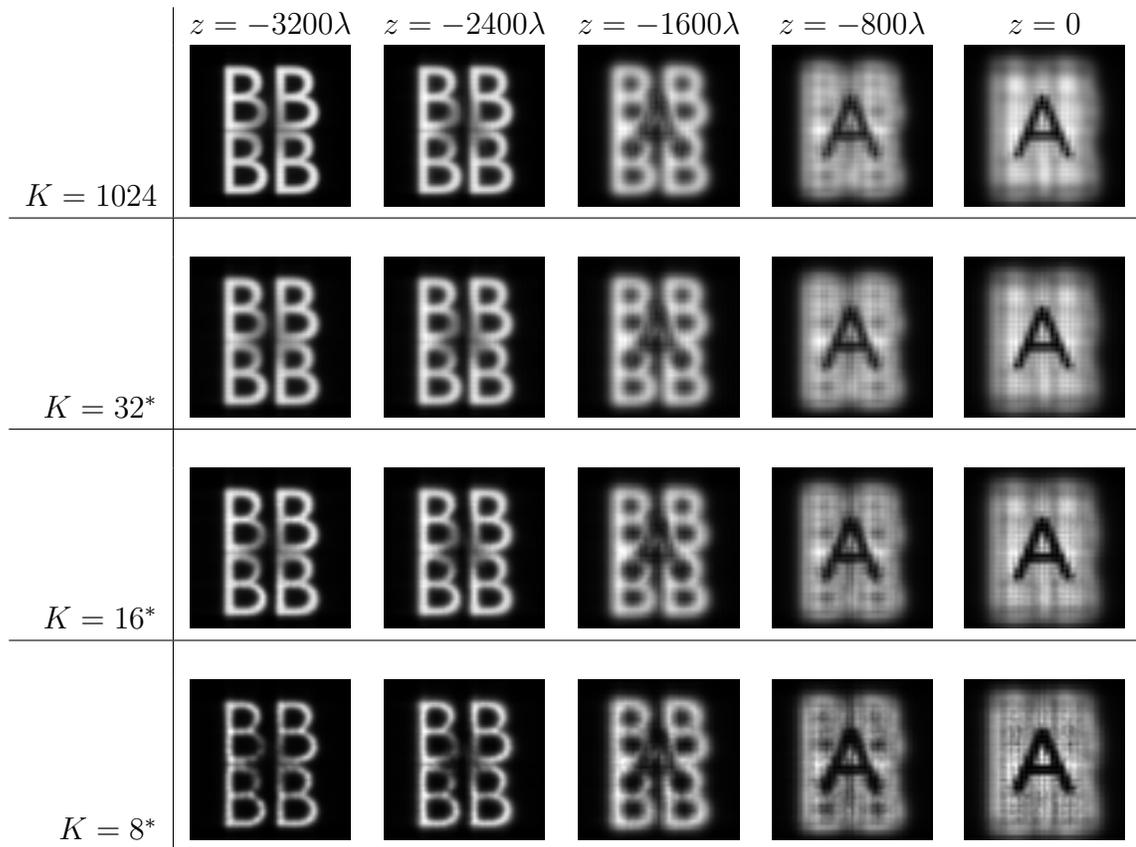
Figure 4.34: The resulting tilt-view images as a function of reducing the number of modes in the mutual intensity via quasi-Schell modes. $K$ for each row after the first row denotes the number of modes kept, while the first row corresponds to the unaltered mutual intensity function and is shown for comparison. The asterisk after the number indicates that a quasi-Schell mode decomposition had been approximated through the weighted low-rank approximation algorithm.

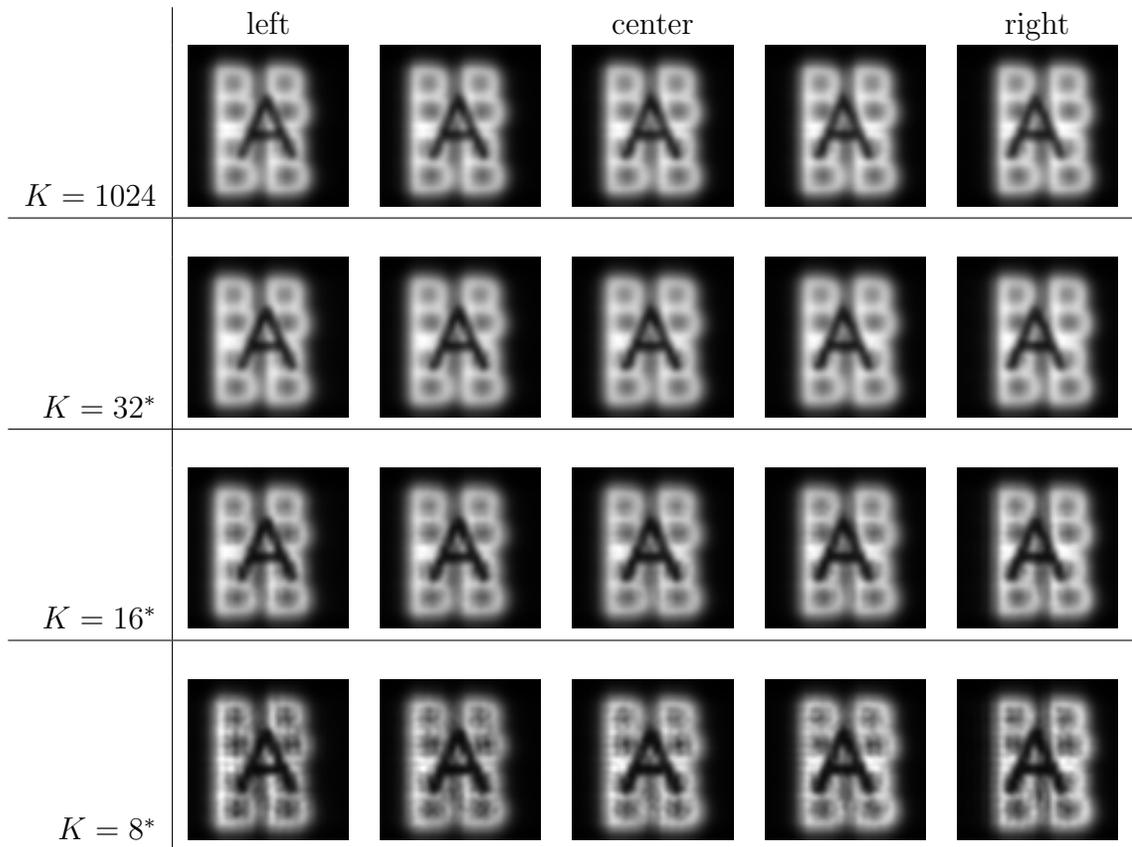two-dimensional intensity function. Therefore, the use of quasi-Schell modes in the generation of arbitrary partially coherent fields is a very useful tool and should be studied more extensively as a new avenue of attack on the problem of illumination generation.

# Chapter 5

# Conclusion

It is evident that partial coherence and its measures are a powerful tool in both analysis and synthesis of three-dimensional illumination patterns. Analysis based on the mutual intensity and phase space representations shows limitations on feasible illumination patterns for existing families of illumination devices. Ray-based devices are affected by the uncertainty principle and thus suffer from overall lack of resolution. Holographic devices suffer from limitations of full coherence and thus also cannot generate specific patterns, even two-dimensional intensity patterns. Volumetric devices effectively create point emitters in space and cannot create astigmatic and occlusion effects. None of the device families can generate arbitrary mutual intensity patterns. These limitations lead naturally to the idea of creating more general patterns of limitation based on a desired mutual intensity, instead of the ray-space, coherent field and point cloud representations used by these device families.

Partially coherent fields can be created through temporal multiplexing of coherence modes of a desired mutual intensity. The light generation hardware would use a spatial light modulator setup to modulate the phase and amplitude of a coherent laser beam, and patterns corresponding to coherence modes would be rapidly iterated. The desired mutual intensity would be computed via various algorithms depending on the application. This effectively moves the application-specific aspect of illumination device design into the software regime, allowing for faster iteration and more flexibility.

One basic application of illumination generation is to reproduce virtually the light emitted from a custom scene. The desired mutual intensity in this case can be computed through simulation of the propagation of light in a simple scene. The simulation is performed by computing a coherence mode for each emitter in the scene and consolidating if the number of emitters is too large. A low rank approximation of the computed mutual intensity can be used to reduce the number of patterns needed to be rapidly interleaved at the SLM. Tests using a simple scene consisting of a plane of emitters followed by a planar occluder shows acceptable results while retaining only a quarter of the original modes.

Another basic application of illumination generation is the generation of some desired voxel intensity pattern in space specified as a series of two-dimensional intensity patterns on transverse planes. A nonlinear conjugate gradients algorithm with global line search has been developed which attempts to minimize the least squares error in intensity across the entire volume. This algorithm was used for an exploration of a three-dimensional intensity pattern created from a Gaussian beam with various amounts of longitudinal compression/expansion, and it demonstrates that partial coherence aids greatly in generating a compressed beam whereas it has much less impact in the case of an expanded (almost propagation invariant) beam. As a further test of the algorithm, a partially coherent beam was generated to attempt to create three specific images at different depths simultaneously using this algorithm, and increasing the number of degrees of freedom through increasing the number of modes results in better performance than increasing the number of degrees of freedom by increasing the number of pixels on the SLM. Results demonstrate that using partial coherence over full coherence results in much higher image quality and that a nonlinear conjugate gradients algorithm can aim specifically for a weighted least squares error, whereas a iterated projections algorithms tends to reduce error more in low-intensity regions.

Lastly, it appears that the computation of what specific sequence of SLM patterns to iterate temporally in the generation of a partially coherent beam is not a trivial problem in that it should be theoretically possible to optimize for specific niceties in the mode patterns such as less wasted optical energy. Furthermore, to reduce the

number of modes needed to create a pattern, it might be possible to use a partially coherent Schell-model source instead of a fully coherent source to illuminate the SLM. This effectively creates a partially coherent field through the accumulation of partially coherent "quasi-Schell modes" instead of fully coherent modes. This trades off off-primal plane resolution for lower modes, instead of increasing error throughout the entire volume. Furthermore, rather incoherent fields can be represented using fewer modes with few losses.

Approaching the illumination problem through the synthesis of partially coherent beams is feasible and should be a fruitful area of research. The algorithm presented for generating a mutual intensity from desired transverse intensity patterns can be improved and studied more for convergence properties, and other algorithms based on different constraints can definitely be derived. Furthermore, derivation of an "optimal" temporal sequence of modes from a desired mutual intensity would be an interesting problem to tackle in terms of efficiency and/or pixel value quantization error minimization. Mode count reduction through the use of partially coherent source illumination and the idea of a quasi-Schell mode decomposition are also viable avenues of research. Lastly, all of these results have been derived through simulation; it would be useful to see how well these algorithms and methods work in a practical setting and to see if there are other optimality criteria needed for better performance. With further research, partial coherence should give the artist, the scientist and the engineer a powerful tool in the creation of illumination.

# Bibliography

[1] M. Levoy, B. Chen, V. Vaish, M. Horowitz, I. McDowall, and M. Bolas, "Synthetic aperture confocal imaging," in *Proc. ACM SIGGRAPH*, 2004.

[2] M. Levoy, Z. Zhang, and I. McDowall, "Recording and controlling the 4d light field in a microscope," *Journal of Microscopy*, vol. 235, pp. 144–162, 2009.

[3] A. Ashkin, "Acceleration and trapping of particles by radiation pressure," *Phys. Rev. Lett.*, vol. 24, no. 4, pp. 156–159, Jan 1970.

[4] A. Ashkin, J. M. Dziedzic, J. E. Bjorkholm, and S. Chu, "Observation of a single-beam gradient force optical trap for dielectric particles," *Opt. Lett.*, vol. 11, no. 5, pp. 288–290, May 1986.

[5] G. Nagel, D. Ollig, M. Fuhrmann, S. Kateriya, A. M. Musti, E. Bamberg, and P. Hegemann, "Channelrhodopsin-1: A light-gated proton channel in green algae," *Science*, vol. 296, no. 5577, pp. 2395–2398, 2002.

[6] G. Nagel, T. Szellas, W. Huhn, S. Kateriya, N. Adeishvili, P. Berthold, D. Ollig, P. Hegemann, and E. Bamberg, "Channelrhodopsin-2, a directly light-gated cation-selective membrane channel," *Proceedings of the National Academy of Sciences*, vol. 100, no. 24, pp. 13 940–13 945, 2003.

[7] E. S. Boyden, F. Zhang, E. Bamberg, G. Nagel, and K. Deisseroth, "Millisecond-timescale, genetically targeted optical control of neural activity," *Nature Neuroscience*, vol. 8, no. 9, pp. 1263–1268, Sep 2005.

[8] X. Li, D. V. Gutierrez, M. G. Hanson, J. Han, M. D. Mark, H. Chiel, P. Hegemann, L. T. Landmesser, and S. Herlitze, "Fast noninvasive activation and inhibition of neural and network activity by vertebrate rhodopsin and green algae channelrhodopsin," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 102, no. 49, pp. 17 816–17 821, 2005. [Online]. Available: http://www.pnas.org/content/102/49/17816.abstract

[9] G. Nagel, M. Brauner, J. F. Liewald, N. Adeishvili, E. Bamberg, and A. Gottschalk, "Light activation of channelrhodopsin-2 in excitable cells of caenorhabditis elegans triggers rapid behavioral responses," *Current Biology*, vol. 15, no. 24, pp. 2279 – 2284, 2005. [Online]. Available: http://www.sciencedirect.com/science/article/pii/S0960982205014077

[10] K. Jain, C. Willson, and B. Lin, "Ultrafast deep uv lithography with excimer lasers," *Electron Device Letters, IEEE*, vol. 3, no. 3, pp. 53 – 55, Mar 1982.

[11] K. Jain, *Excimer Laser Lithography*. Bellingham, WA: SPIE Press, 1990.

[12] B. J. Lin, *Optical Lithography*. Bellingham, WA: SPIE Press, 2009.

[13] A. Gershun, "The light field (translated by P. Moon and G. Timoshenko)," *J. Math. and Physics*, vol. 18, pp. 51–151, 1939.

[14] M. Levoy and P. Hanrahan, "Light field rendering," in *Proc. ACM SIGGRAPH*, 1996.

[15] A. R. L. Travis, "Autostereoscopic 3-d display," *Applied Optics*, vol. 29, no. 29, pp. 4341–4342, Oct 1990.

[16] G. Lippmann, "Epreuves reversible donnant la sensation du relief," *J. Phys.*, vol. 7, pp. 821–825, 1908.

[17] D. Gabor, "Microscopy by reconstructed wave-fronts," *Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences*, vol. 197, no. 1051, pp. 454–487, 1949. [Online]. Available: http://rspa.royalsocietypublishing.org/content/197/1051/454.abstract

[18] B. R. Brown and A. W. Lohmann, "Complex spatial filtering with binary masks," *Appl. Opt.*, vol. 5, no. 6, pp. 967–969, Jun 1966. [Online]. Available: http://ao.osa.org/abstract.cfm?URI=ao-5-6-967

[19] J. P. Waters, "Holographic image synthesis utilizing theoretical methods," *Appl. Phys. Lett.*, vol. 9, pp. 405–407, dec 1966.

[20] B. R. Brown and A. W. Lohmann, "Computer-generated binary holograms," *IBM Journal of Research and Development*, vol. 13, no. 2, pp. 160–168, Mar 1969.

[21] C. Slinger, C. Cameron, and M. Stanley, "Computer-generated holography as a generic display technology," *IEEE Computer*, vol. 38, pp. 46–53, Aug 2005.

[22] J. D. Lewis, C. M. Verber, and R. B. McGhee, "A true three-dimensional display," *Electron Devices, IEEE Transactions on*, vol. 18, no. 9, pp. 724 – 732, Sep 1971.

[23] E. Parker and P. Wallis, "Three-dimensional cathode-ray tube displays," *Electrical Engineers - Part III: Radio and Communication Engineering, Journal of the Institution of*, vol. 95, no. 37, pp. 371–387, Sep 1948.

[24] A. C. Traub, "Stereoscopic display using rapid varifocal mirror oscillations," *Applied Optics*, vol. 6, pp. 1085–1087, Jun 1967.

[25] Z. Zhang and M. Levoy, "Wigner distributions and how they relate to the light field," in *Computational Photography (ICCP), 2009 IEEE International Conference on*, april 2009, pp. 1–10.

[26] Z. Zhang, G. Barbastathis, and M. Levoy, "Limitations of coherent computer generated holograms," in *Digital Holography and Three-Dimensional Imaging*. Optical Society of America, 2011, p. DTuB5. [Online]. Available: http://www.opticsinfobase.org/abstract.cfm?URI=DH-2011-DTuB5

[27] P. Langevin, *Le Radium*, vol. 7, pp. 249–261, 1910.

[28] M. Born, "Electron theory of natural optic rotation processes in isotropic and anisotropic liquids," *Annalen der Physik*, vol. 55, no. 3, pp. 177–240, May 1918.

[29] A. D. Buckingham and J. A. Pople, "Theoretical studies of the kerr effect i: Deviations from a linear polarization law," *Proceedings of the Physical Society. Section A*, vol. 68, no. 10, p. 905, 1955. [Online]. Available: http://stacks.iop.org/0370-1298/68/i=10/a=307

[30] M. Göppert-Mayer, "Über Elementarakte mit zwei Quantensprüngen," *Annalen der Physik*, vol. 401, pp. 273–294, 1931.

[31] W. Kaiser and C. G. B. Garrett, "Two-photon excitation in Ca$F2$: $Eu2+$," *Phys. Rev. Lett.*, vol. 7, no. 6, pp. 229–231, Sep 1961.

[32] L. Mandel and E. Wolf, *Optical Coherence and Quantum Optics*. Cambridge University Press, Sep. 1995. [Online]. Available: http://www.worldcat.org/isbn/0521417112

[33] E. Wolf, "New theory of partial coherence in the space-frequency domain. part i: spectra and cross spectra of steady-state sources," *J. Opt. Soc. Am.*, vol. 72, no. 3, pp. 343–351, 1982. [Online]. Available: http://www.opticsinfobase.org/abstract.cfm?URI=josa-72-3-343

[34] H. Gamo, "Intensity matrix and degree of coherence," *J. Opt. Soc. Am.*, vol. 47, no. 10, pp. 976–976, Oct. 1957.

[35] H. M. Ozaktas, S. Yüksel, and M. A. Kutay, "Linear algebraic theory of partial coherence: discrete fields and measures of partial coherence," *J. Opt. Soc. Am. A*, vol. 19, no. 8, pp. 1563–1571, 2002. [Online]. Available: http://josaa.osa.org/abstract.cfm?URI=josaa-19-8-1563

[36] F. Hlawatsch and G. F. Boudreaux-Bartels, "Linear and quadratic time-frequency signal representations," *IEEE Signal Process. Mag.*, vol. 9, no. 2, pp. 21–67, Apr. 1992.

[37] K.-H. Brenner and J. Ojeda-Castañeda, "Ambiguity function and wigner distribution function applied to partially coherent imagery," *Journal of Modern Optics*, vol. 31, no. 2, pp. 213–223, Feb. 1984.

[38] A. Walther, "Radiometry and coherence," *J. Opt. Soc. Am.*, vol. 58, no. 9, pp. 1256–1259, Sep 1968.

[39] E. W. Marchand and E. Wolf, "Walther's definitions of generalized radiance," *J. Opt. Soc. Am.*, vol. 64, no. 9, pp. 1273–1274, Sep. 1974.

[40] A. T. Friberg, "On the existence of a radiance function for finite planar sources of arbitrary states of coherence," *J. Opt. Soc. Am.*, vol. 69, no. 1, pp. 192–198, Jan. 1979.

[41] J. T. Foley and E. Wolf, "Radiometry as a short-wavelength limit of statistical wave theory with globally incoherent sources," *Opt. Commun.*, vol. 55, no. 4, pp. 236–241, Sep. 1985.

[42] K. Kim and E. Wolf, "Propagation law for Walther's first generalized radiance function and its short-wavelength limit with quasi-homogeneous sources," *J. Opt. Soc. Am. A*, vol. 4, no. 7, pp. 1233–1236, Jul. 1987.

[43] M. J. Bastiaans, "The wigner distribution function applied to optical signals and systems," *Opt. Commun.*, vol. 25, no. 1, pp. 26–30, Apr. 1978.

[44] A. Papoulis, "Ambiguity function in fourier optics," *J. Opt. Soc. Am.*, vol. 64, no. 6, pp. 779–788, Jun. 1974.

[45] K.-H. Brenner, A. W. Lohmann, and J. Ojeda-Castañeda, "The ambiguity function as a polar display of the OTF," *Opt. Commun.*, vol. 44, no. 5, pp. 323–326, Feb. 1983.

[46] E. R. Dowski, Jr. and W. T. Cathey, "Extended depth of field through wave-front coding," *Applied Optics*, vol. 34, no. 11, pp. 1859–1866, Apr. 1995.

[47] J. W. Goodman, *Introduction to Fourier Optics, 3rd. ed.* Greenwood Village, CO: Roberts and Company Publishers, 2004.

[48] E. H. Adelson and J. Y. Wang, "Single lens stereo with a plenoptic camera," *IEEE Trans. PAMI*, vol. 14, no. 2, pp. 99–106, Feb. 1992.

[49] R. Ng, "Fourier slice photography," in *Proc. ACM SIGGRAPH*, 2005.

[50] G. W. Farnell, "On the axial phase anomaly for microwave lenses," *Journal of the Optical Society of America*, vol. 48, no. 9, pp. 643–647, Sep 1958.

[51] Y. Li, "Dependence of the focal shift on fresnel number and f number," *J. Opt. Soc. Am.*, vol. 72, no. 6, pp. 770–774, Jun 1982.

[52] C. J. R. Sheppard and P. T or ok, "Dependence of focal shift on fresnel number and angular aperture," *Optics Letters*, vol. 23, no. 23, pp. 1803–1804, Dec 1998.

[53] T. Bishop, S. Zanetti, and P. Favaro, "Light field superresolution," in *Computational Photography (ICCP), 2009 IEEE International Conference on*, april 2009, pp. 1–9.

[54] T. Georgiev and A. Lumsdaine, "Superresolution with plenoptic camera 2.0," Adobe Systems Incorporated, Tech. Rep., 2009.

[55] A. C. Schell, "The multiple plate antenna," Ph.D. dissertation, Massachusetts Institute of Technology, 1961.

[56] A. S. Marathay, L. Heiko, and J. L. Zuckerman, "Study of rough surfaces by light scattering," *Applied Optics*, vol. 9, no. 11, pp. 2470–2476, Nov 1970.

[57] J. T. Foley and M. S. Zubairy, "The directionality of gaussian schell-model beams," *Optics Communications*, vol. 26, no. 3, pp. 297–300, Sep 1978.

[58] A. W. Lohman, "Three-dimensional properties of wave-fields," *Optik*, vol. 51, no. 2, pp. 105–117, 1978.

[59] R. Piestun, B. Spektor, and J. Shamir, "Wave fields in three dimensions: analysis and synthesis," *J. Opt. Soc. Am. A*, vol. 13, no. 9, September 1996.

[60] A. Walther, "The question of phase retrieval in optics," *Journal of Modern Optics*, vol. 10, pp. 41–49, 1963.

[61] P. De Santis, F. Gori, G. Guattari, and C. Palma, "Synthesis of partially coherent fields," *J. Opt. Soc. Am. A*, vol. 3, no. 8, pp. 1258–1262, 1986. [Online]. Available: http://josaa.osa.org/abstract.cfm?URI=josaa-3-8-1258

[62] M. A. A. Neil, T. Wilson, and R. Juškaitis, "A wavefront generator for complex pupil function synthesis and point spread function engineering," *Journal of Microscopy*, vol. 197, no. 3, pp. 219–223, Mar 2000.

[63] J. N. Mait and K.-H. Brenner, "Dual-phase holograms: improved design," *Applied Optics*, vol. 26, no. 22, pp. 4883–4892, Nov 1987.

[64] C. K. Hsueh and A. A. Sawchuk, "Computer-generated double-phase holograms," *Applied Optics*, vol. 17, no. 24, pp. 3874–3883, Dec 1978.

[65] H. Gross, "Numerical propagation of partially coherent laser beams through optical system," *Optics & Laser Technology*, vol. 29, no. 5, pp. 257–260, July 1997.

[66] C. Rydberg and J. Bengtsson, "Efficient numerical representation of the optical field for the propagation of partially coherent radiation with a specified spatial and temporal coherence function," *J. Opt. Soc. Am. A*, vol. 23, no. 7, pp. 1616–1625, Jul 2006.

[67] M. J. D. Powell, "Convergence properties of algorithms for nonlinear optimization," *SIAM Review*, vol. 28, no. 4, pp. pp. 487–500, 1986. [Online]. Available: http://www.jstor.org/stable/2031100

[68] E. Polak and G. Ribiere, "Note sur la convergence de methodes de directions conjugées," *Revue Française d'Informatique et de Recherche Opérationnelle, Série Rouge*, vol. 3, no. 16, pp. 35–43, 1969.

[69] R. Piestun and J. Shamir, "Synthesis of three-dimensional light fields and applications," *Proceedings of the IEEE*, vol. 90, no. 2, pp. 222 –244, feb 2002.

[70] C. Rydberg and J. Bengtsson, "Numerical algorithm for the retrieval of spatial coherence properties of partially coherent beams from transverse intensity measurements," *Opt. Express*, vol. 15, no. 21, pp. 13 613–13 623, 2007. [Online]. Available: http://www.opticsexpress.org/abstract.cfm?URI=oe-15-21-13613

[71] P. L. Combettes, "Inconsistent signal feasibility problems: least-squares solutions in a product space," *Signal Processing, IEEE Transactions on*, vol. 42, no. 11, pp. 2955 –2966, Nov. 1994.

[72] A. Starikov and E. Wolf, "Coherent-mode representation of gaussian schell-model sources and of their radiation fields," *J. Opt. Soc. Am.*, vol. 72, no. 7, pp. 923–928, Jul 1982.

[73] W.-S. Lu, S.-C. Pei, and P.-H. Wang, "Weighted low-rank approximation of general complex matrices and its application in the design of 2-d digital filters," *IEEE Trans. Circuits Syst. I*, vol. 44, no. 7, Jul 1997.

[74] W.-S. Lu and A. Antoniou, "New method for weighted low-rank approximation of complex-valued matrices and its application for the design of 2-d digital filters," in *Circuits and Systems, 2003. ISCAS '03. Proceedings of the 2003 International Symposium on*, vol. 3, May 2003, pp. III–694 – III–697 vol.3.

[75] J. H. Manton, R. Mahony, and Y. Hua, "The geometry of weighted low-rank approximations," *IEEE Trans. Signal Process.*, vol. 51, no. 2, Feb 2003.

[76] N. Srebro and T. Jaakkola, "Weighted low-rank approximations," in *ICML*, T. Fawcett and N. Mishra, Eds. AAAI Press, 2003, pp. 720–727.

[77] K. Werner and M. Jansson, "Weighted low rank approximation and reduced rank linear regression," in *International Conference on Acoustics, Speech, and Signal Processing*, 2004.