# Design, Realization, and Applications of 3D Multifunctional Nanophotonics

Thesis by
Gregory Roberts

In Partial Fulfillment of the Requirements for the
Degree of
Doctor of Philosophy

## Caltech

CALIFORNIA INSTITUTE OF TECHNOLOGY
Pasadena, California

2024
Defended August 10, 2023

© 2024

Gregory Roberts
ORCID: 0009-0002-0720-3938

# ACKNOWLEDGEMENTS

I will bring everything around with one last mention about my family. The last couple years of my PhD I had the privilege of spending most of my weekends with Erin and Dan where we got up to great adventures - marathons, winery bike tours, half Ironman triathlons, and surfing, to name a few. And everyone else in the family always went to great efforts to visit or host me visiting them throughout my whole time here. Camille could often be found working remotely down in the Newport Beach clubhouse and on multiple occasions Steve would orchestrate a covert surprise visit. Mom and Dad would come regularly as well, always excited to ask how life was going and find new places to get ice cream and coffee. To Mom, Dad, Steve, Erin, Camille, Dan, Sylvia and newest additions Charlie and Rae, thank you for being the most supportive, loving, and inspirational people to be around. I am excited for many adventures to come.

# ABSTRACT

Metaoptics leverages electromagnetic phenomena and the advanced abilities of modern nanofabrication to replicate traditional optical devices in a fraction of the thickness and to realize novel, compact, multifunctional devices with no known bulk equivalent. Motivated by the expanding role of optics in modern technologies, this field has seen a rise in design techniques for realizing increasingly powerful photonic structures. Three-dimensional (3D) devices, with refractive index distributions patterned at subwavelength scales, represent an enormous design space capable of achieving highly efficient, free space, multifunctional structures. By utilizing a gradient-based, iterative optimization algorithm, a technique for nanophotonic inverse design, we demonstrate scattering structures with unique responses to all the fundamental properties of light. The algorithm is constrained such that resulting devices can be made with realistic multilayer fabrication processes. We present dielectric structures that can be placed directly on top of image sensor arrays and sort light to different pixels based on its wavelength, polarization, and angular momentum, thus enabling efficient and exotic camera technologies. The following work contains fabrication and measurement of 3D devices in the mid-infrared, practical evaluations of devices for visible light imaging applications, and visualizations of underlying structure of photonic design optimization problems.

# PUBLISHED CONTENT AND CONTRIBUTIONS

[1]C. Ballew, G. Roberts, and A. Faraon, "Multi-dimensional wavefront sensing using volumetric meta-optics", arXiv preprint arXiv:2301.10346, `10.48550/arXiv.2301.10346` (2023) `10.48550/arXiv.2301.10346`,
G.R. consulted on the design of the device and its evaluation as well as provided code for performing the level-set optimization used in the manuscript. G.R. gave feedback on the writing of the manuscript and analysis of the results.

[2]C. Ballew, G. Roberts, T. Zheng, and A. Faraon, "Constraining continuous topology optimizations to discrete solutions for photonic applications", ACS photonics **10**, 836–844 (2023) `10.1021/acsphotonics.2c00862`,
G.R. consulted on the design of the main algorithm used in the manuscript. G.R. gave feedback on the writing of the manuscript.

[3]G. Roberts, C. Ballew, T. Zheng, J. C. Garcia, S. Camayd-Muñoz, P. W. C. Hon, and A. Faraon, "3D-patterned inverse-designed mid-infrared metaoptics", Nature Communications **14**, 2768 (2023) `10.1038/s41467-023-38258-2`,
G.R. helped conceive the project, carried out optimization, fabrication, and measurement of devices in the manuscript with input from other authors. G.R. prepared the manuscript with input from all authors.

[4]C. Ballew, G. Roberts, S. Camayd-Muñoz, M. F. Debbas, and A. Faraon, "Mechanically reconfigurable multi-functional meta-optics studied at microwave frequencies", Scientific reports **11**, 1–9 (2021) `10.1038/s41598-021-88785-5`,
G.R. consulted on the designs in the manuscripts. G.R. gave feedback on the writing of the manuscript.

[5]P. Camayd-Muñoz, C. Ballew, G. Roberts, and A. Faraon, "Multifunctional volumetric meta-optics for color and polarization image sensors", Optica **7**, 280–283 (2020) `10.1364/OPTICA.384228`,
G.R. carried out device optimizations and production of data for Figure 2. G.R. gave feedback on the writing of the manuscript.

# TABLE OF CONTENTS

# LIST OF ILLUSTRATIONS

# LIST OF TABLES

*Chapter 1*

# INTRODUCTION

The picture most often associated with the word 'optics' is that of a lens, a piece of material that can direct the path of rays of light depending on the shape of its surface. This is due to the phenomenon of refraction, described by Snell's Law [1] which says that at the interface between two materials, a ray of light will change direction depending on its incident angle and the index of refraction of each of the two materials. Depending on the desired optical function, one can tailor the choice of materials as well as the shapes of their surfaces to create refractive optics. This technique is used in almost every imaging technology. Modern mobile phone cameras, for example, form images using a series of low-cost plastic lenses created via injection molding [2], the surface profiles of which are optimized to create high quality images.

In the field of diffractive optics, there is another way to think about the effect of an optical element like a lens. Instead of picturing light traveling through a system as consisting of rays traveling in different directions, one can instead think of light as a wave. It is a wave in the electric and magnetic fields, each of which induce the other to allow light to propagate even free of a physical medium. The oscillating fields can be modeled as complex-valued where the field at each point has both an amplitude and a phase. To connect back to the previous picture, the local gradient in the phase defines the propagation direction of the wavefront or, in other words, the direction we would say a ray of light at that point is traveling [3]. Instead of redirecting this ray of light with refraction, we can, instead, directly manipulate the phase at each point on the wavefront to change its direction. This can be done with the same material that was used to make the refractive lens. The index of refraction specifies how fast light moves in a medium. Light in a larger index of refraction medium moves slower. Since the oscillation frequency is unchanged, the wavelength of light is shorter inside that medium. When the wave travels one wavelength it must have accumulated $2\pi$ phase so a shorter wavelength directly means a faster rate of phase accumulation per unit length. This, then, gives a way to achieve the diffractive lens. Using different thicknesses of material at different points on the wavefront, one can directly transform an incoming planar wavefront to a converging outgoing wave that focuses to a point after some propagation. Diffractive optics operate in this way,

interacting directly with the wavelike nature of light to achieve optical functions.

Metasurfaces build from this idea and can realize diffractive optics in ultra-thin form factors [4]. To achieve this, material is patterned at fine spatial scales below the wavelength of light. Material platforms and physical phenomena differ between approaches, but one example is the use of high index contrast dielectric posts, which confine incoming light tightly in there cores. Even at subwavelength thicknesses these posts can realize different phase delays by simply modifying a geometric parameter such as their radius [5]. These posts can be created with modern nanofabrication techniques, the same ones that are used to make advanced integrated circuits. By modeling the interaction of light with metasurface unit cells, traditionally bulky optics can be made effectively flat.

Metasurfaces provide a beautiful abstraction for optics, they separate the functionality of the optic from the underlying material interaction. With this platform, not only can bulk optics be made thinner and some of them combined into single pieces, but new functionality can be realized. For example, orthogonal polarizations can be controlled independently by the optic by creating birefringent design libraries. Posts that are not azimuthally symmetric can then impart different phase functions to different polarizations [5]. However, despite the fascinating research in this field, metasurfaces have typically been limited to a thin single layer of material. Eventually, this limited design space and propagation length becomes prohibitive, especially when designing multifunctional optics. Even simple lenses are required to be multifunctional in that they need to perform the same focusing function regardless of input wavelength over appreciable bandwidths. Achromatic functionality over broad operating bandwidths is one challenge for the metasurface research community [6]. Further, many optical systems contain numerous elements and perform complex tasks on every degree of freedom of light, not just the spatial, but the spectral and polarization. It is known that for certain functionalities, there are minimum necessary thicknesses [7]. The motivating question in this thesis is what can be done if we add some thickness back to metasurfaces? This may seem an innocuous question, one that is a natural extension of metasurface research. However, it has been mostly unexplored in the realm of free space optics especially when it comes to experimental demonstration. This is primarily due to two major challenges: (1) the design assumptions from metasurfaces quickly disappear with additional layers of patterning (even just by extending to two layers) and 3D metaoptic designs require the simulation and optimization of devices with full field solvers that account for

multiple scattering phenomena thus making their design challenging and computationally costly and (2) the fabrication of more volumetric devices requires multilayer processes with deeply subwavelength resolution and layer-to-layer alignment, the likes of which are standard in industry yet still highly challenging to control in academic cleanrooms.

In the following work, solutions to these two problems are presented and we aim to continue to build foundations in this new field of volumetric metaoptics. In Chapter 2, we lay the groundwork for the first of the two challenges. The photonic design problem of optimizing devices with multiple scattering between their component parts was built up in the research community for on chip devices where light travels along the surface of chips through waveguides and can interact with effectively thick structures that are patterned in just a single layer. We describe the theory behind how this method works and show simplified, example devices that achieve interesting applications to illustrate its power. In Chapter 3, we continue to the fully 3D realm for the design technique. Here, we also show the ability to fabricate and measure devices in the mid-infrared spectrum, thus addressing the second of the two major challenges. In Chapter 4, we explore deeply in simulation the design of volumetric devices for color image sensors that could theoretically improve camera transmission efficiencies by up to a factor of three. We propose a method for prototype fabrication that could be done with the help of an electronics foundry. Finally, in Chapter 5, we connect the photonic device design to optimizations in machine learning by adopting techniques for visualizing the design spaces we are exploring. This lays a framework for how to think about the dual-part problem of photonic optimization, one part being the laws of physics and the other the practical challenges of solving high-dimensional non-convex problems. Throughout this work we explore the design of several nanophotonic devices that can greatly improve the efficiency, compactness, and abilities of modern imaging systems.

## References

[1] M. Born and E. Wolf, "Principles of optics, 1964", New York: Mac-Millan (1959).

[2] A. W. McFarland and J. S. Colton, "Production and analysis of injection molded micro-optic components", Polymer Engineering & Science **44**, 564–579 (2004).

[3] J. W. Goodman, *Introduction to Fourier optics* (Roberts and Company publishers, 2005).

[4] H.-T. Chen, A. J. Taylor, and N. Yu, "A review of metasurfaces: physics and applications", Reports on progress in physics **79**, 076401 (2016).

[5]A. Arbabi, Y. Horie, M. Bagheri, and A. Faraon, "Dielectric metasurfaces for complete control of phase and polarization with subwavelength spatial resolution and high transmission", Nature nanotechnology **10**, 937–943 (2015).

[6]W. T. Chen, A. Y. Zhu, V. Sanjeev, M. Khorasaninejad, Z. Shi, E. Lee, and F. Capasso, "A broadband achromatic metalens for focusing and imaging in the visible", Nature nanotechnology **13**, 220–226 (2018).

[7]D. A. B. Miller, "Why optics needs thickness", Science **379**, 41–45 (2023).

*Chapter 2*

# ADJOINT METHOD FOR ELECTROMAGNETICS

## 2.1 Adjoint Method for Linear Optimization Problems

There are many good references on the adjoint method. The following explanation and derivations are based on [1] and [2]. The starting point for the adjoint method is a physical problem modeled as a linear system with a solution defined by:

$$\underline{\underline{A}}\mathbf{x} = \mathbf{b}$$
$$\underline{\underline{A}} \in \mathbb{C}^{n \times n} \tag{2.1}$$
$$\mathbf{x}, \mathbf{b} \in \mathbb{C}^{n}$$

The input to the system is given by the vector $\mathbf{b}$ and the output, or solution, by the vector $\mathbf{x}$. A solution to the problem can be written as $\mathbf{x} = \underline{\underline{A}}^{-1}\mathbf{b}$. For optimization based design problems, the physical model $\underline{\underline{A}}$ can be modified through some underlying variable that we, as the designer, have control over. Assume this variable is a vector $\mathbf{p}$ such that $\underline{\underline{A}} = \underline{\underline{A}}(\mathbf{p})$ and in turn $\mathbf{x} = \mathbf{x}(\mathbf{p})$.

In an optimization problem, we typically want the output $\mathbf{x}$ to satisfy as best as possible some property, defined by an objective function or figure of merit. For example, we may want to minimize the magnitude of the overlap of the output vector, $\mathbf{x}$, with some reference vector, $\mathbf{r}$, or in other words, minimize the following objective function:

$$g(\mathbf{x}) = ||\mathbf{x}^{\dagger}\mathbf{r}|| \tag{2.2}$$

where $\mathbf{x}^{\dagger}$ denotes the Hermitian adjoint of $\mathbf{x}$, since the above solution $\mathbf{x}$ exists in a complex vector space.

The matrix Equation (2.1) and its solution constitute a physical constraint for the problem and the objective defines the goal of the following optimization problem:

$$\min_{\mathbf{p}^{*}} g(\mathbf{x})$$
$$\text{subject to } \underline{\underline{A}}(\mathbf{p}^{*})\mathbf{x} = \mathbf{b} \tag{2.3}$$

where $\mathbf{p}^{*}$ is an optimal design choice for $\mathbf{p}$.

## 2.2 Adjoint Gradient for Local Optimizers

In many cases, the stated optimization problem in Equation (2.3) is nonlinear and the design vector $\mathbf{p}$ has a large dimension. In cases like this, a typical solution is to use a local optimizer that computes the gradient of the objective function $g$ with respect to the design vector $\mathbf{p}$. The simplest way to do this is via a brute force finite difference calculation. This can be done by perturbing every element of $\mathbf{p}$ individually, solving the linear system for each perturbation, and computing the resulting value of the objective function:

$$\frac{dg}{dp_k} \approx \frac{g(\mathbf{p} + h\hat{\mathbf{k}}) - g(\mathbf{p} - h\hat{\mathbf{k}})}{2h} \tag{2.4}$$

In cases where the linear system is costly to solve and $\mathbf{p}$ is high-dimensional, such as in electromagnetic inverse design, this technique is prohibitively expensive for computing the gradient. It also computes superfluous information by calculating the fields for every individual perturbation. There is a more efficient way to organize the problem of finding the gradient such that the number of times the linear system needs to be solved is constant and does not grow with the number of design parameters. This organization is called the adjoint method, which pre-dates the work in this thesis and has been derived in many places for various problems [1, 2].

First, we use the chain rule to pull apart individual implicit dependencies in the total derivative we are looking for. We note that there is only dependence of the objective function either explicitly on the parameter $\mathbf{p}$ or explicitly on the solution $\mathbf{x}$, which has a dependence on the parameter $\mathbf{p}$ through the solution matrix. Since $\mathbf{x}$ is in general complex, it has two independent variables in its real and imaginary part. This can be taken care of by using $\mathbf{x}$ and $\mathbf{x}^*$ as independent variables. For real-valued $\mathbf{p}$, $\frac{\partial \mathbf{x}^*}{\partial \mathbf{p}} = (\frac{\partial \mathbf{x}}{\partial \mathbf{p}})^*$, so the derivation will focus on just finding a solution for $\frac{\partial \mathbf{x}}{\partial \mathbf{p}}$.

$$\frac{dg}{d\mathbf{p}} = \frac{\partial g}{\partial \mathbf{p}} + \frac{\partial g}{\partial \mathbf{x}}\frac{\partial \mathbf{x}}{\partial \mathbf{p}} + \frac{\partial g}{\partial \mathbf{x}^*}\frac{\partial \mathbf{x}^*}{\partial \mathbf{p}} \tag{2.5}$$

Note that each of these terms are vectors. If $\mathbf{p}$ is $m$-dimensional (i.e., there are $m$ individual design variables), then $\frac{dg}{d\mathbf{p}}$ has dimensions $1 \times m$. Then, $\frac{\partial g}{\partial \mathbf{x}}$ has dimensions $1 \times n$. Therefore, $\frac{\partial \mathbf{x}}{\partial \mathbf{p}}$ is a matrix of size $n \times m$ as it stores the change for each of the $n$ entries of $\mathbf{x}$, given individual changes at each of the $m$ parameter values. The solution $\mathbf{x}$ is given by:

$$\mathbf{x} = \underline{\underline{\mathbf{A}}}(\mathbf{p})^{-1}\mathbf{b} \qquad (2.6)$$

Implicitly differentiating $\underline{\underline{\mathbf{A}}}\mathbf{x} = \mathbf{b}$ with respect to $\mathbf{p}$:

$$\frac{\partial \underline{\underline{\mathbf{A}}}}{\partial \mathbf{p}}\mathbf{x} + \underline{\underline{\mathbf{A}}}\frac{\partial \mathbf{x}}{\partial \mathbf{p}} = \frac{\partial \mathbf{b}}{\partial \mathbf{p}}$$

$$\underline{\underline{\mathbf{A}}}\frac{\partial \mathbf{x}}{\partial \mathbf{p}} = \frac{\partial \mathbf{b}}{\partial \mathbf{p}} - \frac{\partial \underline{\underline{\mathbf{A}}}}{\partial \mathbf{p}}\mathbf{x} \qquad (2.7)$$

$$\frac{\partial \mathbf{x}}{\partial \mathbf{p}} = \underline{\underline{\mathbf{A}}}^{-1}(\frac{\partial \mathbf{b}}{\partial \mathbf{p}} - \frac{\partial \underline{\underline{\mathbf{A}}}}{\partial \mathbf{p}}\mathbf{x})$$

Ultimately, the desired quantity is $\frac{\partial g}{\partial \mathbf{x}}\frac{\partial \mathbf{x}}{\partial \mathbf{p}}$, which looks like:

$$\frac{\partial g}{\partial \mathbf{x}}\frac{\partial \mathbf{x}}{\partial \mathbf{p}} = \frac{\partial g}{\partial \mathbf{x}}[\underline{\underline{\mathbf{A}}}^{-1}(\frac{\partial \mathbf{b}}{\partial \mathbf{p}} - \frac{\partial \underline{\underline{\mathbf{A}}}}{\partial \mathbf{p}}\mathbf{x})]$$

$$\frac{\partial g}{\partial \mathbf{x}}\frac{\partial \mathbf{x}}{\partial \mathbf{p}} = [\frac{\partial g}{\partial \mathbf{x}}\underline{\underline{\mathbf{A}}}^{-1}](\frac{\partial \mathbf{b}}{\partial \mathbf{p}} - \frac{\partial \underline{\underline{\mathbf{A}}}}{\partial \mathbf{p}}\mathbf{x}) \qquad (2.8)$$

$$\frac{\partial g}{\partial \mathbf{x}}\frac{\partial \mathbf{x}}{\partial \mathbf{p}} = [(\underline{\underline{\mathbf{A}}}^{\dagger})^{-1}\frac{\partial g}{\partial \mathbf{x}}^{\dagger}]^{\dagger}(\frac{\partial \mathbf{b}}{\partial \mathbf{p}} - \frac{\partial \underline{\underline{\mathbf{A}}}}{\partial \mathbf{p}}\mathbf{x})$$

where $(\underline{\underline{\mathbf{A}}}^{\dagger})^{-1} = (\underline{\underline{\mathbf{A}}}^{-1})^{\dagger}$ was used in the last step[1].

The 'adjoint problem' appears on the right hand side of (2.8). If we assume that $\boldsymbol{\lambda}$ is the solution of the adjoint problem defined in Equation (2.9), then we have:

$$\underline{\underline{\mathbf{A}}}^{\dagger}\boldsymbol{\lambda}^{\dagger} = \frac{\partial g}{\partial \mathbf{x}}^{\dagger}$$

$$\boldsymbol{\lambda}^{\dagger} = \underline{\underline{\mathbf{A}}}^{\dagger-1}\frac{\partial g}{\partial \mathbf{x}}^{\dagger} \qquad (2.9)$$

$$\frac{\partial g}{\partial \mathbf{x}}\frac{\partial \mathbf{x}}{\partial \mathbf{p}} = \boldsymbol{\lambda}(\frac{\partial \mathbf{b}}{\partial \mathbf{p}} - \frac{\partial \underline{\underline{\mathbf{A}}}}{\partial \mathbf{p}}\mathbf{x})$$

While $\frac{\partial \underline{\underline{\mathbf{A}}}}{\partial \mathbf{p}}$ is in general a large quantity, it can often be sparse. We will see in the electromagnetic equations below, there is only a diagonal, linear dependence on $\mathbf{p}$

---

[1] $(\underline{\underline{\mathbf{A}}}^{\dagger})^{-1} = (\underline{\underline{\mathbf{A}}}^{-1})^{\dagger}$ for any invertible matrix $\underline{\underline{\mathbf{A}}}$. $\underline{\underline{\mathbf{A}}}^{-1}\underline{\underline{\mathbf{A}}} = \underline{\underline{\mathbf{I}}}$ and $\underline{\underline{\mathbf{I}}} = \underline{\underline{\mathbf{I}}}^{\dagger} = (\underline{\underline{\mathbf{A}}}^{-1}\underline{\underline{\mathbf{A}}})^{\dagger} = \underline{\underline{\mathbf{A}}}^{\dagger}(\underline{\underline{\mathbf{A}}}^{-1})^{\dagger}$ which by definition means that $(\underline{\underline{\mathbf{A}}}^{-1})^{\dagger} = (\underline{\underline{\mathbf{A}}}^{\dagger})^{-1}$.

in $\underline{\underline{\mathbf{A}}}$ such that this term effectively reduces to a scalar times an identity on $\mathbf{x}$. For real-valued $g$ and $\mathbf{p}$,

$$\frac{dg}{d\mathbf{p}} = \frac{\partial g}{\partial \mathbf{p}} + 2\Re\left\{\boldsymbol{\lambda}(\frac{\partial \mathbf{b}}{\partial \mathbf{p}} - \frac{\partial \underline{\underline{\mathbf{A}}}}{\partial \mathbf{p}}\mathbf{x})\right\} \tag{2.10}$$

where $\Re$ denotes the real part.

## 2.3 Adjoint Method for Electromagnetics

While Equation (2.10) provides some generality to using the adjoint method, it is most instructive when applied to a specific problem. In this body of work, we care about its application to electromagnetic inverse design. We start, then, with the following Maxwell equations [3]:

$$\nabla \times \mathbf{H} - \epsilon\frac{\partial \mathbf{E}}{\partial t} = \mathbf{J}$$
$$\frac{\partial \mathbf{B}}{\partial t} + \nabla \times \mathbf{E} = 0 \tag{2.11}$$

Taking the partial time derivative of the first equation and the curl of the second equation:

$$\nabla \times \frac{\partial \mathbf{H}}{\partial t} - \epsilon\frac{\partial^2 \mathbf{E}}{\partial t^2} = \frac{\partial \mathbf{J}}{\partial t}$$
$$\nabla \times \frac{\partial \mathbf{B}}{\partial t} + \nabla \times \nabla \times \mathbf{E} = 0 \tag{2.12}$$

Assuming the permeability everywhere is constant ($\mu = \mu_0$), then $\mathbf{B} = \mu_0\mathbf{H}$, the equations can be combined:

$$\mu_0\epsilon\frac{\partial^2 \mathbf{E}}{\partial t^2} + \mu_0\frac{\partial \mathbf{J}}{\partial t} + \nabla \times \nabla \times \mathbf{E} = 0 \tag{2.13}$$

Via separation of variables, the electromagnetic fields ($\mathbf{E}, \mathbf{H}$) and sources ($\mathbf{J}$) are written with a harmonic time dependence given by $e^{i\omega t}$. This leads to:

$$-\omega^2\mu_0\epsilon\mathbf{E} + \nabla \times \nabla \times \mathbf{E} = -i\omega\mu_0\mathbf{J} \tag{2.14}$$

Letting $\epsilon = \epsilon_0\epsilon_r$ and dividing through by $\mu_0$:

$$-\omega^2\epsilon_0\epsilon_r\mathbf{E} + \frac{1}{\mu_0}\nabla \times \nabla \times \mathbf{E} = -i\omega\mathbf{J} \tag{2.15}$$

Instead of immediately working through the full vectorized treatment, we will first consider a simplified class of electromagnetic problems that can be modeled by a scalar version of these equations where two uncoupled polarizations are supported, each one described by a single vector component of **E** or **H**, here labeled as $E_z$ or $H_z$. This can model on-chip photonic structures and free space structures that are extruded infinitely (or for length scales on the order of tens of wavelengths in practice) in one of the in-plane dimensions. For the following examples, we will consider just the $E_z$ polarization.

Equation (2.15) contains a vector equation defined across space. By simplifying to $\mathbf{E} = E\hat{\mathbf{z}}$, we will now write the vector component as $\mathbf{E}_z$ to refer to the value $E_z$ at every point in space. This assumes no variation in the structure or fields in the $z$-direction. It can be checked that under this assumption, the two curl operations in sequence create intermediate $x$- and $y$-polarized components, which are mapped back to only a $z$-polarized field. Then, the equation can be discretized using finite differences for the spatial derivatives and mapped back to Equation (2.1) via:

$$
\begin{aligned}
\mathbf{p} &:= \boldsymbol{\epsilon}_r \\
\underline{\underline{\mathbf{A}}}(\mathbf{p}) &:= -\omega^2 \epsilon_0 \epsilon_r + \frac{1}{\mu_0} \nabla \times \nabla \times \\
\mathbf{x} &:= \mathbf{E}_z \\
\mathbf{b} &:= -i\omega \mathbf{J}_z
\end{aligned}
\tag{2.16}
$$

Since the matrix $\underline{\underline{\mathbf{A}}}$ is real, we can modify the last line of Equation (2.8) to use a transpose instead of the conjugate transpose:

$$
\frac{\partial g}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial \mathbf{p}} = [(\underline{\underline{\mathbf{A}}}^T)^{-1} \frac{\partial g}{\partial \mathbf{x}}^T]^T (\frac{\partial \mathbf{b}}{\partial \mathbf{p}} - \frac{\partial \underline{\underline{\mathbf{A}}}}{\partial \mathbf{p}} \mathbf{x})
\tag{2.17}
$$

Simplifying to the common case where the source region does not overlap with the design region (i.e. $\frac{\partial \mathbf{b}}{\partial \mathbf{p}} = \mathbf{0}$) and further using that $\underline{\underline{\mathbf{A}}}$ is symmetric, which reflects the electromagnetic property of reciprocity in the system (i.e. $\underline{\underline{\mathbf{A}}} = \underline{\underline{\mathbf{A}}}^T$):

$$
\frac{\partial g}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial \mathbf{p}} = -[\underline{\underline{\mathbf{A}}}^{-1} \frac{\partial g}{\partial \mathbf{x}}^T]^T \frac{\partial \underline{\underline{\mathbf{A}}}}{\partial \mathbf{p}} \mathbf{x}
\tag{2.18}
$$

$\frac{\partial \underline{\underline{\mathbf{A}}}}{\partial \mathbf{p}}$ is a large tensor in general, but in many cases may be sparse. In the case of the outlined electromagnetic equations, it simplifies to a scalar $-\omega^2 \epsilon_0$ factor:

$$\frac{\partial g}{\partial \mathbf{x}}\frac{\partial \mathbf{x}}{\partial \mathbf{p}} = \omega^2 \epsilon_0 [\underline{\underline{\mathbf{A}}}^{-1}\frac{\partial g}{\partial \mathbf{x}}^T]^T \mathbf{x} \tag{2.19}$$



Figure 2.1: Illustration of forward and adjoint sources used to compute the gradient for a focusing optimization. (a) The forward source is run to compute the expected output from the device. (b) The adjoint source is derived from the figure of merit and in this case corresponds to exciting a dipole at the desired focal point.

Suppose we would like to optimize that, for a given source input to a device, the intensity at some point on the other side of the device is maximized. This setup is illustrated in Figure 2.1 where the coordinates of the point are marked as $(j, k)$. For a matrix solution, the fields would typically be flattened into a vector whereas here we have noted a coordinate of the target point. These representations can be easily mapped between each other. The figure of merit is:

$$\max_{\epsilon_r} \mathrm{I}_{z,jk} = \max_{\epsilon_r} \mathrm{E}_{z,jk}\mathrm{E}^*_{z,jk}$$

$$g(\mathbf{E}_z) := \mathrm{E}_{z,jk}\mathrm{E}^*_{z,jk}$$

$$\frac{\partial g}{\partial \mathbf{E}_z} = \mathrm{E}^*_{z,jk}\boldsymbol{\delta}_{z,jk} \tag{2.20}$$

$$\frac{\partial g}{\partial \mathbf{E}^*_z} = \mathrm{E}_{z,jk}\boldsymbol{\delta}_{z,jk}$$

where $\boldsymbol{\delta}_{z,jk}$ is a delta function that creates a dipole source at the $(j, k)$ coordinate.

$g(\mathbf{E_z})$ defined in Equation (2.20) is a real-valued objective function. Thus we have:

$$
\begin{aligned}
\frac{dg}{d\epsilon_r} &= \frac{\partial g}{\partial \mathbf{E}_z}\frac{\partial \mathbf{E}_z}{\partial \epsilon_r} + \frac{\partial g}{\partial \mathbf{E}_z^*}\frac{\partial \mathbf{E}_z^*}{\partial \epsilon_r} = \frac{\partial g}{\partial \mathbf{E}_z}\frac{\partial \mathbf{E}_z}{\partial \epsilon_r} + \left(\frac{\partial g}{\partial \mathbf{E}_z}\frac{\partial \mathbf{E}_z}{\partial \epsilon_r}\right)^* \\
&= 2\Re\left\{\frac{\partial g}{\partial \mathbf{E}_z}\frac{\partial \mathbf{E}_z}{\partial \epsilon_r}\right\} = 2\omega^2\epsilon_0\Re\left\{[\underline{\underline{\mathbf{A}}}^{-1}\mathrm{E}_{z,jk}^*\boldsymbol{\delta}_{z,jk}] \odot \mathbf{E}_{z,\mathrm{fwd}}\right\}
\end{aligned}
\tag{2.21}
$$

where $\odot$ denotes a pointwise multiplication of the fields in the design region. In the fully 3D version of the method, this becomes a dot product between the polarized electric field vectors at each point.

Calculating the gradient everywhere in the design region consists of first running the forward simulation where the device is excited with the expected source in practice, $\mathbf{J}_{z,\mathrm{fwd}}$, and the electric fields are recorded in the device, $\mathbf{E}_{z,\mathrm{fwd}} = \mathbf{E}_{z,jk}$. Then, the adjoint source, $\frac{\partial g}{\partial \mathbf{E}_z}$, is simulated using the same solver to get $[\underline{\underline{\mathbf{A}}}^{-1}\frac{\partial g}{\partial \mathbf{E}_z}^T]$. In this case, the adjoint source consists of a dipole excitation at the point we are intending to maximize intensity. The dipole is scaled by the conjugate of the electric field recorded at that point during the forward simulation (i.e., $\mathbf{J}_{z,\mathrm{adj}}$ is set such that $-i\omega\mathbf{J}_{z,\mathrm{adj}} = \mathrm{E}_{z,jk,\mathrm{fwd}}^*\boldsymbol{\delta}_{z,jk}$).[2] Again, the fields are recorded in the device to capture $\mathbf{E}_{z,\mathrm{adj}}$. Then, the gradient is computed as:

$$
2\omega^2\epsilon_0\Re\{\mathbf{E}_{z,\mathrm{adj}} \odot \mathbf{E}_{z,\mathrm{fwd}}\}
\tag{2.22}
$$

## 2.4 Intuitive Electromagnetic Picture

The usage of inverse design generally and specifically in electromagnetics has a rich history preceding this thesis. There is a useful intuitive picture and associated mathematics in Chapter 5 and Appendix A of [4] that is recommended reading and on which this section is heavily based. Further, in that work the derivation is augmented for finding permittivity gradients of general figures of merit based on both the electric and magnetic fields in a fully 3D construction. A particularly useful figure of merit written in electric and magnetic fields is that of coupling input light to some specific output mode like in the case of waveguide mode converters and power splitters, free space to on-chip grating couplers, or grating order optimization to name just a few examples [5–8].

During an optimization, we seek the gradient of a figure of merit based on electric and magnetic fields with respect to the permittivity everywhere in the design region.

---

[2] $\mathrm{E}_{z,jk,\mathrm{fwd}}^*$ is abbreviated to $\mathrm{E}_{z,jk}^*$ in previous equations.

Consider changing the permittivity in one location a small amount such that we assume there is not a significant change in the electric field at that point. Then the additional permittivity effectively induces an additional dipole polarization current that scales linearly with the magnitude of the permittivity change [9]. This current source will now cause a change to the electric and magnetic fields that are being tracked by the figure of merit. Turning on a current source at this point and simulating its effect on the fields everywhere else would indicate if raising or lowering the permittivity at this point was the correct move. As stated before, doing this test point-by-point involves running far too many simulations when the number of points we are considering changing is large (i.e., a sub-wavelength grid of permittivity generally in 3D). Unfortunately, turning on current sources for every point in the design simultaneously and observing the output fields loses information about how each point individually contributes to the output. However, due to reciprocity, the effect of a current source in one location (i.e., in the design region) on the fields in an observation region (i.e., where the figure of merit is defined) can be related to turning a current source on in the observation region and measuring the fields at every point in the design region. In other words, the effect of each current source in the design region can be separately observed by looking at the fields in the device under excitation of a current source from the observation region. Now, we can see an intuitive reason for each part of Equation (2.22). The forward fields at each point, $\mathbf{E}_{z,\mathrm{fwd}}$ describe with what phase the dipole will oscillate when the permittivity is increased a small amount. They are combined with the adjoint fields, $\mathbf{E}_{z,\mathrm{adj}}$ which contain two important features. First, they take advantage of reciprocity to probe how the induced dipole moments all individually modify the fields in the observation region. Second, they contain a phase and amplitude weighting from the forward fields in the observation region. This information describes what fields emanating from the induced dipole moments will constructively interfere with the fields in the observation region to improve the figure of merit.

## 2.5 Example Optimization Problems

To illustrate the method in its simplest form, in this section we show a few photonic optimization problems using the single polarization component $\mathbf{E}_z$ and assuming an extruded design as shown in Figure 2.2. Simulations were carried out with an open source finite difference frequency domain (FDFD) package [10]. In each optimization, the devices are allowed to be made of a continuous permittivity between $\epsilon_r = 1$ and $\epsilon_r = 2.25$ corresponding to designing between vacuum and

a material like silicon dioxide (SiO$_2$). Feasible devices at the scales we are usually interested in are binary consisting of only two materials. At some wavelengths of operation, it may be possible to pattern effective index devices by varying material density at deeply subwavelength scales. While the devices here are intended to demonstrate the inverse design concept and are thus left grayscale, devices in later chapters are optimized to be binary and techniques for achieving this are discussed in those sections.



Figure 2.2: Extruded design space for sample optimizations. The device is intended to operate in free space, but assumed to be extruded infinitely in the $z$-direction. In practice this assumption can be approximately satisfied for the center of the device when extruded a finite amount.

**Focusing Device**

A simple test case of the method is a device that focuses plane wave illumination to a point. This mimics the functionality of a lens for a single angle of incidence. In the extruded design space outlined, the focused spot is actually a focused line infinite in the $z$-direction. Here, several wavelengths are used in the optimization to ensure the resulting design has a broad bandwidth. The total figure of merit is given by the sum over wavelength of the intensity at the center of the focal plane, $\mathbf{r}_{\text{focus}}$:

$$g(\mathbf{E}_z) = \sum_{\lambda_j} \mathbf{I}_z(\mathbf{r}_{\text{focus}}, \lambda_j) = \sum_{\lambda_j} \mathbf{E}_z(\mathbf{r}_{\text{focus}}, \lambda_j)\mathbf{E}_z^*(\mathbf{r}_{\text{focus}}, \lambda_j) \qquad (2.23)$$

**a**



**b**



**c**



Figure 2.3: Simulated broadband focusing device characterization. (a) Intensity in simulation region by wavelength for a plane wave input showing the focus at the dashed blue line. The intensity is normalized to the brightest intensity in all three plots. (b) Optimized permittivity profile between $\epsilon_r = 1$ and $\epsilon_r = 2.25$. (c) Intensity at the focal plane for each wavelength in arbitrary units normalized to the maximum focal intensity across all wavelengths.

The device properties are listed in the following table in units of µm and center wavelength $\lambda_0 = 550$nm:

| Property | Value |
| --- | --- |
| Device width, $w$ | 3.375µm ($6.14\lambda_0$) |
| Device thickness, $t$ | 2.25µm ($4.10\lambda_0$) |
| Focal length, $f$ | 2.25µm ($4.10\lambda_0$) |
| Minimum feature size | 45nm ($0.082\lambda_0$) |

Table 2.1: Focusing and Color Sorting Device Properties

The results of the optimization are plotted in Figure 2.3. Each wavelength is focused to the same spot in both the lateral and axial directions as seen in Figure 2.3a,c. Figure 2.3b shows the resulting permittivity distribution of the device.

**Color Sorting Device**



Figure 2.4: Simulated color sorting device characterization. (a) Intensity in simulation region by wavelength for a plane wave input showing the changing focus position laterally at the dashed blue line. The intensity is normalized to the brightest intensity in all three plots. (b) Optimized permittivity profile between $\epsilon_r = 1$ and $\epsilon_r = 2.25$. (c) Intensity at the focal plane for each wavelength in arbitrary units normalized to the maximum focal intensity across all wavelengths.

The focusing device spectral behavior can be generalized by not enforcing each wavelength gets sent to the same spot laterally. As a preview to the color routing problem that motivates much of the work in this thesis, we show here a device that focuses bands of wavelengths to different spots on the focal plane. The figure of merit from the broadband focusing device can be modified such that each focal

point is a function of wavelength (i.e., $\mathbf{r}_{\text{focus}} = \mathbf{r}_{\text{focus}}(\lambda_j)$). Due to diffraction, shorter wavelengths can focus to tighter spots. For the same power efficiency, shorter wavelengths will have brighter centers and smaller spot sizes compared to longer wavelengths. This can be observed in the broadband focusing device output, in particular in Figure 2.3c. To balance for this effect, while not a perfect normalization in the 2D approximation, intensity for each wavelength is scaled according to the 3D scaling of intensity by wavelength for a diffraction limited spot (i.e., each intensity is normalized by the inverse of the wavelength squared). This leads to the following figure of merit:

$$g(\mathbf{E}_z) = \sum_{\lambda_j} \frac{\mathbf{I}_z(\mathbf{r}_{\text{focus}}(\lambda_j), \lambda_j)}{\mathbf{I}_{\text{max,scale}}(\lambda_j)} = \sum_{\lambda_j} \frac{\mathbf{E}_z(\mathbf{r}_{\text{focus}}(\lambda_j), \lambda_j)\mathbf{E}_z^*(\mathbf{r}_{\text{focus}}(\lambda_j), \lambda_j)}{\mathbf{I}_{\text{max,scale}}(\lambda_j)}$$

$$\mathbf{I}_{\text{max,scale}}(\lambda_j) = \frac{\lambda_0^2}{\lambda_j^2} \tag{2.24}$$

The geometric properties are the same as in Table 2.1 with the same center wavelength, $\lambda_0 = 550$nm. Optimization results are plotted in Figure 2.4. Each wavelength is focused to the same spot axially but different spots laterally depending on the wavelength as seen in Figure 2.4a,c. Figure 2.4b shows the optimized permittivity distribution of the device.

**Stimulated Emission Depletion Microscopy**

In stimulated emission depletion (STED) microscopy, two wavelengths are focused into a sample with different spatial patterns to create a super-resolution imaging system [11]. One of these wavelengths is focused to a tight spot to excite electrons and create a fluorescence signal from the sample. At the same time, another wavelength is focused to a mode with an intensity void in the center. This pattern drives stimulated emission occurring at this depletion wavelength leaving only the center of the pattern to fluoresce. By filtering light at the stimulated emission wavelength on the collection side, the reduced size fluorescence region effectively images the sample at sub-diffraction limited resolution. This involves a sophisticated optical setup [12], where light of different wavelength and spatial mode are multiplexed into the optical path to be condensed onto a sample. Fluorescence is then collected, spectrally filtered, and focused onto a detector. There has been work to create different focusing profiles as a function of wavelength for co-aligned beams such as in easySTED [13] which uses a chromatic segmented waveplate to induce different po-

Figure 2.5: STED device characterization. (a) (top) Intensity in simulation region by wavelength for a plane wave input showing the changing focus shape laterally at the dashed blue line. The intensity is normalized to the brightest intensity in each plot individually. (bottom) Zoomed in view of red dashed box region with same intensity normalization. (b) Optimized permittivity profile between $\epsilon_r = 1$ and $\epsilon_r = 2.25$. (c) Intensity at the focal plane for each wavelength in arbitrary units normalized to the maximum focal intensity across all wavelengths.

larization states for the STED beam and the excitation beam. This causes the STED beam to focus with an intensity null in the center while the excitation beam focuses to a spot when passed through the objective lens. Shown here is an alternative approach where we shape different wavelengths into different focusing profiles in a single, compact optic by optimizing for wavelength-dependent scattering behavior under a single polarization excitation.

The figure of merit in the optimization is more complex in this case and tries to

not only maximize intensity in the center for the illumination and fluorescence wavelengths but also minimize the depletion intensity at this same point. Further, the depletion intensity needs to rise quickly on either side of the central focus. To achieve this, we use the following figure of merit:

$$g_{\text{sted,left}}(\mathbf{E}_z) = \frac{\mathbf{I}_z(\mathbf{r}_{\text{left}}, \lambda_{\text{sted}}) - \alpha \mathbf{I}_z(\mathbf{r}_{\text{center}}, \lambda_{\text{sted}})}{\mathbf{I}_{\text{max,scale}}(\lambda_{\text{sted}})}$$

$$g_{\text{sted,right}}(\mathbf{E}_z) = \frac{\mathbf{I}_z(\mathbf{r}_{\text{right}}, \lambda_{\text{sted}}) - \alpha \mathbf{I}_z(\mathbf{r}_{\text{center}}, \lambda_{\text{sted}})}{\mathbf{I}_{\text{max,scale}}(\lambda_{\text{sted}})}$$

$$g_{\text{illum}}(\mathbf{E}_z) = \frac{\mathbf{I}_z(\mathbf{r}_{\text{center}}, \lambda_{\text{illum}}) - \beta \mathbf{I}_z(\mathbf{r}_{\text{right}}, \lambda_{\text{illum}}) - \beta \mathbf{I}_z(\mathbf{r}_{\text{left}}, \lambda_{\text{illum}})}{\mathbf{I}_{\text{max,scale}}(\lambda_{\text{illum}})} \quad (2.25)$$

$$g_{\text{fluor}}(\mathbf{E}_z) = \frac{\mathbf{I}_z(\mathbf{r}_{\text{center}}, \lambda_{\text{fluor}}) - \beta \mathbf{I}_z(\mathbf{r}_{\text{right}}, \lambda_{\text{fluor}}) - \beta \mathbf{I}_z(\mathbf{r}_{\text{left}}, \lambda_{\text{fluor}})}{\mathbf{I}_{\text{max,scale}}(\lambda_{\text{fluor}})}$$

$$g(\mathbf{E}_z) = g_{\text{sted,left}}(\mathbf{E}_z) + g_{\text{sted,right}}(\mathbf{E}_z) + g_{\text{illum}}(\mathbf{E}_z) + g_{\text{fluor}}(\mathbf{E}_z)$$

where $\alpha = 1.5$ and $\beta = 0.5$ were used. For the STED beam, it is important to have high intensity next to the imaging location for depletion, but also important that the intensity is low in the imaging location so that there is a fluorescence signal there. Without a high enough suppression of intensity through the weighting factor $\alpha$, the optimization may settle on an undesirable solution with a wide focal spot for the STED beam that peaks in the imaging location but is also high on either side of it.

The stimulation/depletion, illumination, and fluorescence wavelengths are assumed to be $\lambda_{\text{sted}} = 560\text{nm}$, $\lambda_{\text{illum}} = 480\text{nm}$, and $\lambda_{\text{fluor}} = 500\text{nm}$, respectively. Device properties are listed in the following table in units of µm and center wavelength $\lambda_0 = 513\text{nm}$:

| Property | Value |
|---|---|
| Device width, $w$ | 11.25µm ($21.93\lambda_0$) |
| Device thickness, $t$ | 2.25µm ($4.39\lambda_0$) |
| Focal length, $f$ | 4.5µm ($8.77\lambda_0$) |
| Minimum feature size | 45nm ($0.088\lambda_0$) |

Table 2.2: STED Device Properties

Figure 2.5a depicts the difference in focal plane intensity for each of the wavelengths in the system. The illumination wavelength excites fluorescence in the center of the focal plane and the fluorescence wavelength is imaged from this same spot. In contrast, the STED wavelength focuses to a donut pattern. This wavelength

suppresses fluorescence signal on either side of the central point. The bottom of Figure 2.5a looks more closely at the focal pattern. While only the lateral donut shape was optimized for explicitly, the STED beam also forms a donut shape axially, the 3D version of which is called a bottle shape. This is often desirable in STED microscopy to also improve the resolution axially in addition to the lateral benefit. Figure 2.5b shows the optimized permittivity and Figure 2.5c plots along the blue cut line from Figure 2.5a to show a more detailed profile of the intensity for each wavelength at the focal plane.

## 2.6   Conclusion

This chapter laid the foundation for using the adjoint method for local optimization of linear systems. Further, we specialized to the case of electromagnetic inverse design enabled by efficient gradient calculation through the adjoint method. The technique is flexible and with it, we can optimize for devices that may be difficult to achieve with other design methods especially in the form factors we are often considering. This flexibility will be further demonstrated in later chapters. In free space, however, one of the largest challenges is achieving physical realizations of the optimization results. For this, we need 3D nanofabrication and the ability to constrain optimizations to respect realistic fabrication tolerances. In Chapter 3, we will address these challenges.

## References

[1]G. Strang, *Computational science and engineering*, Vol. 791 (Wellesley-Cambridge Press Wellesley, 2007).

[2]S. G. Johnson, "Notes on adjoint methods for 18.335", Introduction to Numerical Methods (2012).

[3]J. C. Maxwell, "VIII. A dynamical theory of the electromagnetic field", Philosophical transactions of the Royal Society of London, 459–512 (1865).

[4]O. D. Miller, *Photonic design: From fundamental solar cell physics to computational inverse design* (University of California, Berkeley, 2012).

[5]C. M. Lalau-Keraly, S. Bhargava, O. D. Miller, and E. Yablonovitch, "Adjoint shape optimization applied to electromagnetic design", Optics express **21**, 21693–21701 (2013).

[6]C. Dory, D. Vercruysse, K. Y. Yang, N. V. Sapra, A. E. Rugar, S. Sun, D. M. Lukin, A. Y. Piggott, J. L. Zhang, and M. Radulaski, "Inverse-designed diamond photonics", Nature communications **10**, 3309 (2019).

[7] J. Lu and J. Vučković, "Objective-first design of high-efficiency, small-footprint couplers between arbitrary nanophotonic waveguide modes", Optics express **20**, 7221–7236 (2012).

[8] D. Sell, J. Yang, S. Doshay, R. Yang, and J. A. Fan, "Large-angle, multifunctional metagratings based on freeform multimode geometries", Nano letters **17**, 3752–3757 (2017).

[9] J. D. Jackson, *Classical electrodynamics*, 1999.

[10] T. W. Hughes, I. A. D. Williamson, M. Minkov, and S. Fan, "Forward-mode differentiation of Maxwell's equations", ACS Photonics **6**, 3010–3016 (2019).

[11] S. W. Hell and J. Wichmann, "Breaking the diffraction resolution limit by stimulated emission: stimulated-emission-depletion fluorescence microscopy", Optics letters **19**, 780–782 (1994).

[12] G. Vicidomini, P. Bianchini, and A. Diaspro, "STED super-resolved microscopy", Nature methods **15**, 173–182 (2018).

[13] M. Reuss, J. Engelhardt, and S. W. Hell, "Birefringent device converts a standard scanning microscope into a STED microscope that also maps molecular orientation", Optics express **18**, 1049–1058 (2010).

*Chapter 3*

# 3D-PATTERNED INVERSE-DESIGNED MID-INFRARED METAOPTICS

This chapter is directly reproduced from [1], a paper published as a key part of the work in this thesis. As described in the abstract below, this work shows the ability to manipulate light based on all its fundamental degrees of freedom: spectral, polarization, and spatial. In it, we tackle the challenges of design, fabrication, and measurement of these complex 3D inverse designed devices.

## 3.1 Abstract

Modern imaging systems can be enhanced in efficiency, compactness, and application through the introduction of multilayer nanopatterned structures for manipulation of light based on its fundamental properties. High transmission multispectral imaging is elusive due to the commonplace use of filter arrays which discard most of the incident light. Further, given the challenges of miniaturizing optical systems, most cameras do not leverage the wealth of information in polarization and spatial degrees of freedom. Optical metamaterials can respond to these electromagnetic properties but have been explored primarily in single-layer geometries, limiting their performance and multifunctional capacity. Here we use advanced two-photon lithography to realize multilayer scattering structures that achieve highly nontrivial optical transformations intended to process light just before it reaches a focal plane array. Computationally optimized multispectral and polarimetric sorting devices are fabricated with submicron feature sizes and experimentally validated in the mid-infrared. A final structure shown in simulation redirects light based on its angular momentum. These devices demonstrate that with precise three-dimensional nanopatterning, one can directly modify the scattering properties of a sensor array to create advanced imaging systems.

## 3.2 Introduction

Nanophotonics synthesizes the study of light-matter interaction with the precise, repeatable techniques of nanofabrication. For example, dielectric metasurfaces are arrays of subwavelength scatterers that apply a spatially varying phase, polarization, or amplitude response to an incoming wavefront [2]. The local control is related to

the specific shape of each scatterer which can be chosen to compactly replicate and combine functionalities of common optical components like lenses, beamsplitters, polarizers, and waveplates or realize more novel devices such as those used for visible color routing at the pixel level [3]. For metasurfaces, the absence of substantial inter-element electromagnetic coupling is often leveraged for ease of design, but this simplification also limits the available degrees of freedom. Ultimately, we would like to tailor unique scattering behaviors for wavefronts with different spectral, spatial, and polarization properties. To do this, we can expand the design space to volumetric devices where a material is patterned at subwavelength resolution in three dimensions.

Three-dimensional (3D) devices take advantage of a larger set of optical modes to achieve unprecedented performance in a variety of tasks, but require an efficient gradient-based optimization algorithm based on full-wave electromagnetic simulation. Searching the high-dimensional space of permittivity profiles, typically for a local optimum to an electromagnetic merit function, is referred to as inverse design [4–7]. In this area, quasi-2D on-chip photonic devices have been explored extensively where patterning in the direction of light propagation is achieved in a single fabrication layer [8–10]. The fully 3D design paradigm for free space applications is yet to emerge in the infrared and visible spectra, mostly due to the increased fabrication complexity of volumetric devices. However, early works in this area utilizing one- and two-layer processes for optical applications or many-layer microwave prototypes have shown the utility of moving to thicker devices [11–13]. In this work, we optimized a two-photon polymerization (TPP) lithography process to create multilayer structures at optical wavelengths. This technique has been employed in the past for fabricating refractive, diffractive, gradient index, and extruded 2D inverse-designed optical components [14–17]. By exploiting TPP flexibility for 3D patterning at subwavelength resolution, we experimentally demonstrated multiple inverse-designed, multilayer photonic devices with applications to advanced imaging in the mid-infrared band (3-6µm).

Compact imaging systems utilize wavelength- and polarization-selective elements to characterize fundamental properties of wavefronts. Color imaging in consumer cameras follows this approach where absorptive filters are placed on top of collections of pixels to sense three or four spectral overlaps. The classic arrangement, referred to as a Bayer pattern, consists of a red, blue, and two green filters arranged 2x2 in a square [18]. Filtering schemes like this come at a cost of transmission efficiency be-

cause they absorb all light outside of their passband leading to average transmission values of approximately 33% under uniform spectral illumination. Solutions to this problem have converged on the concept of color routing where scattering structures accept light incident on a group of pixels and redirect each wavelength band to a different pixel [3, 13, 19–21]. In this manuscript, we demonstrate an efficient, multilayer inverse-designed device in the mid-infrared for accomplishing this task and further augment it to sense linear polarization. Beyond multispectral imaging, the geometry of splitting light at the the focal plane, depicted in Figure. 3.1a, can be tailored to efficiently decode other electromagnetic properties. Designing at the pixel level modularizes the optical system, allowing focal plane arrays equipped with arrangements of scattering structures to control the imaging properties of the camera. Figure 3.1b-d indicates the breadth of devices in this manuscript.

### 3.3 Results

**Multispectral and Polarization Sorting**

The first application we explored is combined multispectral and polarization imaging. Absorption spectra in the mid-infrared, part of the molecular fingerprinting region [22], correlate strongly to distinct chemical species. Among many areas of interest, this can be used for environmental monitoring [23, 24] and biomedical imaging [25, 26]. Solutions such as multiplexed filters in the mid-infrared suffer from low overall transmission efficiency [27]. They also lack a straightforward path towards multifunctionality that may be critical for a given application. In remote thermal monitoring, for example, multispectral and polarization filtering can be used in tandem to distinguish radiated and reflected light reducing instances of thermal blindness [28]. To address these challenges, we designed and fabricated a multilayer color-routing device with additional linear polarization discrimination.

The optimization goal, stated in `Eq. 1`, is constructed to sort three spectral bands from $3.7 - 5\mu m$ and distinguish between linear polarization for the middle band. The device dimensions are $30.15\mu m$ x $30.15\mu m$ x $18\mu m$ (6.6 x 6.6 x 4.0 $\lambda_{mid}^3$), broken into six $3\mu m$ thick layers, compact enough to be tiled on a high resolution focal plane array.

$$\max_{\boldsymbol{\epsilon} \in \{\epsilon_{\min}, \epsilon_{\max}\}^N} g(\mathbf{E}) = \sum_\lambda S((\sum_p \sum_q \kappa(q, p, \lambda) \frac{I_p(\mathbf{r}_q, \lambda)}{I_{\max}(\lambda)}); k)$$

$$I_p(\mathbf{r}_q, \lambda) = ||\mathbf{E}_p(\mathbf{r}_q, \lambda)||^2$$

(3.1)

Figure 3.1: Conceptual depiction of devices in this work. **(a)** 2D cross section schematic of camera with inverse designed scattering elements placed on top of photosensitive elements at the focal plane of the imaging lens. Green elements sort by color and blue elements sort by polarization, shown in more detail in (b, c). **(b)** Rendering of multispectral and linear polarization device that sorts three bands of wavelengths with the middle band further split on polarization. **(c)** Rendering of full Stokes polarimetry device that sorts four analyzer Jones vectors to different quadrants. **(d)** Rendering of angular momentum splitting device that sorts combinations of orbital angular momentum ($l$) and spin ($s$) degrees of freedom.

Electromagnetic inverse design that utilizes the mathematical adjoint method for calculating gradients with respect to material permittivity, aims to efficiently op-

Figure 3.2: Fabrication and measurement results of multispectral and linear polarization sorting device. **(a)** (top) Simulated device sorting spectrum showing both relative sorting efficiency ($S$) and net transmission ($T$) to focal plane normalized to pinhole transmission. For quadrant $k$, $S_{Q_k} = \frac{T_{Q_k}}{\sum_{i=0}^{3} T_{Q_i}}$. (bottom) Intensity images accounting for the expected imaging lens numerical aperture (NA=0.67) and showing the focal spot moving as a function of wavelength. Each numbered plot corresponds to the labeled dashed vertical line in the spectrum above it. Intensity units are arbitrary, but comparable between all plots in (a). Different maximum values on the colorbars here and in other figures are labeled and utilized for better visibility of the plotted intensity features. **(b)** Same plots as in (a) for vertical polarization input. **(c, d)** Experimental comparison plots to (a, b) respectively with standard deviation (SD) error bars. Wavelength coverage differences between simulation and experiment are due to the limited tuning range of the QCL used experimentally. **(e)** Schematic and associated SEM images of fabricated devices. The rightmost device was printed with one quarter missing to show internal structure. Scale bars: 5µm, 5µm (inset 2µm), 5µm from left to right.

timize merit functions like this, where device performance is phrased in terms of electric and magnetic fields in an observation region [4–7]. Here, the electric field intensity at the center of each quadrant is maximized for correct wavelengths and po-

larizations, and minimized for incorrect ones through choice of sign in the $\kappa(q, p, \lambda)$ weighting function where $p$ indexes the linear polarization and $q$ indexes the quadrant with center $\mathbf{r}_q$. The first summation targets broadband performance by including closely spaced wavelengths in each band to effectively optimize the device across a continuum. The purpose of the softplus function, $S$, is described in the methods alongside other optimization figures of merit for this chapter. This optimization function is nonlinear over the high-dimensional (~$10^4$-dimensional for devices in this work) permittivity tensor, composed of deeply subwavelength volumetric units (voxels). It is optimized via gradient descent enabled by the well-known adjoint method [4, 8, 9]. Combining the electric fields in the device from adjoint simulations with those from the expected illumination, in this case broadband linearly polarized plane waves, the gradient is computed in a fixed number of simulations independent of the number of design voxels. Fabrication constraints were incorporated for layering, feature size control, and binarization using averaging, lateral maximum blurring, and sigmoid projection filters, respectively [29].

The optimization results are shown in Figure 3.2a,b, where three sorting bands are present with the middle band focal spot conditioned on linear polarization. Following this result, the device was fabricated using the Nanoscribe Photonic Professional GT, where subwavelength features in the mid-infrared are readily created in a proprietary IP-Dip polymer with low loss from roughly 3.5-5.5µm [30]. The real and imaginary refractive indices of this polymer averaged over the design wavelength range are accounted for in the device optimization and presented simulation results. Using a photolithography-based liftoff procedure, a series of 30µm diameter circular aluminum apertures were fabricated on a sapphire substrate. Apertures, also included in the optimization, restrict the illumination to single devices for imaging and experimental power calibration. The Nanoscribe was aligned to write devices directly on top of the apertures. Figure 3.2e shows scanning electron microscope (SEM) images of fabricated devices. Each design was illuminated by a quantum cascade laser (QCL) with a beam waist on the order of the device size defocused such that the apertures were overfilled and sampled a roughly flat amplitude and phase section of the diverging beam. This is intended to mimic the plane wave input used for device optimization. The QCL can be tuned spectrally to probe the device at different wavelengths and the addition of linear polarizers and waveplates were used to control the input polarization. Various focal planes of the device were imaged by a zinc selenide (ZnSe) aspheric lens onto a focal plane array (see Figure 3.5). The QCL used in the experiment had a limited wavelength tuning range from

3.95μm to 5.05μm, which is why the plots in Figure 3.2c,d do not cover the full simulated spectra.

Figure 3.2c,d contains the experimental spectral and polarization sorting efficiency, overall focal transmission, and focal spot intensities to compare to simulation. Sorting efficiency measures the ratio of total focal plane signal incident on a given quadrant. In the simulated sorting efficiency, the middle band under horizontal (vertical) polarization has a width of 420nm (430nm) with a center wavelength of 4.34μm (4.35μm) with respect to its crossover points with the upper and lower bands. By comparison, in the experimental sorting efficiency, the middle band under horizontal (vertical) polarization has a width of 390nm (420nm) with a center wavelength of 4.42μm (4.43μm) with respect to its crossover points with the upper and lower bands, exhibiting a 20nm (4.71%) smaller average width and an 80nm (1.84%) average redshift. Taking into account this redshift and considering data over an equivalent total bandwidth, the peak sorting efficiencies for the three bands under horizontally (vertically) polarized illumination were 0.78, 0.63, 0.73 (0.78, 0.63, 0.74) in simulation and 0.47, 0.38, 0.46 (0.45, 0.39, 0.44) experimentally for lower, middle, and upper bands, respectively.

We suspect the reduced contrast in experiment is due to imaging aberrations, experimental beam non-idealities and fabrication errors, which include device shrinkage and feature size mismatch from proximity effects and resolution limits [31]. Transmission is measured as power through the device printed on top of a 30μm aperture that reaches the focal plane versus power through an empty 30μm aperture. We speculate the fluctuations in the experimental transmission around 4.25μm and 4.4μm could be due to minor laser power fluctuations around its transition between two QCL modules, small beam shifts between the device and pinhole normalization measurements, or differing amounts of ambient carbon dioxide ($CO_2$) absorption between measurements given the long path optical path length and strong $CO_2$ absorption near 4.25μm [32]. In the focal plane images in Figure 3.2c,d, one can see the focused spot move between the quadrants as the wavelength changes demonstrating the splitting functionality with the middle band sorted to opposite corners depending on its linear polarization.

**Full Stokes Polarimetry**

For the second application, we investigated full Stokes imaging polarimetry, where one characterizes not only the linear polarization amplitudes, but also the phase

Figure 3.3: Fabrication and measurement results of Stokes polarimetry device. **(a)** (top) Polarization contrast ($C$) in simulation quantifying the transmission ($T$) into the desired quadrant for a given analyzer state versus transmission into the same quadrant for the orthogonal state. For input $k$, $C_k = \frac{T_{|S_k\rangle \to Q_k} - T_{|\hat{S_k}\rangle \to Q_k}}{T_{|S_k\rangle \to Q_k} + T_{|\hat{S_k}\rangle \to Q_k}}$. (bottom) Transmission into the desired quadrants for the analyzer states (solid) and their orthogonal complements (dashed). **(b)** Simulated focal intensity images ($\lambda = 4.5\mu$m) accounting for imaging lens numerical aperture (NA=0.67) for the various input states where the top row contains analyzer states and the bottom row contains orthogonal states. Intensity units are arbitrary, but comparable between all plots in (b). **(c)** Comparison plots of contrast and transmission for the experimental results with analyzer states shown with open circles and orthogonal states shown with stars in the transmission plot (SD error bars). The experimental plot region is different than the simulation one in both the x- and y-axes. This region is marked with a dashed box on the simulation plot in (a). **(d)** Experimental focal intensity images ($\lambda = 4.5\mu$m) showing a bright quadrant for each analyzer state and the same quadrant dark for the complementary orthogonal state. Intensity units are arbitrary, but comparable between all plots in (d). **(e)** Schematic and associated SEM images of fabricated devices. The rightmost device was printed with one quarter missing to show internal structure. Scale bars: 5μm (inset 2μm), 5μm, 5μm from left to right.

relationship between them and the degree of polarization. This rich information is widely applicable, including in areas of biomedical imaging and diagnosis [33],

depth-based imaging and facial recognition [34, 35], atmospheric monitoring [36], and bio-inspired polarization based navigation [37]. In polarimetric imaging, the input state is cast in terms of a four-dimensional vector containing its Stokes parameters, which together specify the orientation, handedness, and degree of polarization. Complete reconstruction of this state is done through at least four independent measurements. Measurements can be multiplexed in time using a rotating waveplate [38] or in space by dividing up the area on one or more focal plane arrays [39]. The analogous geometry to using absorptive filters for color imaging is the division of focal plane (DoFP) technique where pixels are grouped together with each responsible for analyzing a specific polarization component. Many implementations use micropolarizer elements as filters [40], thus limiting the transmission efficiency of the camera to 50% by rejecting orthogonally polarized light to each filter. Lost transmission can be recovered using pixel-sized metasurface lenses that apply different phase masks to two orthogonal polarizations. For example, six projections done pairwise onto orthogonal polarization basis states directly measure the four Stokes parameters [41]. However, these six measurements contain redundant information which reduces camera resolution or degrades signal-to-noise ratio compared to a four-measurement device with the same overall size. Recently, it was shown that a metasurface grating could project incident light onto four equally spaced analyzer states with each projection belonging to a different order [42]. This approach requires propagation to spatially separate each order and is inherently chromatic due to grating dispersion. We adopted benefits and addressed shortcomings of both approaches by employing the modularity of a pixel-level design for adaptation of any camera sensor to full polarimetric imaging and utilizing a minimal four-state projection for maximal compactness. Using only four measurements is a 33% improvement in required chip area or, alternatively, a commensurate resolution or signal-to-noise ratio enhancement. As an added benefit, inverse design provides a path towards broadband polarimetry, which is difficult to achieve with metasurface and waveplate based systems due to their inherent chromatic dispersion.

We optimized a device of size 30µm x 30µm x 18µm in six 3µm layers for this purpose, with the optimization figure of merit adapted to focus four analyzer polarization states to different quadrants and reject their orthogonal states to those same quadrants. Further, we augmented the experimental system to probe arbitrary polarization states for different wavelengths depicted in Figure 3.5, 3.6. The simulation and experimental results are presented in Figure 3.3a-d and fabricated devices are shown in Figure 3.3e. Performance is quantified with two metrics. First, for each

quadrant, the contrast, $C \in [-1, 1]$, is the transmission for an analyzer state versus its orthogonal state: $C = \frac{T_{\text{analyzer}} - T_{\text{orthogonal}}}{T_{\text{analyzer}} + T_{\text{orthogonal}}}$. The optimization solution performs better for the three elliptical polarizations compared to the circular polarization state in this case, likely due to a lack of degrees of freedom. In supplementary Figure 3.7, we show a thicker 12-layer device in simulation with improved contrast of the circular polarization state. Similar to the multispectral device, there is a reduced contrast experimentally which we attribute again to fabrication and experimental imperfections. Second, transmission is quantified for each analyzer state, which, as shown in the Polarimetry Contrast Bounds supplementary section, is limited to 50% in a perfect device due to required vector overlaps between analyzer states [43]. We note that this is not a limit on total device transmission, but simply a requirement of linearity. Observing the focal plane images in Figure 3.3b,d demonstrates the polarization sorting capability of the device. The most telling indication of desired behavior is seen by observing the orthogonal state inputs where the device can theoretically fully extinguish transmission to a quadrant. By comparison to the analyzer state, the same quadrant under each orthogonal state is dark, which is supported quantitatively with specific transmission and contrast values in Figure 3.3a,c. Due to experimental and fabrication non-idealities, the measured device exhibits lower splitting contrast compared to simulation. For example, at $\lambda = 4.5\mu m$, the simulated device achieved contrasts of 0.46, 0.82, 0.84, and 0.83 for analyzer states $S_0, S_1, S_2$, and $S_3$, respectively. In comparison, for these same four states at $\lambda = 4.5\mu m$, the measured device achieved contrasts of 0.25, 0.39, 0.45, and 0.37. Practically, reconstruction of an arbitrary input polarization state can be done via a calibration procedure to account for imperfect contrast. An example calibration considering the simulated device behavior is demonstrated in the supplementary and analysis of reconstruction accuracy of pure and mixed polarization states is shown in Figure 3.13 and Figure 3.14, respectively [44].

**Angular Momentum Sorting**

A third device that we explored only in simulation sorts on the spatial degree of freedom. One property of wavefronts with spatial structure is their orbital angular momentum (OAM). Beams with discrete OAM values are modeled as Laguerre-Gaussian (LG) modes, which comprise a set of orthogonal spatial modes in the paraxial wave equation [45]. These modes are candidates for free space optical communication networks where information can be multiplexed on both the OAM and spin degrees of freedom [46]. Isolated devices that efficiently demultiplex

Figure 3.4: Simulation performance for angular momentum sorting device.

Figure 3.4: **(a)** Schematic of device and focal plane quadrants. **(b)** Contrast for sorting each state ($C \in [-1, 1]$) defined by the transmission of a state into the desired quadrant versus the transmission into the rest of the focal plane. For source $k$, $C_k = \frac{T_{|S_k\rangle \to Q_k} - \sum_{i \neq k} T_{|S_k\rangle \to Q_i}}{T_{|S_k\rangle \to Q_k} + \sum_{i \neq k} T_{|S_k\rangle \to Q_i}}$. **(c)** Transmission spectrum for ($l = -2$, $s = +1$) input with desired quadrant transmission in blue. Transmission is normalized by power through the device aperture with no device present. Inset: Intensity at focal plane (arbitrary units, but comparable to other intensity plots in figure). **(d-f)** Same plots as (c), but for ($l = -1$, $s = +1$), ($l = +1$, $s = -1$), ($l = +2$, $s = -1$), respectively.

different angular momentum values in free space [47] or from fibers [48] are essential to high bandwidth communication links. Further, an efficient sorting device can also be used to reciprocally generate different OAM states when specific quadrants are illuminated in the focal plane. Additional communication bandwidth is achievable with further spatial multiplexing of angular momentum beamlets [49], where the receiver requires an array of devices with similar geometry to those shown in this manuscript. In addition, spatially resolved OAM information can enhance contrast in imaging systems due to the asymmetric phase with azimuthal angle [50] and can find applications in high bandwidth holographic optical encryption systems [51]. Applicable to either the isolated or array geometry, we consider a routing structure sensitive to combinations of four OAM states and two spins in the form of circular polarization handedness. Figure 3.4a illustrates the optimized angular momentum sorting device, consisting of 8 design layers, each 2.4µm thick and a 30.15µm x 30.15µm lateral aperture. Sorting contrast by input $k$, $C_k \in [-1, 1]$, is shown in Figure 3.4b with an average value of 0.57 over the four angular momentum states at wavelength $\lambda = 4.47$µm. Contrast is defined as transmission into the desired quadrant versus elsewhere in the focal plane for source $S_k$ and target quadrant $Q_k$, specifically $C_k = \frac{T_{Q_k} - \sum_{i \neq k} T_{Q_i}}{T_{Q_k} + \sum_{i \neq k} T_{Q_i}}$. Each combination of OAM and spin is efficiently focused to a different quadrant as seen in Figure 3.4c-f. Transmission values for a beamsplitting and subsequent filtering scheme as opposed to routing would be limited to 25% for each state, so the proposed device roughly doubles the signal-to-noise ratio of detection. The response of the device to excitations with OAM and spin values different from the design points is analyzed in Figure 3.15, 3.16.

## 3.4 Discussion

In this work, we demonstrate multilayer, inverse-designed nanophotonic structures capable of augmenting both the performance and multifunctionality of imaging systems. Using the same configuration of lenses inside a typical camera and replacing

the scattering element on top of the focal plane array, this technology enables cameras sensitive to angular momentum, polarization state, arbitrary spectral signatures, or combinations of multiple electromagnetic properties. There are exciting avenues for exploration in the mid-infrared where fabrication is accessible via TPP tools such as the Nanoscribe. We envision targeting specific narrow absorption bands for applications in chemical and biomedical imaging and tiling different types of splitting elements in the same array. Moving forward, we can think of these elements as part of a computational imaging system where we design efficient reconstruction problems by utilizing direct control over the scattering properties of an array of elements in the optical path [52].

By utilizing a well-optimized fabrication procedure and additional design rules, TPP fabrication can be pushed to the near-infrared range [17]. Scaling to longer wavelengths in the infrared requires a polymer transparent beyond 5.5µm or the use of a material inversion technique [53] and parallel writing strategies for feasible fabrication times [54]. Currently, necessity of industry-level fabrication procedures are a barrier to demonstrating volumetric inverse design for visible wavelengths. Typical integrated circuits, like those found in modern computer processors, consist of greater than ten layers of precisely aligned subwavelength structures with respect to visible wavelengths [55]. By replacing metals with transparent optical materials in silicon-based CMOS processes, these fabrication techniques can realize the types of structures shown in this work at an industrial scale. Currently, cost, complexity, and availability of these fabrication methods limit the exploration of multilayer photonic devices in academic and prototype settings. Advances in accessible multilayer fabrication is a worthwhile endeavor to unlock a broad spectrum of imaging applications [56]. Beyond replacing traditional absorptive-based Bayer filters with color routing structures, optimized devices targeting structural color will impact reflective display technology [57] and efficient, spectrally-selective waveguide couplers will improve performance of augmented reality displays [58]. We believe there is large, untapped potential for 3D, inverse-designed photonics in both research and commercial settings. The present work is a substantial step towards the realization of these complex devices for real-world applications.

## 3.5 Methods

**Simulation**

Simulations and inverse design optimizations of structures were carried out using Lumerical (ANSYS, Inc.) finite-difference time-domain (FDTD) Maxwell equations

solver. Working in the time domain with pulsed sources, the broadband response of the device to forward and adjoint sources can be computed in single simulations. During design, we used a total field scattered field (TFSF) source to create a finite-sized, normally incident plane wave input to the device, depicted in Figure 3.8. The simulation boundary conditions were perfectly matched layers (PML) to create the effect of an infinite simulation domain for the isolated structure.

When verifying and reporting simulation performance in the manuscript, we doubled the mesh resolution in the device region, increased the simulation mesh accuracy everywhere from 2 to 4 (with 1 being the least accurate and 8 considered the most accurate), and used a defocused Gaussian source that matched more closely to the experiment compared to the plane wave optimization source. We used a Gaussian beam with a waist radius of $w_0 = 12.5\mu m$ and a defocus amount at the aperture opening of $z = 500\mu m$ such that the beam was diverging. For the angular momentum device which we did not validate experimentally, we used focused Laguerre-Gaussian modes as defined below with their waist positioned at the device face for both optimization and evaluation. We simulated devices on top of $Al_2O_3$ substrates using material permittivity values from literature [59]. Further, for the multispectral and Stokes polarimetry we included the $30\mu m$ diameter aperture in the optimization and evaluation using a perfect electrical conductor (PEC).

**Simulation Resources**

Device optimizations are run on a high performance computing cluster. During each optimization iteration, multiple forward and adjoint simulations are run to compute gradient information for the multi-objective problems specified below. For reference, we used 10 computing nodes for the multispectral and Stokes polarimetry devices and 12 computing nodes for the angular momentum device. Each node was allocated 8 Intel CPU cores (mix of Skylake 2.1 GHz and Cascadelake 2.2 GHz processors), to run simulations in parallel. With this, the optimizations complete in approximately 30-40 hours of compute time depending on the specific device thickness in the paper. The thicker Stokes polarimetry device shown in the supplemental took roughly 78 hours to complete.

**Optimization Figures of Merit and Weighting**

**Multispectral and Linear Polarization**

$$\max_{\boldsymbol{\epsilon} \in \{\epsilon_{\min}, \epsilon_{\max}\}^N} g(\mathbf{E}) = \sum_\lambda S((\sum_p \sum_q \kappa(q, p, \lambda) \frac{I_p(\mathbf{r}_q, \lambda)}{I_{\max}(\lambda)}); k)$$

$$I_p(\mathbf{r}_q, \lambda) = ||\mathbf{E}_p(\mathbf{r}_q, \lambda)||^2$$

(3.2)

where $S(x; k) = \frac{\ln(1 + e^{kx})}{k}$ is the softplus function, which ensures positivity of all figures of merit. For large $x$, which corresponds to focusing of a given wavelength primarily in the correct quadrant, this function acts in a linear regime with the onset of this regime controlled by the value $k$. In the opposite extreme, when a given wavelength is primarily in the undesired quadrants, $x$ is negative and the figure of merit tapers to zero ensuring its gradient contribution is large in the dynamic performance weighting scheme. We use $k = 2$ for this optimization.

The polarization index, $p$, separates performance for plane wave excitations with different linear polarizations. The weighting $\kappa(q, p, \lambda)$ directs bands of wavelengths evenly spaced between $3.7 - 5\mu m$ to different quadrants, $q$, with the middle band sent to one of two locations based on linear polarization and the other two bands operating independent of polarization. Wavelengths outside of a given band for a target quadrant are punished through a negative weight using $\kappa(q, p, \lambda)$.

$$\kappa(q, p, \lambda) = \begin{cases} 1, & \text{if } (\lambda, p) \text{ desired for quadrant } q \\ -\alpha, & \text{if } \lambda \text{ within } \Delta\lambda = \beta\frac{\Lambda}{3} \text{ of quadrant } q\text{'s target band} \\ 0, & \text{otherwise} \end{cases}$$

(3.3)

Here, $\Lambda$ is the full bandwidth of the optimization, in our case $1.3\mu m$. For our optimization, we used $\beta = 0.5$, which controls how far spectrally into the neighboring bands we explicitly reject intensity in a given quadrant. Finally, $\alpha$ sets the punishment weighting term for out-of-band light versus desired in-band light. We used $\alpha = 0.75$.

The normalization term $I_{\max}(\lambda)$ accounts for the fact that for a fixed power, the intensity at the center of a focal spot will scale with wavelength. We use the scaling $I_{\max}(\lambda; l, f) = l^2/(f^2\lambda^2)$ for device lateral size, $l$, and focal length, $f$.

**Full Stokes Polarimetry**    In the Stokes polarimetry device, we searched for devices with high contrast in intensity in each quadrant when illuminated with an analyzer

state versus the state orthogonal to it. For each analyzer state ($a$) and its orthogonal state ($\bar{a}$), we measured the intensity in the middle of the quadrant:

$$
\begin{aligned}
g_{a,q} &= \frac{I_{a,q}(\mathbf{r}_q, \lambda)}{I_{\max}(\lambda)} \\
g_{\bar{a},q} &= \frac{I_{\bar{a},q}(\mathbf{r}_q, \lambda)}{I_{\max}(\lambda)} \\
I_{a,q} &= ||\mathbf{E}_{a,q}(\mathbf{r}_q, \lambda)||^2 \\
I_{\bar{a},q} &= ||\mathbf{E}_{\bar{a},q}(\mathbf{r}_q, \lambda)||^2
\end{aligned}
\tag{3.4}
$$

We would like the analyzer value to be large and the other orthogonal value to be small, so we combined them together with the following product figure of merit:

$$
g_q(\lambda) = g_{a,q}(\lambda) * (1 - g_{\bar{a},q}(\lambda))
\tag{3.5}
$$

Further, given the bound of 50% transmission for an analyzer state to a quadrant for a perfect device, we only optimized the parallel intensity up until that point. The net figure of merit, then, is:

$$
\max_{\boldsymbol{\epsilon} \in \{\epsilon_{\min}, \epsilon_{\max}\}^N} g(\mathbf{E}) = \sum_{\lambda} \sum_{q} g_q(\lambda)
\tag{3.6}
$$

**Angular Momentum**

$$
\max_{\boldsymbol{\epsilon} \in \{\epsilon_{\min}, \epsilon_{\max}\}^N} g(\mathbf{E}) = \sum_{q} \sum_{s} \sum_{\lambda} \kappa(q, s) \frac{I_s(\mathbf{r}_q, \lambda)}{I_{\max}(\lambda)}
\tag{3.7}
$$

$$
I_s(\mathbf{r}_q, \lambda) = ||\mathbf{E}_s(\mathbf{r}_q, \lambda)||^2
$$

The weighting function $\kappa(q, s)$ controls whether focusing into a given quadrant, $q$, is desirable for a given mode source $s$. For a total number of optimization iterations, $M$, we define $\kappa(q, s; m)$ parameterized by iteration number, $m$, as

$$
\kappa(q, s; m) = \begin{cases} 1, & \text{if } q = s \\ -w_{\text{end}}, & \text{if } m \geq m_{\text{end}}, q \neq s \\ -(w_{\text{start}} + (w_{\text{end}} - w_{\text{start}})\frac{m}{m_{\text{end}}-1}), & \text{if } m < m_{\text{end}}, q \neq s \end{cases}
\tag{3.8}
$$

We chose $w_{\text{start}} = \frac{1}{3}$, $w_{\text{end}} = 1$, $m_{\text{end}} = 90$, and $M = 300$. The first case corresponds to one quadrant being excited for a given source, so it receives a positive weight. The second two cases describe a linear ramp of negative weight for rejecting intensity into incorrect quadrants for a certain number of iterations after

which point a constant negative weight is used for the rest of the optimization. With the dynamic performance weighting function described below, it is important that individual figures of merit stay positive for the entire optimization. The ramping of the negative weighting helps in this regard and our specific optimization maintains positive individual figures of merit throughout. We emphasize this is not a guarantee of the weighting scheme above, but a fortunate instance where it worked for our optimization. In other optimizations, we explicitly ensured the individual figures of merit remained positive.

**Dynamic Weighting** All optimizations are multi-objective and require balancing many individual figures of merit. One way of achieving a balance and preventing certain figures of merit from dominating the optimization solution is by using a dynamic performance-based weighting scheme. As certain figures of merit start performing better than others, their relative optimization weight decreases. For $N_{FOM}$ individual figures of merit, the $j^{th}$ figure of merit with current performance $f_j$ (with net performance defined as $\sum_j f_j$) is weighted:

$$w_j = \frac{2}{N_{\text{FOM}}} - \frac{f_j^p}{\sum_n f_n^p} \tag{3.9}$$

For the multispectral and linear polarization optimization, $f_j$ corresponded to each wavelength figure of merit. In other words, $f_j = f_\lambda = S((\sum_p \sum_q \kappa(q, p, \lambda)\frac{I_p(\mathbf{r}_q, \lambda)}{I_{\max}(\lambda)}); k)$. In the Stokes polarimetry case, $f_j$ corresponded to each product figure of merit for combinations of quadrant and wavelength. Specifically, each $g_q(\lambda)$ was an individual figure of merit in the weighting scheme. Finally, for the angular momentum optimization, each quadrant figure of merit corresponded to an $f_j$. In other words, $f_j = f_q = \sum_s \sum_\lambda \kappa(q, s)\frac{I_s(\mathbf{r}_q, \lambda)}{I_{\max}(\lambda)}$.

We chose $p = 2$ for all optimizations. This weighting scheme relies on positive individual figures of merit to function properly. With the above formula, weights can become negative when a given figure of merit is far ahead of the others. In these cases, we simply shifted and rescaled to ensure all weights were greater than or equal to 0. After computing individual gradients for each figure of merit, $\frac{\partial f_j}{\partial \vec{\epsilon}}$, the net gradient was computed using this weighting scheme as $\sum_j w_j \frac{\partial f_j}{\partial \vec{\epsilon}}$.

The gradients derived during the optimization were based on a focusing figure of merit (norm electric field squared at a point) and use dipole excitations as adjoint sources. Typically, this is a good proxy for transmission into a given quadrant, which

is ultimately what we care to achieve with devices placed on top of focal plane arrays. For some optimizations, the transmission-measured performance is used instead of the intensity-measured performance as the input to the dynamic weighting function. This reduced dependencies on exact normalizations of intensity with power and took into account electric field profiles that looked different than a traditional focusing profile or were shifted from the center of the quadrant.

**Optimization Fabrication Constraints**

Projection filters were used to push the optimization solution towards devices that respected certain fabrication constraints [29]. These filters are differentiable functions applied in sequence to a design variable in order to create a device variable. The device variable describes the structure being optimized and that will eventually be fabricated. We found the gradient of the figure of merit with respect to the device variable and then backpropagated that gradient to the design variable using the chain rule. The design variable was then stepped in the gradient direction.

For binarization, we used a sigmoid filter of the form $f(\rho_k) = \frac{\tanh(\beta\eta)+\tanh(\beta(\rho_k-\eta))}{\tanh(\beta\eta)+\tanh(\beta(1-\eta))}$. The strength, $\beta$, was increased over a series of 10 epochs, each with 30 iterations, starting at $\beta = 0.0625$ and doubling each epoch. The center point was fixed at $\eta = 0.5$. For layering, permittivity values were averaged vertically over the layer thickness at each lateral point. This corresponds to averaging the two-dimensional gradient of each slice in a given layer during backpropagation.

For minimum feature size, we used an averaging blurring function that tapered from the center pixel, specifically:

$$f(\rho_k) = \frac{1}{N} \sum_{r_j \in \Omega_k} \frac{b_r + 1 - r_j}{b_r} \rho_j \tag{3.10}$$

where $N = \sum_{r_j \in \Omega_k} \frac{b_r+1-r_j}{b_r}$ for normalization [60]. $b_r$ is the blur radius based on a circle inscribed on a square with sides equal to the desired minimum feature size of 750nm, such that $b_r = 0.5 * \sqrt{2} * 750\text{nm} = 530\text{nm}$ ($r_j$ is the distance of voxel $j$ from voxel $k$). We computed this sum out to $r_{j,\max} = b_r$. This filter tends to increase the density value of voxels nearby a solid one. This encourages a minimum feature size in the solid domain. However, the drawback of this method is that it does not guarantee a minimum feature size and we do end up with features smaller than the minimum value in our final designs. This can be improved through filters that blur even further out, directly fixing the design grid to be that of minimum feature size increments, or use of a level set procedure at the end of the optimization with a

feature size constraint. The other disadvantage of this method is it only attempts to control feature size in the solid domain and does not address minimum gap sizes in the void domain.

Some fabrication constraints are difficult to encapsulate in a filtering function. Final designs needed to consist of a single piece of material for fabrication. Further, with the Nanoscribe, an enclosed void cannot be realized because the liquid polymer would have no way to escape during development. During the optimization, every 8 iterations, each design layer was patched to ensure it was a single piece of material and bridges formed in this step were restricted from changing for the next 8 iterations until the patching occurred again. The bridges were chosen via a shortest path (Dijkstra's) graph algorithm with the cost equal to the amount a given voxel would need to move to become fully solid. Islands of material were connected via a greedy minimum spanning tree approach with net bridge costs used as the edge weight for connecting two islands together. The density value of the inserted bridge was then set to 0.75 (fully solid corresponds to a density of 1) everywhere along its path. This works to ensure solid connectivity and often the void connectivity is maintained by chance throughout the optimization. Final patches were done after the optimization to create void connectivity if it had not happened naturally. These were usually minimal changes that did not have large effects on the device performance. Nevertheless, the reported simulation results in this manuscript used the fully patched designs that were fabricated.

**General Optimization Details**

For each design, the pre-filtered density is initialized to a midpoint uniform value of 0.5 ($\rho \in [0, 1]$). For some designs, additional built-in structural robustness is added by enforcing a solid material border around the edges of the device. This has little effect on the final optimization result. Only one initial seed is used for each optimization and results tend to not be plagued by poor-performing local minima. Empirically, when optimizing with a relatively low index contrast as is the case in this work, designs tend to converge to well-performing solutions without needing to run large numbers of individual optimizations. However, augmenting the methods presented here to perform more of a global search on the design space could potentially result in even higher performance. The number of layers was fixed for each design and chosen empirically to achieve high photonic performance while still maintaining fabricability. Further, the layer height was not part of the optimization process. All layers were optimized simultaneously where each layer gradient was

computed as the average gradient over all of the vertical slices in the layer. While not shown here via direct comparison to single layer designs, the necessity of multilayer structures for these types of sorting tasks has been evaluated in the past including our prior work on this topic [13]. Further, Figure 3.7 demonstrates the utility of adding thickness via additional layers to a device optimization region. It shows improvement of the sorting contrast for circular polarization in the polarimetry device.

**Device Fabrication**

The fabrication procedure for the measured devices is shown in Figure 3.9. Devices were printed directly on top of 30µm apertures defined via a photolithography-based lifotff procedure. Apertures were 150nm thick and controlled the illumination on the device. They were also used for measuring beam power through blank apertures to normalize net transmission of the device. Negative-tone photoresist, AZ nLoF 2070 (MicroChemicals GmbH), was patterned with photolithography to create a variety of apertures as well as alignment marks for the optical setup. Following oxygen and argon direct plasma cleaning to remove undesired residual photoresist left after development, 150nm of aluminum (Al) was deposited via electron beam evaporation. The apertures were lifted off in acetone and the substrate was cleaned in IPA followed by DI water. IP-Dip resist was dropped onto the substrate surface for direct write lithography using the Nanoscribe Photonic Professional GT. In the Nanoscribe, the apertures were located by moving the stage after the substrate surface was found by the microscope. By turning on the laser below the polymerization threshold such that the microscope could still image a fluorescence signal from the laser focus, we aligned the center of the printing axes to the aperture center. After writing, the devices were developed in propylene glycol methyl ether acetate (PGMEA) for 20 minutes and rinsed in two successive IPA baths for 3 minutes each. The surface of the substrate was dried with a gentle nitrogen stream. We found that critical point drying was not necessary for the integrity of our structures through the drying process. For imaging in the scanning electron microscope, a 5nm coating of platinum (Pt) was sputtered onto the surface to reduce charging effects due to the insulating polymer.

**Optical Experimental Setup**

The optical setup shown in Figure 3.5 illuminated the devices through the sapphire substrate with a diverging Gaussian beam across the pinhole aperture on which the

the device was printed. Polarization in the case of the multispectral device was changed via a half-wave plate and linear polarizer where the former served to rotate more power into the polarization less overlapped with the laser output mode. In the case of the Stokes measurement setup, a combination of a linear polarizer, half-wave plate, and quarter-wave plate were used to achieve desired input states to the device. These wave plates applied a wavelength-dependent retardance, which was taken into account in their chosen rotation. To ensure the polarization states were correct, we used a traditional method for reconstruction of the polarization state consisting of a quarter-wave plate followed by a Wollaston prism (polarizing beam splitter). The prism split orthogonal linear polarizations into two different angles which were imaged onto a power meter. By rotating the quarter-wave plate to three known rotations and measuring the power in each angle, we reconstructed what the input polarization state must have been. The reconstructed state overlaps are shown in Figure 3.6c for each probing wavelength of the Stokes measurement setup.

**Imaging, Focal Length, and Chromatic Dispersion**    The imaging objective, L2 in Figure 3.5, was translated in the axial direction to measure different planes moving back from the substrate surface. First, we observed and verified in simulation the presence of significant focal shift with respect to the focal length of the device with wavelength due to chromatic aberration of the aspheric lens. Modeling the optical setup in simulation using the Stellar Software Beam4 open source ray tracing program, we found this shift to be 31.4µm between $\lambda = 3.95$µm and $\lambda = 5$µm. Experimentally, we measured this effect by imaging the diffraction pattern of an empty 30µm aperture on the substrate surface. By tuning the focus until the diffraction pattern disappeared, we could ensure we were focused on the aperture surface for a given wavelength. We adjusted the axial position of the lens for different laser wavelengths until we were at the surface of the aperture and noted the micrometer position in order to characterize the chromatic focal shift of the imaging lens. Experimentally, we computed this dispersion to be 46µm for the same range of $\lambda = 3.95$µm and $\lambda = 5$µm. If we assume this can be used as a calibration of the micrometer on the stage, then each marking corresponds to 0.68 µm tick$^{-1}$.

Using stage markings, we found the best focal plane of the multispectral device for $\lambda = 3.95$µm to be located at 63 ticks from the substrate surface. For a total device height of 19.5µm and designed focal length of 25µm, we expected the focus to be located at 44.5µm off the substrate surface. Applying the calibration above of 0.68 µm tick$^{-1}$, we estimate the measured focal plane to be located at 43µm from

the surface ($f = 23.5$µm for assumed device height of 19.5µm), which is close to the design and within reasonable inaccuracies of the above calibration and small errors in printed device height.

For the multispectral device, we took 15 measurements evenly spaced between 3.95µm and 5µm. Since the chromatic dispersion is not equal across this whole range, we broke the range into two parts and linearly interpolated the axial position of the imaging objective to probe the same focal plane for each wavelength. Between 3.95µm and 4.48µm, we interpolated over 26 ticks corresponding to 17.7µm and between 4.48µm and 5µm, we interpolated over 20 ticks corresponding to 13.7µm.

For the Stokes polarimetry device, we measured at three distinct wavelengths, 4.18µm, 4.5µm, and 4.8µm. We directly set the focal lengths based on the empirical axial position of a blank 30µm aperture on the substrate for each wavelength. We used the same focal length of 63 ticks corresponding to around 43µm. The device height in this case was designed to be 19.8µm, so this corresponded to $f = 23.2$µm.

**Transmission Normalization**    The imaging of device focal planes onto the camera needed to be calibrated with a net device transmission. We quantified the transmission of the device by using an empty circular aperture on the substrate with the same diameter as the one sectioning off illumination to the device. Using a power meter, we measured the power through the aperture without the device versus the power through the aperture with the device. We tracked the beam center on the camera and the method depended on properly centering the beam on both apertures. Further, we assumed that laser power was not fluctuating significantly in time and the beam center was not shifting upon successive wavelength tuning due to thermal effects and a changing laser mode. Using this method, we saw consistent and expected transmission values through the device with only minor fluctuations for the multispectral device around 4.4µm, which may have been due to invalidation of assumptions previously stated. This wavelength is close to the crossover between two modules in the QCL covering different spectral ranges and we speculate the power may be less stable here compared to other wavelengths. The measured transmission was assumed to be contained in the camera image of the focal plane and its surrounding area. We then assumed that the transmission corresponding to a patch of the camera image was equal to the ratio of its intensity to the total intensity multiplied by the net measured transmission. Transmission into the focal plane, for example, was computed by multiplying the measured total transmission value by the ratio of intensity

in the focal plane to the intensity in the focal plane and surrounding area. Camera images were taken of the focal plane for each wavelength and a background of the camera (with the laser emission off) was taken immediately afterward. Taking the background immediately after each measurement reduced error in the background drifting over the course of the long experimental procedure from temperature drifts in the room or the camera housing itself. These background images were subtracted from the camera images.

**Stokes State Creation and Verification**

Each polarization state was generated through choice of rotation of a linear polarizer, a half-wave plate, and a quarter-wave plate pictured in Figure 3.6a. The wave plates are chromatic components with a retardance defined for a given wavelength. Using the provided retardance data from Thorlabs for each component, we computed the effect a rotated component will have on each input wavelength. By optimizing the choice of angles of these three components, we can generate all of the desired input polarization states to the device. To verify the correctness of each state, we used the setup in Figure 3.6b consisting of a quarter-wave plate and Wollaston prism. The Wollaston prism splits the input polarization into its x- and y-polarization components, each of which are imaged onto and measured by a power meter. Under different rotations of the quarter-wave plate, the magnitude of x- and y-polarized components will change as a function of the input state. By measuring these components under three rotations and using the specified retardance values of this quarter-wave plate as a function of wavelength, we reconstructed the input state. We used rotation values of 0°, 22°, and 44° rotations of the fast axis with respect to the x-polarization direction. Plotted in Figure 3.6c are the magnitudes of the vector overlaps of the reconstructed and the desired Jones state for each wavelength. Note the y-axis on the plot begins at 0.9.

## 3.6   Supplementary Information

**Two-Photon Polymerization (TPP) Accuracy**

Fabrication via TPP is a flexible and powerful method, but also has known challenges in printing accuracy [31]. We observe shrinkage of the structure, which is dependent on the height of the layer from the substrate. Material printed on the bottom layer is not able to shrink from its printed size because it is physically adhered to the substrate. The topmost layer is roughly 90% of the desired lateral size and the bottom layer is close to the expected size. We also observe dilation of the smallest

features in the design. Designs were compensated for this effect by pre-eroding features in the STL file before printing. Finally, the Nanoscribe had a mismatch between the feature size in each lateral direction. This is not a limitation of TPP, but likely the result of astigmatism in the optical alignment of our specific tool.

**Laguerre Gaussian Modes for Angular Momentum Splitter**

A spatially varying field can carry orbital angular momentum (OAM). Discrete values of OAM, $l$, can be found in the Laguerre-Gaussian orthonormal basis for solutions of the paraxial wave equation [45]. We used a simplified set with $p = 0$, such that each mode was defined at its waist ($z = 0$) with spatial profile in cylindrical coordinates:

$$u(r, \phi, z = 0) = (\frac{r\sqrt{2}}{w_0})^{|l|} e^{\frac{-r^2}{w_0^2}} e^{-il\phi} \tag{3.11}$$

where $w_0$ is the waist radius of the beam. We chose $w_0 = 8.5\mu m$ to ensure the mode was confined to the device. Transmission plots shown are geometrically normalized against the transmission of this beam through the device aperture with no device present. We can further assign a spin angular momentum of the mode by choosing the handedness of its circular polarization. The following pairs of OAM values $l$ and spin values $s$ were used in the optimization: $(l, s) = (-2, 1), (-1, 1), (1, -1), (2, -1)$. These states were assigned to quadrants starting with the top right (blue) and moving counterclockwise (green, red, magenta).

**12-Layer Stokes Polarimetry Device**

The polarimetry device in the main text consists of six $3\mu m$ layers and struggles to achieve equal contrast for all four analyzer states with the circular polarization state lagging the others. We speculate this may be due to lack of degrees of freedom in the thickness of device. As a comparison, we optimize a thicker device consisting of twelve $3\mu m$ layers to see if the solution will display better contrast for all analyzer states. In Figure 3.7 we show the comparison of the thicker device to the original. While the quadrant transmission per analyzer state is slightly reduced, the contrast metric is improved for the circular polarization state without sacrificing the other analyzer state contrasts.

**Polarimetry Splitting Bounds**

We can model the Stokes polarimetry device as a linear system that projects an input Jones state describing the x- and y-polarized electric field components onto several

analyzer states. The Jones polarization is a two-dimensional complex vector. The four analyzer states for our device are specifically chosen Jones vectors. In Figure 3.10, analyzer states correspond to $|v_i\rangle$, where $N = 4$ for the device in the paper. We assume the device outputs into four spatially distinct modes $|w_k\rangle$, such that we take them to be orthogonal ($\langle w_i|w_k\rangle = \delta_{ik}$). Specifically, we model each output mode as a focused spot in a different quadrant of the focal plane and thus we assume the lack of spatial overlap implies orthogonality to a good approximation. The functionality of the device is described by an operator $\hat{Q}$ where projection of an input state on each analyzer direction modulates the amplitude of an outgoing mode. We write

$$\hat{Q} = \sum_i \alpha_i \,|w_i\rangle \,\langle v_i| \tag{3.12}$$

Without loss of generality, we assume $\alpha_i$ is real. Any complex phase can be included in output mode $|w_i\rangle$.

**Maximum transmission into each analyzer state**

Next, we assume for simplicity that all states have the same projection efficiency, such that $\alpha_i = \alpha$. The transmission bound will differ from the following if each state does not split at the same projection efficiency. Consider an arbitrary state $|a\rangle$ and its orthogonal complement $|\bar{a}\rangle$. The action of $\hat{Q}$ on $|a\rangle$ is

$$\hat{Q} \,|a\rangle = \alpha \sum_i |w_i\rangle \,\langle v_i|a\rangle \tag{3.13}$$

Taking the vector magnitude squared of the resulting state

$$\langle a|\hat{Q}^\dagger\hat{Q}|a\rangle = \alpha^2 \sum_{i,j} \langle w_i|w_j\rangle \,\langle v_j|a\rangle \,\langle a|v_i\rangle \tag{3.14}$$

Since $\langle w_i|w_k\rangle = \delta_{ik}$, the double sum reduces to

$$\langle a|\hat{Q}^\dagger\hat{Q}|a\rangle = \alpha^2 \sum_i \langle v_i|a\rangle \,\langle a|v_i\rangle = \alpha^2 \sum_i |\,\langle a|v_i\rangle\,|^2 \tag{3.15}$$

Following this pattern, we also have

$$\langle \bar{a}|\hat{Q}^\dagger\hat{Q}|\bar{a}\rangle = \alpha^2 \sum_i \langle v_i|\bar{a}\rangle \,\langle \bar{a}|v_i\rangle = \alpha^2 \sum_i |\,\langle \bar{a}|v_i\rangle\,|^2 \tag{3.16}$$

Due to energy conservation, we cannot have gained any magnitude through applying $\hat{Q}$ on the state so $\langle a|\hat{Q}^\dagger\hat{Q}|a\rangle \leq 1$ and $\langle \bar{a}|\hat{Q}^\dagger\hat{Q}|\bar{a}\rangle \leq 1$. Summing these together, we get

$$\langle a|\hat{Q}^\dagger\hat{Q}|a\rangle + \langle \bar{a}|\hat{Q}^\dagger\hat{Q}|\bar{a}\rangle = \alpha^2 \sum_i (|\,\langle a|v_i\rangle\,|^2 + |\,\langle \bar{a}|v_i\rangle\,|^2) \leq 2 \tag{3.17}$$

Because the Jones vector space is two-dimensional, $|a\rangle$ and $|\bar{a}\rangle$ form an orthonormal basis, so by definition $(|\langle a|v_i\rangle|^2 + |\langle \bar{a}|v_i\rangle|^2) = 1$. Thus, the sum simply becomes

$$\langle a|\hat{Q}^\dagger \hat{Q}|a\rangle + \langle \bar{a}|\hat{Q}^\dagger \hat{Q}|\bar{a}\rangle = N\alpha^2 \leq 2 \tag{3.18}$$

If we assume $\alpha$ is the largest it can be, then $\alpha^2 = \frac{2}{N}$. For $N = 4$ as is the case for the device in this manuscript, $\alpha^2 = 0.5$. Thus, the maximum transmission we can achieve for each analyzer state into its output mode is 0.5.

**Minimum overlap between analyzer states**

Given a maximum transmission efficiency of 0.5 for each analyzer state, we can set a minimum overlap, $\beta$, for Jones vector analyzer states used in the splitter. While the choice is not unique, a maximally spaced set of vectors will have a common mutual overlap. Assume for our set of analyzer states,

$$|\langle v_i|v_j\rangle|^2 = \begin{cases} 1 \text{ if } i = j \\ \beta^2 \text{ if } i \neq j \end{cases} \tag{3.19}$$

Sending in an analyzer state to the device

$$\hat{Q}|v_k\rangle = \alpha \sum_i |w_i\rangle \langle v_i|v_k\rangle \tag{3.20}$$

Taking the magnitude like before and using the orthogonality of the $|w_i\rangle$ states

$$\langle v_k|\hat{Q}^\dagger \hat{Q}|v_k\rangle = \alpha^2 \sum_i \langle v_i|v_k\rangle \langle v_k|v_i\rangle = \alpha^2 \sum_i |\langle v_k|v_i\rangle|^2 \tag{3.21}$$

Using the common overlap between states in the analyzer set and requiring that by energy conservation this magnitude squared is bound by 1,

$$\langle v_k|\hat{Q}^\dagger \hat{Q}|v_k\rangle = \alpha^2(1 + (N-1)\beta^2) \leq 1 \tag{3.22}$$

The relation between $\alpha$ and $\beta$, then is given by

$$\alpha^2 \leq \frac{1}{1 + (N-1)\beta^2} \tag{3.23}$$

Suppose we specialize to the case where the transmission is maximized into each analyzer state ($\alpha^2 = \frac{2}{N}$) and we have no lost transmission for any given analyzer

state through the system ($\langle v_k | \hat{Q}^\dagger \hat{Q} | v_k \rangle = 1$). Then,

$$\alpha^2(1 + (N-1)\beta^2) = 1$$
$$\frac{2}{N}(1 + (N-1)\beta^2) = 1$$
$$1 + (N-1)\beta^2 = \frac{N}{2} \quad \text{(3.24)}$$
$$(N-1)\beta^2 = \frac{N-2}{2}$$
$$\beta^2 = \frac{N-2}{2(N-1)}$$

Note the case of $N = 2$ requires no overlap between the vectors with $\beta^2 = 0$ and $\alpha^2 = \frac{2}{N} = 1$ because that matches the dimensionality of the Jones vector space. However, from two measurements, we cannot reconstruct the full Stokes vector where in order to do so we need at least $N = 4$. As stated before, for $N = 4$, $\alpha^2 = 0.5$ at best and with no lost transmission for the analyzer states, $\beta^2 = \frac{1}{3}$.

**Polarimetry Contrast Bounds**

The contrast figure of merit for the Stokes polarimetry device is independent of overall transmission. For a given quadrant corresponding to analyzer state $|v_i\rangle$ and orthogonal complement $|\bar{v}_i\rangle$, the contrast is related to the analyzer transmission $T_{\text{analyzer}}$ and orthogonal transmission $T_{\text{orthogonal}}$ to the quadrant as $C = \frac{T_{\text{analyzer}} - T_{\text{orthogonal}}}{T_{\text{analyzer}} + T_{\text{orthogonal}}}$. In order to get a contrast of $C = 1$, we need to be able to completely extinguish light in the analyzer quadrant for the orthogonal state.

**Analyzer state transmission to all quadrants**

We first show that a given analyzer state must necessarily appear in more than just the desired quadrant. Following from the notation above, the action of the device on an analyzer state, $|v_k\rangle$ is given by

$$\hat{Q} |v_k\rangle = \sum_i \alpha_i |w_i\rangle \langle v_i | v_k \rangle \quad \text{(3.25)}$$

We ask how much overlap does this have with one of the output modes $|w_j\rangle$ not corresponding to the analyzer quadrant (i.e. $i \neq j$).

$$\langle w_j | \hat{Q} | v_k \rangle = \sum_i \alpha_i \langle w_j | w_i \rangle \langle v_i | v_k \rangle = \alpha_j \langle v_j | v_k \rangle \quad \text{(3.26)}$$

where we used $\langle w_j | w_i \rangle = \delta_{ij}$ to eliminate the sum. However, as we showed above, with four analyzer states, $\langle v_j | v_k \rangle \neq 0$ even for $j \neq i$. So there is energy in the other quadrants according to the splitting efficiency of the $j^{th}$ analyzer state and the overlap between the $j$ and $k$ analyzer states.

**Extinguishing orthogonal state to analyzer quadrant**

We now check if an orthogonal state can be completely extinguished to the analyzer quadrant, which will determine if we can achieve a contrast of $C = 1$. When we send in the orthogonal state to a given analyzer, $|\bar{v}_k\rangle$, the device output is given by

$$\hat{Q}\,|\bar{v}_k\rangle = \sum_i \alpha_i\,|w_i\rangle\,\langle v_i|\bar{v}_k\rangle \tag{3.27}$$

Since it is true that $\langle v_k|\bar{v}_k\rangle = 0$ by definition, the sum is reduced to

$$\hat{Q}\,|\bar{v}_k\rangle = \sum_{i\neq k} \alpha_i\,|w_i\rangle\,\langle v_i|\bar{v}_k\rangle \tag{3.28}$$

Now, we ask how much overlap does this have with the output mode corresponding to this analyzer quadrant, $|w_k\rangle$, since we are interested in seeing if this overlap can be zero.

$$\langle w_k|\hat{Q}|\bar{v}_k\rangle = \sum_{i\neq k} \alpha_i\,\langle w_k|w_i\rangle\,\langle v_i|\bar{v}_k\rangle = 0 \tag{3.29}$$

where $\langle w_k|w_i\rangle = \delta_{ki}$ is only nonzero for $i = k$, but the sum explicitly ranges over values of $i \neq k$. Thus, we can extinguish a quadrant completely for a given orthogonal state and a contrast of 1 is theoretically achievable even if we transmit all incident light through the device to the focal plane.

**Polarimetry Analyzer States**

The choice of analyzer states that fits the above criteria is not unique, but will correspond to a tetrahedron with points lying on the Poincaré sphere. First, we choose evenly spaced pure polarization states in Stokes space and then evaluate their mutual overlaps in Jones space. One state is fixed in Stokes space to be right circular polarization (RCP), which is encoded as $\begin{bmatrix} 1, & 0, & 0, & 1 \end{bmatrix}$. This choice is arbitrary and different starting states will generate equally suitable sets of analyzer states. Staying on the Poincaré sphere surface means the first entry is fixed to 1 (from here, we write the vector in terms of $S_1$, $S_2$, and $S_3$). The other three states should lie on a circle with a fixed polar angle from this first state such that all mutual overlaps are the same. For polar angle $\theta$ and azimuthal angle $\phi$, these states can be parameterized $\begin{bmatrix} \sin\theta\cos\phi, & \sin\theta\sin\phi, & \cos\theta \end{bmatrix}$. To evenly spread out these states azimuthally, the spacing should be $\Delta\phi = \frac{2\pi}{3}$. We make the non-unique choice to set the first $\phi = 0$. The first two states on the circle, then are $\begin{bmatrix} \sin\theta, & 0, & \cos\theta \end{bmatrix}$ and $\begin{bmatrix} \sin\theta\cos\frac{2\pi}{3}, & \sin\theta\sin\frac{2\pi}{3}, & \cos\theta \end{bmatrix}$. Evaluating the dot product between any of the

states on the circle and the right circular polarization state yields $\cos\theta$. The first two states on the circle have a dot product of $\sin^2\theta\cos\frac{2\pi}{3} + \cos^2\theta$. Equating these two values generates the relation:

$$\sin^2\theta\cos\frac{2\pi}{3} + \cos^2\theta = \cos\theta \tag{3.30}$$

Solving for $\cos\theta$ gives $\cos\theta = -\frac{1}{3}$. Completing the tetrahedron, the final Stokes states (rounded to the thousands place) are:

$$
\begin{aligned}
&\begin{bmatrix} 1, & 0, & 0, & 1 \end{bmatrix} \\
&\begin{bmatrix} 1, & -0.471, & 0.816, & -0.333 \end{bmatrix} \\
&\begin{bmatrix} 1, & 0.943, & 0, & -0.333 \end{bmatrix} \\
&\begin{bmatrix} 1, & -0.471, & -0.816, & -0.333 \end{bmatrix}
\end{aligned} \tag{3.31}
$$

Converting these states to Jones vectors, the analyzer states we used (rounded to the thousands place) are given by:

$$
\begin{aligned}
&\begin{bmatrix} 0.707, & -0.707j \end{bmatrix} \\
&\begin{bmatrix} 0.514, & 0.794 + 0.324j \end{bmatrix} \\
&\begin{bmatrix} 0.986, & 0.169j \end{bmatrix} \\
&\begin{bmatrix} 0.514, & -0.794 + 0.324j \end{bmatrix}
\end{aligned} \tag{3.32}
$$

The squared overlap magnitudes between any of these states, $\beta^2 = \frac{1}{3}$ as desired for equally split analyzer states.

**Device Index of Refraction Profiles**

Optimized index of refraction profiles for the multispectral and angular momentum sorting devices are shown in Figure 3.11 and those for the Stokes polarimetry device from the main text and the one from the supplement with more layers are shown in Figure 3.12.

**Polarimetry Reconstruction**

The following section shows how the polarimetry device presented in the main text can be used to recover the Stokes parameters of arbitrarily polarized inputs. This addresses interpretation of quadrant outputs when the excitation is different than the four analyzer states used in the design. It further addresses the ability of the device

to utilize the four measurements to recover the degree of polarization for partially polarized light. This exploration is done in simulation, but the same calibration and reconstruction procedure can be used experimentally as well.

**Reconstruction Method**

The problem of converting the signal in each of the four quadrants into the incident polarization state can be phrased as follows:

$$\underline{\underline{M}}\mathbf{S} = \mathbf{T} \tag{3.33}$$

where $\underline{\underline{M}}$ is the forward model that maps the Stokes vector, $\mathbf{S}$, to the observed quadrant transmissions, $\mathbf{T}$. We utilize the common definition of the Stokes parameters:

$$\mathbf{S} = \begin{bmatrix} S_0 \\ S_1 \\ S_2 \\ S_3 \end{bmatrix} = \begin{bmatrix} E_x^2 + E_y^2 = E_{45}^2 + E_{-45}^2 = E_R^2 + E_L^2 \\ E_x^2 - E_y^2 \\ E_{45}^2 - E_{-45}^2 \\ E_R^2 - E_L^2 \end{bmatrix} \tag{3.34}$$

where $E_x$, $E_y$, $E_{45}$, and $E_{-45}$ are projections onto horizontal, vertical, 45-degree, -45-degree linear polarizations, respectively and $E_R$ and $E_L$ are projections onto right- and left-circular polarizations, respectively. To calibrate the device, we input each of these individual polarization components and observe the transmission into each of the four quadrants. Then, we form:

$$\underline{\underline{M}}\underline{\underline{\sigma}} = \underline{\underline{\tau}}$$

$$\underline{\underline{\sigma}} = \begin{bmatrix} \mathbf{S}_x & \mathbf{S}_y & \mathbf{S}_{45} & \mathbf{S}_{-45} & \mathbf{S}_R & \mathbf{S}_L \end{bmatrix} \in \mathbb{R}^{4x6}$$

$$\underline{\underline{\tau}} = \begin{bmatrix} \mathbf{T}_x & \mathbf{T}_y & \mathbf{T}_{45} & \mathbf{T}_{-45} & \mathbf{T}_R & \mathbf{T}_L \end{bmatrix} \in \mathbb{R}^{4x6} \tag{3.35}$$

$$\underline{\underline{M}} \in \mathbb{R}^{4x4}$$

where $\mathbf{S}_x = \begin{bmatrix} 1 & 1 & 0 & 0 \end{bmatrix}^\dagger$, $\mathbf{S}_y = \begin{bmatrix} 1 & -1 & 0 & 0 \end{bmatrix}^\dagger$, $\mathbf{S}_{45} = \begin{bmatrix} 1 & 0 & 1 & 0 \end{bmatrix}^\dagger$, $\mathbf{S}_{-45} = \begin{bmatrix} 1 & 0 & -1 & 0 \end{bmatrix}^\dagger$, $\mathbf{S}_R = \begin{bmatrix} 1 & 0 & 0 & 1 \end{bmatrix}^\dagger$, $\mathbf{S}_L = \begin{bmatrix} 1 & 0 & 0 & -1 \end{bmatrix}^\dagger$ and $\mathbf{T}_\alpha$ are the four quadrant transmissions under excitation by the the $\mathbf{S}_\alpha$ state. We solve for $\underline{\underline{M}}$ by taking the pseudo-inverse of $\underline{\underline{\sigma}}$ and applying it on the right side, $\underline{\underline{M}} = \underline{\underline{\tau}}\underline{\underline{\sigma}}^\dagger$. Then, we form the solution or reconstruction matrix by taking the inverse of $\underline{\underline{M}}$, such that given a set of measurements $\mathbf{T}$, we compute the Stokes parameters as $\mathbf{S} = \underline{\underline{M}}^{-1}\mathbf{T}$. We note this calibration could alternatively be done with the four analyzer states used in the design and we expect the results would be similar.

**Reconstructing Pure Polarization States**

The reconstruction method applied to pure polarization states is shown in Figure 3.13 for different amounts of added noise in the transmission measurements to simulate different signal-to-noise ratios in the sensor detection. For $p$ added noise, we add a normally distributed random variable with a mean of 0 and a standard deviation equal to $p * T_{\text{avg}}$ where $T_{\text{avg}}$ is the mean transmission across the four quadrant transmissions. As can be seen for increasing noise, the $S_3$ parameter is the most susceptible to a reduced signal-to-noise ratio. This is likely due to the circular polarization analyzer state exhibiting the lowest contrast and the $S_3$ Stokes parameter being a direct measure of the handedness of the circular polarization in the input.

**Reconstructing Mixed Polarization States**

The use of four projective measurements means information about partially polarized input states is contained in the quadrant transmissions. To test our ability to recover this property, we consider the situation where the polarization vector input into the device is randomly changing. We input a series of random polarization states into the device, and average the resulting quadrant transmission values for each quadrant. From these averaged transmission values, we reconstruct the Stokes vector in the same way as above. This reconstructed vector is compared to the averaged Stokes vectors for all the states input into the device. The degree of polarization of the light is computed as $p = \frac{\sqrt{S_1^2 + S_2^2 + S_3^2}}{S_0}$.

Figure 3.14 shows the results of reconstructing mixed polarization states. As the number of averaged states increases, the degree of polarization starts dropping. When noise is added per averaged state (using the same type of distribution as above), the squared error for the reconstruction is highest for the smaller number of averaged states. As this number of states increases, the fluctuating noise term starts averaging to zero, thus decreasing the overall effect of noise on the reconstruction.

**Angular Momentum Sorting Device Outside of Design Points**

Figure 3.15 and Figure 3.16 demonstrate the behavior of the angular momentum sorting device for different values of spin and OAM, respectively, than the design states. In an optical communication application, controlling the behavior of the device at these alternate points will depend on the amount noise present and mode distortion between communication links. However, in an advanced imaging con-

text where information about the scene is inferred through the spatially resolved projection of the input onto different angular momentum states, the response of the device to other mode inputs needs to be at least characterized if not explicitly designed for the given application. As a note, the optimization technique used here was not directed to explicitly minimize or control the behavior of the device under these other excitations. By adding more simulations to each iteration to capture the effect of illuminating with these other modes, we can compute a gradient that either enables control over the quadrant these other modes couple to or extinguishes their transmission.

**Illumination with Different Spin Values**

In Figure 3.15, we observe the device behaves similarly upon a flip in the handedness of the circular polarization for each angular momentum state. This can be seen through similar contrast and transmission profiles albeit at lower overall values. Thus, the optimization solution for the device relied primarily on the different OAM values for splitting and does not have strong polarization discriminating behavior.

**Illumination with Different OAM Values**

In Figure 3.16, we observe the device output changes drastically when illuminated with different OAM values. Most of the light for each of the four states goes to the quadrants designed for the original higher design OAM values (i.e., $l = -2, +2$). This is the reason for the negative contrast in the other two quadrants. Further, overall transmission values are significantly reduced with the higher transmission occurring for OAM values closer to the design points (i.e., $l = -3, +3$).

**3.7 Supplementary Figures**

Figure 3.5: Optical setup for characterization of multispectral and polarimetry devices. **(a)** Configuration for imaging of device focal plane and power normalization. Without the mirror in place, the lens images focal planes of the device onto the camera. Normalization of the device transmission is done with the mirror and power meter path of the setup. For these measurements, net power through an empty aperture is used to normalize net power through an aperture of the same size with the device on top of it. The power meter is aligned to the beam center, which is aligned to the pinhole centers during measurement. QCL: MIRcat-QT Mid-IR Quantum Cascade Laser (DRS Daylight Solutions); HWP1: Thorlabs WPLH05M-4500, Low-Order 4.5µm Half-Wave Plate; HWP2: Thorlabs WPLH05M-5300, Zero-Order 5.3µm Half-Wave Plate; QWP 1: Thorlabs WPLQ05M-4500, Low-Order 4.5µm Quarter-Wave Plate; LP: Thorlabs WP25M-IRA, Wire Grid Polarizer; L1: Thorlabs AL72525-E1, ZnSe Aspheric Lens, NA=0.42; L2: Thorlabs AL72512-E1, ZnSe Aspheric Lens, NA=0.67; Camera: Electrophysics PV320L IR Camera. **(b)** Configuration for verifying the polarization states used to test the Stokes polarimetry device. The second quarter-wave plate is moved to three distinct positions and the power in each linear polarization component separated by the Wollaston prism is recorded. QWP2: Thorlabs WPLQ05M-3500, Low-Order 3.5µm Quarter-Wave Plate; WP: Thorlabs WPM10, Wollaston Prism.

Figure 3.6: Stokes state creation and verification. **(a)** Polarization states are created through choice of angles of the linear polarizer, half-wave plate, and quarter-wave plate ($\theta_1$, $\theta_2$, and $\theta_3$). **(b)** Each state is verified by measuring the horizontal and vertical polarization component magnitudes output from the Wollaston prism after the state passes through a quarter-wave plate under three different rotations, $\phi$. **(c)** Plot of measured Jones vector overlap for the 4 analyzer states and their 4 orthogonal complements for each measurement wavelength used in the experiment.

Figure 3.7: Simulation performance for Stokes polarimetry device with additional degrees of freedom compared to Stokes polarimetry device from the main text. **(a)** Polarization contrast ($C$) and transmission ($T$) for device from the main text showing low contrast for the circular polarization analyzer state. For input $k$, $C_k = \frac{T_{|S_k\rangle \to Q_k} - T_{|\hat{S_k}\rangle \to Q_k}}{T_{|S_k\rangle \to Q_k} + T_{|\hat{S_k}\rangle \to Q_k}}$. **(b)** Polarization contrast and transmission for device with additional degrees of freedom showing high contrast for all four analyzer states at the cost of slightly reduced analyzer state transmission. **(c)** Focal intensity images for device from the main text with the top row showing the analyzer states and the bottom row showing their orthogonal complements. Intensity units are arbitrary but comparable between all plots in (c). The focal plane size is same as device aperture (30μm x 30μm). **(d)** Focal intensity image comparison for the device with additional degrees of freedom. Intensity units are arbitrary but comparable between all plots in (d). The focal plane size is the same as device aperture (30μm x 30μm).

Figure 3.8: Schematic of simulation geometry for optimization and evaluation. **(a)** Simulation geometry for optimization of the multispectral and Stokes polarimetry devices using a plane wave excitation. The angular momentum devices are optimized using focused angular momentum states with different circular polarization handedness for spin. **(b)** Evaluation geometry for the multispectral and Stokes polarimetry devices where the plane wave excitation is replaced with a defocused Gaussian source intended to match with the experimental source. Angular momentum devices are evaluated with the same sources as used for optimization.

Figure 3.9: Schematic of fabrication process. **(a)** Fabrication starts with a sapphire substrate (Al$_2$O$_3$, C-plane (0001), double side polished, 2-inch diameter, 0.5mm thickness). **(b)** Using a negative tone photoresist, apertures are patterned onto the substrate using photolithography. **(c)** After direct oxygen and argon plasma cleaning to remove undesired residual resist on the substrate, 150nm of Al is deposited on top using an electron beam evaporator. **(d)** The liftoff procedure is finished in acetone to remove remaining photoresist followed by cleaning in IPA and then DI water. **(e)** The IP-Dip resist from Nanoscribe is dropped onto the substrate. **(f)** Alignment is done by keeping the laser power below polymerization threshold and using fluorescence from its focused spot to align to the aperture centers for printing. **(g)** Development in propylene glycol methyl ether acetate (PGMEA) for 20 minutes followed by two three-minute rinses in IPA reveals the final device.



Figure 3.10: Conceptual diagram of the Stokes polarimetry device. The device acts to project an input state onto four outgoing states depending on its overlap with each analyzer state.

Figure 3.11: Index of refraction profiles for multispectral (b) and angular momentum (c) devices. Dark colors are IP-Dip polymer and light areas are void. **(a)** Schematic of geometry showing the location of each labeled layer. Each layer is $3\mu m$ thick for the multispectral device and 2.4μm thick for the angular momentum device. **(b)** Multispectral and linear polarization device index profile with each layer 30.15μm x 30.15μm. **(d)** Angular momentum device index profile with each layer 30.15μm x 30.15μm.

Figure 3.12: Index of refraction profiles for Stokes polarimetry device from main text (b) and Stokes polarimetry device with more degrees of freedom from supplement (c). Dark colors are IP-Dip polymer and light areas are void. **(a)** Schematic of geometry showing location of each labeled layer. Each layer is 3µm thick. **(b)** Stokes polarimetry device from main text index profile with each layer 30µm x 30µm. **(c)** Stokes polarimetry device with additional degrees of freedom index profile with each layer 30µm x 30µm.

Figure 3.13: Stokes polarimetry reconstruction in simulation using the device from the main text under random pure polarization inputs at a wavelength of 4.5µm. **(a)** Reconstructed state locations shown on Poincaré sphere on the left and comparison of the reconstructed Stokes parameters to the actual ones shown on the right with the associated squared error (green dashed line, right y-axis). **(b)** Same plots as (a) but with 5% added noise. **(c)** Same plots as (a) but with 10% added noise.

Figure 3.14: Stokes polarimetry reconstruction in simulation using the device from the main text under random mixed polarization inputs at a wavelength of 4.5µm. **(a)** Reconstructed state locations broken down by Stokes parameter as well as degree of polarization compared to actual ones shown on the left with the associated squared error per Stokes parameter shown on the right. **(b)** Same plots as (a) but with 5% added noise. **(c)** Same plots as (a) but with 10% added noise.

Figure 3.15: Simulation of angular momentum sorting device for the same OAM, but different spin values than the design states. **(a)** Schematic of device and focal plane quadrants. **(b)** Contrast for sorting each state ($C \in [-1, 1]$) defined by the transmission of a state into the desired quadrant versus the transmission into the rest of the focal plane. For source $k$, $C_k = \frac{T_{|S_k\rangle \to Q_k} - \sum_{i \neq k} T_{|S_k\rangle \to Q_i}}{T_{|S_k\rangle \to Q_k} + \sum_{i \neq k} T_{|S_k\rangle \to Q_i}}$. **(c)** Transmission spectrum for ($l = -2$, $s = -1$) input with desired quadrant transmission in blue. Transmission is normalized by power through the device aperture with no device present. Inset: Intensity at focal plane (arbitrary units, but comparable to other intensity plots in figure). **(d-f)** Same plots as (c), but for ($l = -1$, $s = -1$), ($l = +1$, $s = +1$), ($l = +2$, $s = +1$), respectively.

Figure 3.16: Simulation of angular momentum sorting device for the same spin values, but different OAM than the design points. **(a)** Schematic of device and focal plane quadrants. **(b)** Contrast for sorting each state ($C \in [-1, 1]$) defined by the transmission of a state into the desired quadrant versus the transmission into the rest of the focal plane. For source $k$, $C_k = \frac{T_{|S_k\rangle \to Q_k} - \sum_{i \neq k} T_{|S_k\rangle \to Q_i}}{T_{|S_k\rangle \to Q_k} + \sum_{i \neq k} T_{|S_k\rangle \to Q_i}}$. **(c)** Transmission spectrum for ($l = -4$, $s = +1$) input with desired quadrant transmission in blue. Transmission is normalized by power through the device aperture with no device present. Inset: Intensity at focal plane (arbitrary units, but comparable to other intensity plots in figure). **(d-f)** Same plots as (c), but for ($l = -3$, $s = +1$), ($l = +3$, $s = -1$), ($l = +4$, $s = -1$), respectively.

**References**

[1]G. Roberts, C. Ballew, T. Zheng, J. C. Garcia, S. Camayd-Muñoz, P. W. C. Hon, and A. Faraon, "3D-patterned inverse-designed mid-infrared metaoptics", Nature Communications **14**, 2768 (2023).

[2]H.-T. Chen, A. J. Taylor, and N. Yu, "A review of metasurfaces: physics and applications", Reports on progress in physics **79**, 076401 (2016).

[3]M. Miyata, N. Nemoto, K. Shikama, F. Kobayashi, and T. Hashimoto, "Full-color-sorting metalenses for high-sensitivity image sensors", Optica **8**, 1596–1604 (2021).

[4]O. D. Miller, *Photonic design: From fundamental solar cell physics to computational inverse design* (University of California, Berkeley, 2012).

[5]S. Molesky, Z. Lin, A. Y. Piggott, W. Jin, J. Vucković, and A. W. Rodriguez, "Inverse design in nanophotonics", Nature Photonics **12**, 659–670 (2018).

[6]M. P. Bendsøe and N. Kikuchi, "Generating optimal topologies in structural design using a homogenization method", Computer methods in applied mechanics and engineering **71**, 197–224 (1988).

[7]J. S. Jensen and O. Sigmund, "Topology optimization for nano-photonics", Laser & Photonics Reviews **5**, 308–321 (2011).

[8]L. Su, A. Y. Piggott, N. V. Sapra, J. Petykiewicz, and J. Vuckovic, "Inverse design and demonstration of a compact on-chip narrowband three-channel wavelength demultiplexer", Acs Photonics **5**, 301–305 (2018).

[9]C. M. Lalau-Keraly, S. Bhargava, O. D. Miller, and E. Yablonovitch, "Adjoint shape optimization applied to electromagnetic design", Optics express **21**, 21693–21701 (2013).

[10]P. I. Borel, A. Harpøth, L. H. Frandsen, M. Kristensen, P. Shi, J. S. Jensen, and O. Sigmund, "Topology optimization and fabrication of photonic crystal structures", Optics express **12**, 1996–2001 (2004).

[11]D. Sell, J. Yang, S. Doshay, R. Yang, and J. A. Fan, "Large-angle, multifunctional metagratings based on freeform multimode geometries", Nano letters **17**, 3752–3757 (2017).

[12]M. Mansouree, H. Kwon, E. Arbabi, A. McClung, A. Faraon, and A. Arbabi, "Multifunctional 2.5 D metastructures enabled by adjoint optimization", Optica **7**, 77–84 (2020).

[13]P. Camayd-Muñoz, C. Ballew, G. Roberts, and A. Faraon, "Multifunctional volumetric meta-optics for color and polarization image sensors", Optica **7**, 280–283 (2020).

[14]T. Gissibl, S. Thiele, A. Herkommer, and H. Giessen, "Two-photon direct laser writing of ultracompact multi-lens objectives", Nature photonics **10**, 554–560 (2016).

[15] S. Thiele, C. Pruss, A. M. Herkommer, and H. Giessen, "3D printed stacked diffractive microlenses", Optics Express **27**, 35621–35630 (2019).

[16] C. R. Ocier, C. A. Richards, D. A. Bacon-Brown, Q. Ding, R. Kumar, T. J. Garcia, J. Van De Groep, J.-H. Song, A. J. Cyphersmith, and A. Rhode, "Direct laser writing of volumetric gradient index lenses and waveguides", Light: Science & Applications **9**, 1–14 (2020).

[17] C. Roques-Carmes, Z. Lin, R. E. Christiansen, Y. Salamin, S. E. Kooi, J. D. Joannopoulos, S. G. Johnson, and M. Soljačić, "Toward 3D-Printed Inverse-Designed Metaoptics", ACS Photonics (2022).

[18] B. E. Bayer, "Color imaging array", United States Patent 3,971,065 (1976).

[19] S. Nishiwaki, T. Nakamura, M. Hiramoto, T. Fujii, and M.-a. Suzuki, "Efficient colour splitters for high-pixel-density image sensors", Nature Photonics **7**, 240–246 (2013).

[20] N. Zhao, P. B. Catrysse, and S. Fan, "Perfect RGB-IR Color Routers for Sub-Wavelength Size CMOS Image Sensor Pixels", Advanced Photonics Research **2**, 2000048 (2021).

[21] E. Johlin, "Nanophotonic color splitters for high-efficiency imaging", Iscience **24**, 102268 (2021).

[22] A. Schliesser, N. Picqué, and T. W. Hänsch, "Mid-infrared frequency combs", Nature photonics **6**, 440–449 (2012).

[23] J. Hodgkinson, R. Smith, W. O. Ho, J. R. Saffell, and R. P. Tatam, "Non-dispersive infra-red (NDIR) measurement of carbon dioxide at 4.2 $\mu$m in a compact and optically efficient sensor", Sensors and Actuators B: Chemical **186**, 580–588 (2013).

[24] S. Kang, Z. Qian, V. Rajaram, S. D. Calisgan, A. Alù, and M. Rinaldi, "Ultra-narrowband metamaterial absorbers for high spectral resolution infrared spectroscopy", Advanced Optical Materials **7**, 1801236 (2019).

[25] A. B. Seddon, "Mid-infrared (IR)–A hot topic: The potential for using mid-IR light for non-invasive early detection of skin cancer in vivo", physica status solidi (b) **250**, 1020–1027 (2013).

[26] M. J. Baker, J. Trevisan, P. Bassan, R. Bhargava, H. J. Butler, K. M. Dorling, P. R. Fielden, S. W. Fogarty, N. J. Fullwood, and K. A. Heys, "Using Fourier transform IR spectroscopy to analyze biological materials", Nature protocols **9**, 1771–1791 (2014).

[27] A. Wang and Y. Dan, "Mid-infrared plasmonic multispectral filters", Scientific reports **8**, 1–7 (2018).

[28] H. Zhao, Y. Li, G. Jia, N. Li, Z. Ji, and J. Gu, "Comparing analysis of multispectral and polarimetric imaging for mid-infrared detection blindness condition", Applied Optics **57**, 6840–6850 (2018).

[29]F. Wang, B. S. Lazarov, and O. Sigmund, "On projection methods, convergence and robust formulations in topology optimization", Structural and multidisciplinary optimization **43**, 767–784 (2011).

[30]D. B. Fullager, G. D. Boreman, and T. Hofmann, "Infrared dielectric response of nanoscribe IP-dip and IP-L monomers after polymerization from 250 cm-1 to 6000 cm-1", Optical Materials Express **7**, 888–894 (2017).

[31]X. Zhou, Y. Hou, and J. Lin, "A review on the processing accuracy of two-photon polymerization", Aip Advances **5**, 030701 (2015).

[32]F. Ottonello-Briano, C. Errando-Herranz, H. Rödjegård, H. Martin, H. Sohlström, and K. B. Gylfason, "Carbon dioxide absorption spectroscopy with a mid-infrared silicon photonic waveguide", Optics Letters **45**, 109–112 (2020).

[33]B. Kunnen, C. Macdonald, A. Doronin, S. Jacques, M. Eccles, and I. Meglinski, "Application of circularly polarized light for non-invasive diagnosis of cancerous tissues and turbid tissue-like scattering media", Journal of biophotonics **8**, 317–323 (2015).

[34]K. P. Gurton, A. J. Yuffa, and G. W. Videen, "Enhanced facial recognition for thermal imagery using polarimetric imaging", Optics letters **39**, 3857–3859 (2014).

[35]A. Kadambi, V. Taamazyan, B. Shi, and R. Raskar, "Polarized 3d: High-quality depth sensing with polarization cues", in Proceedings of the ieee international conference on computer vision (2015), pp. 3370–3378.

[36]N. J. Pust and J. A. Shaw, "Digital all-sky polarization imaging of partly cloudy skies", Applied optics **47**, H190–H198 (2008).

[37]D. Wang, H. Liang, H. Zhu, and S. Zhang, "A bionic camera-based polarization navigation sensor", Sensors **14**, 13006–13023 (2014).

[38]H. G. Berry, G. Gabrielse, and A. E. Livingston, "Measurement of the Stokes parameters of light", Applied optics **16**, 3200–3205 (1977).

[39]J. S. Tyo, D. L. Goldstein, D. B. Chenault, and J. A. Shaw, "Review of passive imaging polarimetry for remote sensing applications", Applied optics **45**, 5453–5469 (2006).

[40]J. Bai, C. Wang, X. Chen, A. Basiri, C. Wang, and Y. Yao, "Chip-integrated plasmonic flat optics for mid-infrared full-Stokes polarization detection", Photonics Research **7**, 1051–1060 (2019).

[41]E. Arbabi, S. M. Kamali, A. Arbabi, and A. Faraon, "Full-Stokes imaging polarimetry using dielectric metasurfaces", Acs Photonics **5**, 3132–3140 (2018).

[42]N. A. Rubin, G. D'Aversa, P. Chevalier, Z. Shi, W. T. Chen, and F. Capasso, "Matrix Fourier optics enables a compact full-Stokes polarization camera", Science **365**, eaax1839 (2019).

[43]H. Zhang, C. W. Hsu, and O. D. Miller, "Scattering concentration bounds: brightness theorems for waves", Optica **6**, 1321–1327 (2019).

[44]*Materials, methods, and additional text are available in the supplementary.*

[45]L. Allen, M. W. Beijersbergen, R. J. C. Spreeuw, and J. P. Woerdman, "Orbital angular momentum of light and the transformation of Laguerre-Gaussian laser modes", Physical review A **45**, 8185 (1992).

[46]A. E. Willner, K. Pang, H. Song, K. Zou, and H. Zhou, "Orbital angular momentum of light for communications", Applied Physics Reviews **8**, 041312 (2021).

[47]J. Wang, J.-Y. Yang, I. M. Fazal, N. Ahmed, Y. Yan, H. Huang, Y. Ren, Y. Yue, S. Dolinar, and M. Tur, "Terabit free-space data transmission employing orbital angular momentum multiplexing", Nature photonics **6**, 488–496 (2012).

[48]N. Bozinovic, Y. Yue, Y. Ren, M. Tur, P. Kristensen, H. Huang, A. E. Willner, and S. Ramachandran, "Terabit-scale orbital angular momentum mode division multiplexing in fibers", science **340**, 1545–1548 (2013).

[49]H. Ren, X. Li, Q. Zhang, and M. Gu, "On-chip noninterference angular momentum multiplexing of broadband light", Science **352**, 805–809 (2016).

[50]L. Torner, J. P. Torres, and S. Carrasco, "Digital spiral imaging", Optics express **13**, 873–881 (2005).

[51]X. Fang, H. Ren, and M. Gu, "Orbital angular momentum holography for high-security encryption", Nature Photonics **14**, 102–108 (2020).

[52]Z. Lin, C. Roques-Carmes, R. Pestourie, M. Soljačić, A. Majumdar, and S. G. Johnson, "End-to-end nanophotonic inverse design for imaging and polarimetry", Nanophotonics **10**, 1177–1187 (2021).

[53]N. Tétreault, G. von Freymann, M. Deubel, M. Hermatschweiler, F. Pérez-Willard, S. John, M. Wegener, and G. A. Ozin, "New route to three-dimensional photonic bandgap materials: silicon double inversion of polymer templates", Advanced Materials **18**, 457–460 (2006).

[54]V. Hahn, P. Kiefer, T. Frenzel, J. Qu, E. Blasco, C. Barner-Kowollik, and M. Wegener, "Rapid assembly of small materials building blocks (voxels) into large functional 3D metamaterials", Advanced Functional Materials **30**, 1907795 (2020).

[55]L. Hong, H. Li, H. Yang, and K. Sengupta, "Fully integrated fluorescence biosensors on-chip employing multi-functional nanoplasmonic optical structures in CMOS", IEEE Journal of Solid-State Circuits **52**, 2388–2406 (2017).

[56]R. Fatemi, C. Ives, A. Khachaturian, and A. Hajimiri, "Subtractive photonics", Optics Express **29**, 877–893 (2021).

[57]Y. Zhao, Y. Zhao, S. Hu, J. Lv, Y. Ying, G. Gervinskas, and G. Si, "Artificial structural color pixels: A review", Materials **10**, 944 (2017).

[58]Y.-H. Lee, T. Zhan, and S.-T. Wu, "Prospects and challenges in augmented reality displays.", Virtual Real. Intell. Hardw. **1**, 10–20 (2019).

[59]E. D. Palik, *Handbook of optical constants of solids*, Vol. 3 (Academic press, 1998).

[60]J. K. Guest, "Topology optimization with multiple phase projection", Computer Methods in Applied Mechanics and Engineering **199**, 123–135 (2009).

*Chapter 4*

# VISIBLE LIGHT COLOR SPLITTERS

## 4.1 Image Sensor Architecture

Cameras are a vital part of modern life, integrated into a variety of consumer, scientific and industrial electronics including cell phones, laptops, augmented and virtual reality systems, remote sensing devices, and robotics and vision-based autonomous systems. Most cameras share a common anatomy with a series of lenses forming an image onto a 2D grid of photosensitive elements, or pixels. Even lens-less computational imaging systems rely on an image sensor array for detection [1]. Typically, each pixel is sensitive to a band of wavelengths depending on many factors, but strongly influenced by the underlying detection mechanism and material composition. While the detection bandwidth varies, image sensors are often only weakly dependent on wavelength and operate as grayscale detectors. For a scene illuminated by broadband light, spatially resolved spectral data is obtained passively by placing filters above individual pixels. Commercial sensors often contain absorptive filters based on organic or pigment dyes [2, 3]. The color filtering effect has also been achieved using photonic structures like, for example, metal hole arrays or wire grids that have enhanced transmission for bands that can be spectrally tuned by via the metal/air geometry [4–7].

In complementary metal–oxide–semiconductor (CMOS) image sensors, a type of active-pixel sensor consisting of an array of photodetectors and amplifying MOS field effect transistors (MOSFETs) [8–10], the trend has been towards reductions in pixel size [11] leading to an increasing need for high transmission efficiency for maximizing signal to noise especially for low light and high-speed imaging conditions. Many improvements have been made over the years to improve efficiency including increasing the fractional area of the light sensitive region of the pixel and creating back-side illuminated arrays [12–14].

## 4.2 Nanophotonic Color Splitting

One of the largest potential improvements for increasing transmission efficiency in color image sensors is capturing spectral information without the need for absorptive and/or reflective filtering mechanisms. Color filtering, while efficient at capturing tailor-able spectral band overlaps at each pixel, rejects all signal outside of each

sensing band. A famous color filtering pattern, the Bayer pattern [15], is a quad of color filters (red, blue, and two green). The average transmission efficiency of this arrangement is bound to be 33% at a maximum. A solution to this low efficiency is to create a device that can sit on top of the same four pixels and route different bands of light to different detectors. This can capture all of the same color information while not sacrificing any overall signal. There has been extensive research in this area [16], only a few of which are highlighted here. Initial work used simple wavelength scale silicon nitride (SiN) structures embedded in silicon dioxide ($SiO_2$) that exhibited non uniform scattering by wavelength [17]. Although transmission was high and color could be reconstructed as the solution to a simple linear system, the chromatic effect was weak, thus strongly trading off signal and color information. Other work has employed metasurfaces with spatially multiplexed or engineered chromatic phase masks to realize different focusing profiles as a function of wavelength [18, 19]. Single-layer metasurface color routers have also been discovered via inverse design techniques where the desired chromatic output is specified and the whole design space is searched for a high performing solution. For example, a genetic algorithm can carry out a global search on single layer devices where the binary pattern encoding the presence of one of two materials is evolved heuristically. The population of devices is not constrained to behave like a metasurface, allowing a broader range of scattering phenomena that can lead to higher performance color routing [20]. Even with this additional design freedom, the problem is complex and the single layer device performance is often lacking in either overall transmission efficiency, color discrimination or both. Using the color routing device as a replacement for the microlens array, but keeping the absorptive color filters, one can recover robust color discrimination while also increasing transmission [21]. However, theoretically if the routing device is working well, there should be no need for any additional filtering positioning this solution as a great near-term commercial option with the promise of better color routing devices in the future with even higher transmission efficiencies.

Increasing the degrees of freedom available in the design has been shown to enable highly efficient solutions to the color routing problem. Adding patterned layers to the metasurface design space, such that light can be controlled and sorted over a longer propagation length, has gained increasing interest from the community. This 3D design space, where each layer can be patterned with subwavelength features, even in the binary material constraint, is exponentially larger than the single-layer one. Further, each simulation takes longer to run since devices are thicker in the ax-

ial direction and multiple scattering effects make accurate approximate solvers more difficult to construct. To optimize in this massive design space in a reasonable computational budget, many approaches use efficient gradient computations combined with local optimizers to evolve designs on single trajectories toward high performing solutions. When equipped with exceedingly small feature sizes ($10 - 20$nm laterally and 40nm axially) over a $0.64 \times 0.64 \times 2$µm thickness (i.e., 50 design layers), a periodic array of devices illuminated with a plane wave can exhibit nearly perfect routing efficiency [22]. These devices are not simply replacements for the color filter array, but also replace multiple elements like the microlens array, antireflective coating, and could be designed to additionally operate as an infrared cutoff filter, which is found in many CMOS image sensors to reduce infrared light that may be absorbed by the pixels.

## 4.3   3D Visible Light Designs

**Isolated Optimization**

Color routing devices naturally live in an array, but as was the case for the mid-infrared we study first some of their behavior in an isolated setting. Using a design space of titanium dioxide ($TiO_2$) and $SiO_2$, two CMOS-compatible materials, and embedding the structure in a background of $SiO_2$, we aim to find a structure that consists of just 5 patterned material layers while restricting to a feature size on the order of 60nm. The assumed refractive indices of $TiO_2$ and $SiO_2$ were 2.6 and 1.5, respectively. As was the case before, the feature size arises from a filtering kernel, in this case a continuous and differentiable approximation to a morphological dilation operation detailed in the supplementary material of the published version of this work [23]. This dilation filter works on the higher index material phase in the design and is not simultaneously controlled in the lower index phase. There are a handful of features in the final design underneath a strict 60nm feature size limit, which if removed have been seen in simulation to only slowly degrade device performance. In other words, the function of the device does not hinge completely on the presence of these features. The next subsection features an optimization using a fixed design grid where the minimum design voxel size is equal to the feature size in both domains to guarantee no feature is below the critical dimension in both materials. Assuming the average wavelength to be $\lambda_0 = 550$nm, the device properties for the isolated device optimization are given in Table 4.1.

The optimization goal, as outlined in [23] is to focus each band of wavelengths to the centers of different quadrants at the focal plane. This is done via a figure

| Property | Value |
|---|---|
| Device lateral dimension, $w$ | $2\mu m \times 2\mu m$ ($3.64 \times 3.64\lambda_0$) |
| Device thickness, $t$ | $2\mu m$ ($3.64\lambda_0$) |
| Number of layers (layer thickness), N | 5 (400nm) |
| Focal length, $f$ | $1.5\mu m$ ($2.73\lambda_0$) |
| Minimum feature size, $\delta$ | $\approx 60nm$ ($\approx 0.11\lambda_0$) |

Table 4.1: Isolated Device Properties

of merit that maximizes intensity at each quadrant center for the different spectral bands. Given the presence in typical image sensors of two 'green' pixels, we opted to showcase the flexibility of the design technique and sort the middle band by linear polarization. This concept would augment a typical camera sensor with a function not easily achievable with just absorptive filters. To achieve this effect via filtering, the sensor would require additional micro-polarizer elements, which would further reduce overall transmission efficiency while adding fabrication complexity.

The optimized device is shown in Figure 4.1a,b with a rendered, simulated view of the focal plane under broadband, uniform x-polarized plane wave illumination in Figure 4.1c. The sorting efficiency is plotted in Figure 4.1d where each curve corresponds to quadrant transmission referenced to the whole input aperture. Color contrast is defined here as the difference in transmission to the target quadrant for a given wavelength and the highest transmission to any of the other three quadrants ($C_{\text{color}} = T_{\text{target}} - \max_{k \neq \text{target}} T_k$). Polarization contrast is defined similarly, except only for the middle band as the difference between the transmission to the two polarization sensitive 'green' quadrants. These contrast metrics are averaged spectrally and plotted as a function of the number of layers used in the design in Figure 4.1e. For each point in the plot, the optimization was repeated using a different number of patterned layers (and consequently different overall thickness) to show the ability of the algorithm to improve contrast with increased degrees of freedom and propagation length. The focal length and layer thickness was kept constant, so the single layer device is 400nm thick. Increasing the size of the design space by adding patterned layers improves performance under this optimization technique. It is important to keep in mind the local nature of the optimization method when drawing conclusions from results like these. The optimized devices only provide lower bounds on performance and results in the literature suggest that global or more rational design techniques outperform the results shown here for a single patterned layer albeit with slightly different materials, device geometries, and target function

[19, 20, 24]. Global techniques will be discussed and explored more later in the chapter. For now, we note that while the thinner devices may be under-performing with this design method that it is an attractive feature of the method that often simply adding design-able thickness to the device can steadily add performance. This is not a general truth, but when it is, one can directly trade off performance for fabrication complexity and cost.

The optimizations and evaluations in these simulations were carried out using a total field scattered field (TFSF) plane wave source. It is a common simulation technique for simulating the effect of a theoretically infinite illumination source in a finite simulation region with pefectly matched layer (PML) boundary conditions. Putting a terminated plane wave source in the optimization region would lead to non-ideal illumination since the edges of the source would cause diffraction artifacts in the simulation. The challenge of the TFSF source, however, is that it is difficult to normalize against when computing transmission since the source is technically infinite in extent, but the device is finite. The power used for transmission calculations is normalized against the injected power in the lateral extent of the TFSF source. However, to maintain the assumption the source appears to be an infinite plane wave, the device is able to pull in power from a wider aperture than just its geometrical cross section [25]. Thus, additional power can enter the focal plane through this process, which means the transmission values can sum up to greater than unity. This is still physical since in isolation, the device would scatter in this way. However, in an array, the power outside of the device geometrical area would ideally be scattering to the neighboring device focal plane. In the next section, we address this shortcoming and other practicalities moving towards making industry-ready devices.

Figure 4.1: Result of 5-layer isolated device optimization. (a) Exploded view of the five lithographic layers required to form the final device. (b) An example device layer showing the designed inclusions of $TiO_2$ in the background $SiO_2$. (c) Depiction of normally incident plane wave illumination exhibiting one of two linear polarizations and rendering of apparent device focal plane color under broadband excitation. (d) Sorting efficiency for each quadrant as a function of input wavelength for x-polarized illumination. Each colored curve corresponds to the transmission spectrum from light incident on the device to the quadrant. The 'blue' curve corresponds to the quadrant targeted towards sensing the shortest wavelengths. The middle band is broken into both green and yellow depending on the quadrant targeted for each linear polarization. Here, the 'green' curve is the desired quadrant for x-polarized illumination. (e) Efficiency and contrast plotted as function of the number of layers used in the design. Each point corresponds to running the optimization procedure for a different number of layers (with fixed focal length and layer thickness). As can be seen, at least up to 5 layers, adding more degrees of freedom in this way steadily improves the device performance. Adapted with permission from [23] ©The Optical Society. Figure data from author and figure artwork by Sarah Camayd-Muñoz with rewritten caption from original work.

**Replicated Boundary Optimization**

Color routing devices tiled on a CMOS image sensor can be designed with a locally periodic assumption. In other words, the same design can be used across small areas of the sensor array. Moving spatially on the sensor corresponds to receiving light from the lens at a different chief ray angle. Since each optimized device has a nonzero angular bandwidth of operation, the same device can be used for small spatial shifts on the sensor, corresponding to small input angle deviations. In contrast, a large spatial shift requires re-designing the device around a new center angle. The use of a periodic boundary condition in the simulation, however, should be treated with care especially when used during the optimization procedure. In periodic simulations, both the device and the sources are replicated periodically. The plane wave excitation illuminates the entire array and light transmitted to the focal plane of one device could have been incident originally on a different device in the array. There is no protection against this type of spatial crosstalk and in fact it can be encouraged by the optimization if it leads to a higher transmission solution. The adjoint sources are periodically replicated and so the device has the option to increase its efficiency by coupling light out to any focal plane, not just its own.

To account for the presence of a neighboring device in the optimization while still avoiding the perils of a fully periodic simulation region, we can include a small slice of the neighboring device on all borders. Neighboring devices are assumed to be the same as the optimization target. Via a perfect electrical conductor (PEC) aperture, we section off illumination to only excite a single device and not directly its neighbors. The illumination is Gaussian with a waist radius of 1.0164µm, corresponding to an approximation for an Airy spot for a wavelength of 550nm from a lens with an f-number of 2.2. Since the neighboring device is not directly illuminated, it can also enter into the optimization of the central device because it can serve to help isolate light that might otherwise escape out the device sides. However, the neighbors are constrained to be the same device as the central one. Thus, we can think of the optimization of the central device as having two purposes: (1) funnel and sort the light directly incident on the top to the correct quadrant in the focal plane and (2) have a structure at its border such that when its tiled in an array it retains high efficiency in the first objective. The second purpose is optimized indirectly by using a target figure of merit defined by the first objective and by including gradient information from the border regions in the optimization. This is depicted in Figure 4.2. The gradient of the objective function is computed in the central device region as well as the border regions. Then, since the external

border of the central device will change when its borders are evolved, the gradient from outside the device is folded back onto the appropriate region of the central device. For example, the right-hand external bordering region gradient is added to the left-hand border of the central device. This is done for the entire device border.



Figure 4.2: Device bordering in optimization and remapping of gradient information. (a) The central device is surrounded by partial replications of itself of size 0.255µm on all sides. (b) Gradient information computed outside the device region can be remapped onto the central device to compute a complete gradient for the central device that accounts for its array nature. Each lettered label shows how the external gradient is remapped onto the central device to be summed with the gradient at those locations.

The design space for the device is similar to the isolated case with some minor differences. Index of refraction of $SiO_2$ and $TiO_2$ are assumed to be 1.5 and 2.4, respectively, with the background still being made of $SiO_2$ in the simulation. The feature size is fixed to 51nm by reinterpolating the gradient onto a grid with this resolution. This is in contrast to the previous design where the feature size was encouraged via a feature blurring filter. A filtering approach sacrifices less degrees of freedom to achieve a certain feature size because it allows features to be placed with finer resolution than the minimum dimension. It further does not restrict features to be only multiples of the minimum dimension. However, the core drawback is a lack of guarantee of the final feature size especially in both material domains. To circumvent this issue, the design below simply restricts the optimization to a grid with the resolution of the minimum feature size. The properties used in the design are (with respect to $\lambda_0 = 550$nm) given in Table 4.2.

The optimization, as demonstrated in the mid-infrared section, seeks to maximize

| Property | Value |
|---|---|
| Device lateral dimension, $w$ | 2.04µm × 2.04µm (3.71 × 3.71$\lambda_0$) |
| Device thickness, $t$ | 2.04µm (3.71$\lambda_0$) |
| Number of layers (layer thickness), N | 10 (204nm) |
| Focal length, $f$ | 1.53µm (2.78$\lambda_0$) |
| Minimum feature size, $\delta$ | 51nm (0.93$\lambda_0$) |

Table 4.2: Replicated Device Properties

intensity in the center of each quadrant depending on the input wavelength. Since we are evaluating the performance of the device for a standard CMOS image sensor, we drop the multifunctionality of sorting the middle band by polarization. The optimization problem, then, is defined over each quadrant, $q$, and polarization, $p$:

$$\max_{\boldsymbol{\epsilon} \in \{\epsilon_{\min}, \epsilon_{\max}\}^N} g(\mathbf{E}) = \sum_\lambda \sum_p \sum_q w_q(\lambda) \frac{I_p(\mathbf{r}_q, \lambda)}{I_{\max}(\lambda)}$$

$$I_p(\mathbf{r}_q, \lambda) = ||\mathbf{E}_p(\mathbf{r}_q, \lambda)||^2 \tag{4.1}$$

$$w_q(\lambda) = e^{\frac{-(\lambda - \nu_q)^2}{\alpha \beta_q(\lambda)^2}}$$

where $\alpha = \frac{1}{\ln 2}$ and $\nu_q$ = 452nm, 541nm, 648nm for the 'blue', 'green', and 'red' quadrants, respectively. For 'blue' and 'red', $\beta_q(\lambda)$ = 23.7nm and 30.5nm, respectively. For 'green', $\beta_q(\lambda)$ = 23.7nm if $\lambda <$ 541nm and 30.5nm if $\lambda \geq$ 541nm. The gradients are further weighted by a performance-based weighting function, similar to the mid-infrared optimizations. For each wavelength, the performance is computed as a function of wavelength based on the transmission, $T$, into each quadrant, $q$:

$$f(\lambda) = \sum_p \sum_q w_q(\lambda) T_{p,q}(\lambda) \tag{4.2}$$

The gradient for each quadrant, polarization, and wavelength is multiplied by a performance based weight, $\gamma(\lambda)$, that depends only on the wavelength (where the number of discrete wavelengths in the optimization is given by $N$):

$$\gamma(\lambda) = \frac{2}{N} - \frac{f(\lambda)^2}{\sum_\nu f(\nu)^2} \tag{4.3}$$

Similar to before, if this is negative at any point it is renormalized by first subtracting the minimum value and then dividing by the sum of the resulting array.

A plot of the optimized refractive index for the device is plotted in Figure 4.3 with the darker regions showing $TiO_2$ inclusions in the white $SiO_2$ background. The device results are plotted for both the optimization condition and the array condition, the difference being the size of the border is set equal to a full device when simulating the array. This is depicted in Figure 4.4a while the focal plane arrangement to identify with colored transmission curves is in 4.4b. The transmission spectra are shown for both input polarizations in 4.4c,d. First, the difference between the array and optimization conditions are shown, with considerable peak performance difference between the two. We suspect the agreement will improve with a thicker device border used in the optimization, although at the cost of total optimization time. The array condition is shown in more detail in 4.4d where one can see a small polarization effect in the two green quadrants by observing the difference in the 'green' and 'magenta' curves for the middle band. If the two 'green' quadrants are pooled together, the dashed curve demonstrates the device does not have a strong overall polarization preference.



Figure 4.3: Index of refraction for 10-layer replicated border device. (a) Diagram showing location of the layers in the optimized design with respect to the input illumination and the focal plane. (b) Optimized index profile for each layer with the darker regions indicating $TiO_2$ and the white regions indicating $SiO_2$.

Figure 4.4: Transmission spectra from 10-layer replicated border device optimization. (a) Device configurations for evaluation with the optimization condition shown on the left and the array condition shown on the right. The optimization condition has a 0.255µm slice of the neighboring device included on all borders while the array condition replicates a full copy of the neighboring device on all borders. These correspond to the partially transparent and fully solid lines in (c) below. (b) The focal plane arrangement for the transmission spectra plotted. Where the 'magenta' is absent, the two 'green' pixel transmissions have been summed together. (c) Transmission spectra under Gaussian illumination normalized to transmission of beam through empty aperture. Each polarization excitation is shown on top and bottom. The difference in the solid and partially transparent lines shows the remaining inaccuracy in the partially replicated device border optimization condition to the more complete array condition. (d) Further breakdown of the transmission by quadrant for the array condition showing a small amount of polarized behavior between the two 'green' quadrants. The vertical dotted lines indicate the band crossover points.

The device reflection is shown in Figure 4.5. On the left side of Figure 4.5a is the simulation configuration for measuring reflection off the device. It contains backscattering that is a combination of that from the aperture and the device since the input illumination overfills the aperture area. Thus, on the right side of Figure 4.5a is a normalization simulation to show how much light would reflect off the aperture without any device present. The transmission spectra for each polarization is shown in Figure 4.5b with the normalization plots as well. From this it can be seen that reflection off the top of the device surface is on average 11.9% for both input polarizations. This value is low considering that without a device present, 8.5% of the input illumination reflects off the aperture.

Light making it through the device does not exclusively reach its focal plane and some of it is scattered to neighboring focal planes. In Figure 4.6a, the simulation condition is shown alongside the focal plane arrangement for the scattering plots. For each polarization, the scattering plot in Figure 4.6b shows different transmission channels. Transmissions are averaged across each band, cutoffs for which can be seen in Figure 4.4c,d. In the focal plane of the central device, each quadrant's spectral response is indicated by the relative radii of the colored circles (i.e., the top right quadrant mostly contains the 'blue' band of wavelengths). Arrows pointing from the central focal plane indicate the relative magnitude of transmission scattering to that neighboring focal plane. The scattering by band tends to be slightly higher for neighboring focal planes nearest to the target quadrant for that band. Transmission by band to each focal plane is also indicated.

Figure 4.5: Reflection spectra from 10-layer replicated border device optimization. (a) Configurations corresponding to plots in (b). The full reflection from the device in the array configuration as depicted by the solid purple line while the normalization spectra in the dotted purple and green lines depict the amount of transmission and reflection with no device present. The reflection with no device is particularly useful because it shows how much light from the overfilled beam reflects back off the aperture directly. (b) Reflection plots for horizontal and vertical input polarizations, respectively.

Figure 4.6: Scattering response from 10-layer replicated border device optimization. (a) Device configuration on left and focal plane arrangement corresponding to the plots in (b). (b) Scattering for horizontal (left side) and vertical (right side) input polarizations. The center illustrates the device focal plane with the relative amount of each color of light in each quadrant depicted by scaling the radius of the circle proportional to the average transmission in each band. The length of the arrows pointing to each neighboring focal plane is scaled proportional to the average transmission in each band scattering to that focal plane. The average transmission by band in each focal plane is also written.

While the aperture will not be present in the final device, its usage allows us to get a sense of the quantitative view of the device performance. With an understanding of the various scattering channels and overall device efficiency, the optimization can be repeated without the aperture. Overall, the performance in the design is promising in both the optimization and array condition. Performance can be further improved in the array condition by replicating more of the device boundary during the optimization. This will increase the time needed to run the optimization, but will more accurately model the device in its eventual array configuration.

## 4.4   CMOS Foundry Fabrication Process

Prototyping free space inverse designed devices for visible and near infrared wavelengths is challenging due to the required fabrication complexity. While assumed feature sizes and material refractive indices are compatible with modern industrial fabrication competencies used to manufacture advanced integrated circuits, experimental realization of designs with precise layer-to-layer alignment with limited process control and equipment in an academic cleanroom setting is challenging. In electronics, CMOS multi project wafer (MPW) services allow academic and fabless companies to access industrial scale nanofabrication at lower costs by sharing the chip area and thus mask set cost with a large number of other projects. While an equivalent MPW does not currently exist for making 3D nanophotonic devices, there has been work showing the ability to make optical devices out of electronics chips. Integrated circuits consist of a layer of devices (i.e., a transistor layer) connected by a large 3D interconnect matrix of metallic wires, typically primarily made of copper. Even older technology nodes contain on the order of 10 interconnect layers with the thinnest feature size wires closest to the transistors and bulkiest ones on top for routing longer range connections in the chip and connecting to external wire bonding pads. These wires have previously been used to create integrated plasmonic filters [26]. In more recent work, metal interconnects were etched away to realize dielectric structures in a process the authors call subtractive photonics [27].

We propose a postprocessing method that would allow multilayer, dielectric photonic devices to be fabricated from CMOS MPW services. For mechanical robustness during postprocessing and for design simplification, we consider extruded 2D devices that operate in free space, but are modeled to behave as if they are infinite in one direction. In fully 3D devices, design rules will limit the types of shapes that can be reliably fabricated. For 2D designs, we need to consider only the minimum allowable interconnect width and pitch. The proposed method, sketched in Figure

Figure 4.7: CMOS foundry MPW postprocessing scheme showing the side and top views of the chip at each step. (a) Start of the process showing a silicon substrate, device layer, multiple layers of interconnect, via and capping layers, background oxide, and wire bonding pads. (b) Using the pads as alignment marks, an etch mask is placed on the backside of the substrate. (c) A series of wet etchants are used to remove the oxide and diffusion barrier surrounding the copper leaving a 3D copper structure. (d) Oxide of desired refractive index is backfilled into the design to surround the copper structure. It is planarized after backfilling so the top part of the interconnect is exposed. (e) The silicon substrate is wet or dry etched and the etch mask is then stripped. This leaves the bottom of the device exposed so that it can be measured in transmission. (f) The copper is wet etched from the oxide leaving a dielectric/air device.

4.7 would allow these extruded 2D designs to be prototyped. First, pads on the top of the device would be used to coarsely align an etch mask on the backside of the substrate that will be used later to open a window through the otherwise absorbing silicon. Then, the background oxide, capping/etch stop layers, and diffusion barrier materials surrounding the copper wires can be etched away. The diffusion barrier is usually made from tantalum (Ta) or tantalum nitride (TaN) and is used to prevent

copper from diffusing into the oxide surrounding the wires [28]. While these barriers are thin, they are also conductive and will affect the scattering properties and loss inside the device. Wet etchants for Ta and TaN often also etch oxide or other materials inside the stack. While not necessary in the proposed method here, other work appears to have been able to remove the diffusion barrier material without etching the oxide [27]. After this step, a freestanding copper structure remains. The via layers on either end of the extruded copper sections in this scaffold need to be used to provide mechanical stability. Then, the desired design material can be backfilled into the structure, assumed below to be $SiO_2$. The etch mask deposited in the first step can be used to etch a window in the silicon substrate. The copper wires exposed on both sides of the structure can be removed with a copper wet etchant. The remaining device is composed now of the backfilled oxide material and air with air on either side for optical measurements.

## 4.5   Foundry Compatible Extruded 2D Devices

Based on the proposed postprocessing technique, we aim to design devices that respect the constraints of a candidate MPW process in simulation. We work in the extruded 2D design space to simplify the application of the fabrication restrictions to satisfying the layer thicknesses and feature pitch in both material phases. This would also allow interconnect material to be used for support structures and mechanical stability in the postprocessing method.

### Local Optimization

We first use gradient-based optimization techniques to find solutions for splitting three wavelength bands on the focal plane. The device is made of $SiO_2$ and air with refractive indices 1.5 and 1, respectively, with a background index in the simulation region of 1. One of the more difficult aspects of utilizing the interconnects is the coarser feature size and pitch relative to the wavelength compared to what has been used in previous designs. This will be different depending on the technology node, but thinner wires have larger resistance and more closely spaced wires have higher capacitance, both of which reduce operation speed. Interconnect layers further from the transistors are often used to route over further distances on the chip, so they tend to be restricted to thicker sizes and pitches. The device parameters, referenced to $\lambda_0 = 550nm$, are given in Table 4.3.

The layer closest to the transistors is the thinnest with the smallest feature size and all other layers have coarser dimensions. Features are constructed on a fixed grid with

| Property | Value |
|---|---|
| Device width, $w$ | 3.6µm ($6.55\lambda_0$) |
| Device thickness, $t$ | 2.5µm ($4.55\lambda_0$) |
| Number of designable layers, N | 7 |
| Layer thicknesses | 130, 220nm |
| Spacer thicknesses | 175nm |
| Focal length, $f$ | 1.8µm ($3.27\lambda_0$) |
| Minimum feature size, $\delta$ | 90, 100nm ($0.16, 0.18\lambda_0$) |

Table 4.3: Foundry Restricted Device Parameters

spacing equal to the minimum feature size, thus guaranteeing no features in either material domain come out smaller than the minimum. Devices are optimized using TFSF plane wave illumination, but evaluated using an apertured (PEC) Gaussian input with a waist radius of 1.29µm. Transmissions are normalized to transmission through this aperture with no device present. For the same reason as in the bordered device case, this illumination ensures a more interpretable transmission metric.

The optimization is carried out for both single, $E_z$, and dual, $(E_z, H_z)$, polarization conditions. The optimization problem, for polarizations, $p$, and focal locations, $q$, is:

$$
\max_{\boldsymbol{\epsilon} \in \{\epsilon_{\min}, \epsilon_{\max}\}^N} g(\mathbf{E}) = \sum_\lambda \sum_p \sum_q w_q(\lambda) \frac{I_p(\mathbf{r}_q, \lambda)}{I_{\max}(\lambda)}
$$
$$
I_p(\mathbf{r}_q, \lambda) = ||\mathbf{E}_p(\mathbf{r}_q, \lambda)||^2
$$

(4.4)

where $w_q(\lambda) = 1$ if $\lambda$ is in the band for location $q$ and 0 otherwise. Bands are evenly split in [400nm, 700nm] and optimizations are run with two unique mappings of band to focal spot. Specifically, in one configuration, the spots are ordered 'blue', 'red', 'green' (Configuration A) and in the other configuration, the spots are ordered 'blue', 'green', 'red' (Configuration B). The performance weighting scheme for the gradients is applied first based on the intensity figure of merit as a function of optimized wavelength in a given band for each polarization separately. Then, the weighting scheme is performed again for polarization depending on the average performance for each polarization.

### $E_z$ Polarization, Configuration A

In the first design, we consider just $E_z$ polarization and a configuration of focal spots that places 'red' between 'blue' and 'green.' Figure 4.8a,b show the transmission

spectra of the optimized device under both polarizations. While not optimized for $H_z$, we still see color splitting behavior show up when illuminating with this polarization. Scattering channels indicate, the key for which are shown in Figure 4.8c alongside the optimized refractive index profile, that there is fairly even loss out the side and to reflection with some amount of reflection due to the beam spilling over the sides of the aperture. While a lot of light makes it to the focal plane, the band definitions are not sharp. For example, at the dotted vertical red line in Figure 4.8a for $E_z$ polarization, the transmission to the 'red' spot is 48.1% while it has only cut off to 19.5% to the 'green' spot. Empirically, we find restriction to coarse feature sizes often associated with diminished ability to realize sharp band cutoffs. Further, in these optimizations, we only optimized for high intensity in the correct focal spot for a given wavelength and did not use an augmented figure of merit to further reduce intensity for incorrect wavelengths to those same focal spots, which may show improvements in this contrast metric. The field intensities are also plotted showing the directing of different colors of light to different focal spots.

**($E_z$, $H_z$) Polarization, Configuration A**

Next, we consider adding the $H_z$ polarization to the optimization. Results of the optimization in Figure 4.9 follow the same format to Figure 4.8. The $H_z$ polarized output now has higher peak transmission at the cost of lowering the $E_z$-polarized transmission.

**($E_z$, $H_z$) Polarization, Configuration B**

We consider the effect of the configuration of focal spots on the dual-polarized design, here optimizing for the 'blue', 'green', 'red' arrangement as indicated in Figure 4.10c. In this configuration, the overall sorting contrast appears to improve in the $H_z$ polarization at the expense of the contrast and overall transmission for $E_z$ polarization. The greater color crosstalk in Figure 4.9 for $E_z$ polarization was between the 'green' and 'red' bands. Spatially, these focal spots are placed next to each other on the focal plane while the 'blue' and 'green' spots are maximally spaced. Moving the 'green' spot to the middle of the focal plane now places it next to the 'blue' spot which seems to open up a larger crosstalk channel between the two bands. This can be seen in the field intensity plots in Figure 4.10a.

Figure 4.8: Optimization results for $E_z$ polarization, configuration A. (a) Transmission and field intensity plots for $E_z$ polarization normalized to transmission through an empty aperture with the same Gaussian source. The left plot shows the transmission to each focal quadrant and the middle plot shows the transmission to various scattering channels. Panel (c) shows the location of transmission monitors corresponding to curves in these plots. On the right is the field magnitude corresponding to the dotted vertical lines in the transmission plots. Each plot is individually normalized by its maximum intensity. (b) Same plots as in (a), but for $H_z$ polarization. (c) The left figure shows the simulation geometry and serves as a key for various scattering channels corresponding to the transmission plots. On the right is the optimized device refractive index profile.

Figure 4.9: Optimization results for ($E_z$, $H_z$) polarization, configuration A. (a) Transmission and field intensity plots for $E_z$ polarization normalized to transmission through an empty aperture with the same Gaussian source. The left plot shows the transmission to each focal quadrant and the middle plot shows the transmission to various scattering channels. Panel (c) shows the location of transmission monitors corresponding to curves in these plots. On the right is the field magnitude corresponding to the dotted vertical lines in the transmission plots. Each plot is individually normalized by its maximum intensity. (b) Same plots as in (a), but for $H_z$ polarization. (c) The left figure shows the simulation geometry and serves as a key for various scattering channels corresponding to the transmission plots. On the right is the optimized device refractive index profile.
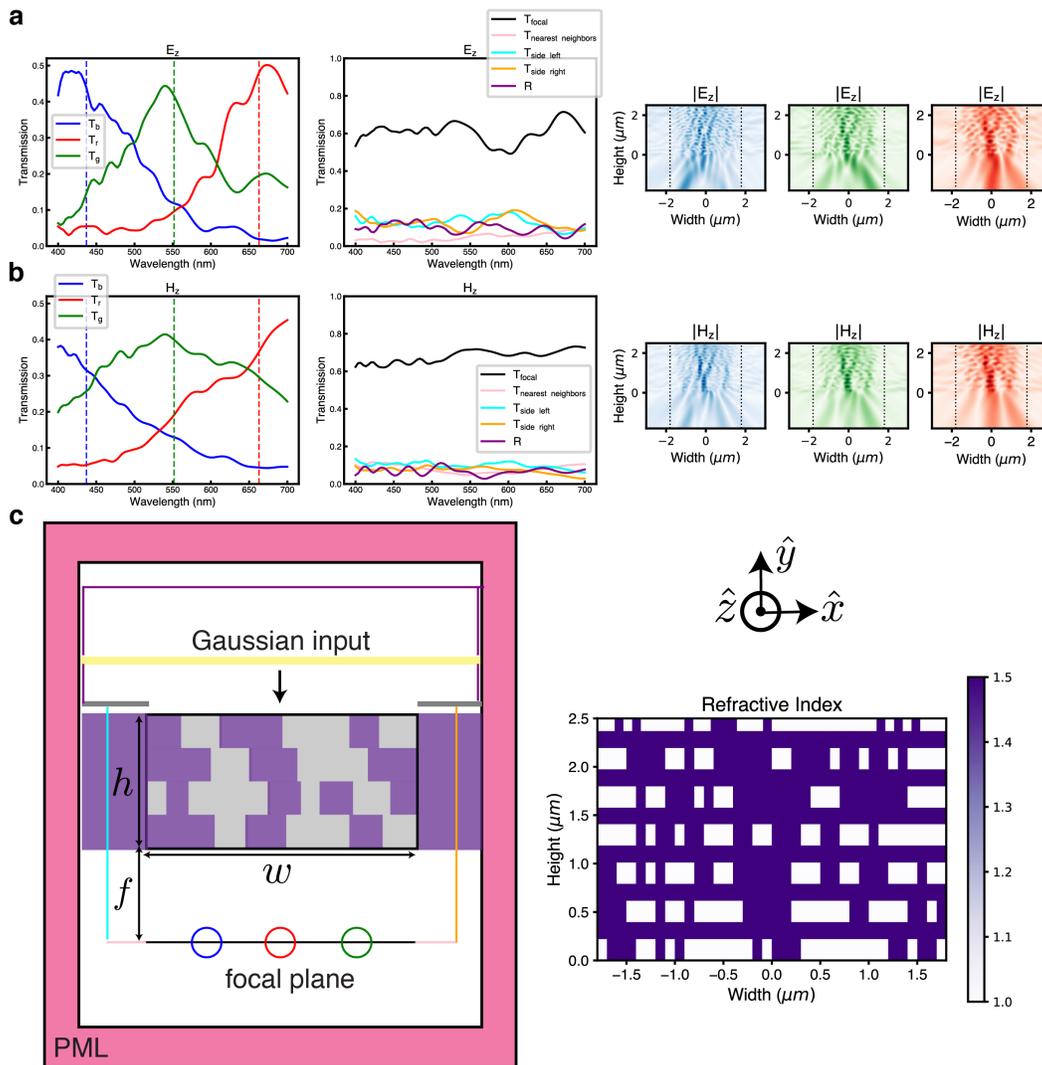
Figure 4.10: Optimization results for ($E_z$, $H_z$) polarization, configuration B. (a) Transmission and field intensity plots for $E_z$ polarization normalized to transmission through an empty aperture with the same Gaussian source. The left plot shows the transmission to each focal quadrant and the middle plot shows the transmission to various scattering channels. Panel (c) shows the location of transmission monitors corresponding to curves in these plots. On the right is the field magnitude corresponding to the dotted vertical lines in the transmission plots. Each plot is individually normalized by its maximum intensity. (b) Same plots as in (a), but for $H_z$ polarization. (c) The left figure shows the simulation geometry and serves as a key for various scattering channels corresponding to the transmission plots. On the right is the optimized device refractive index profile.
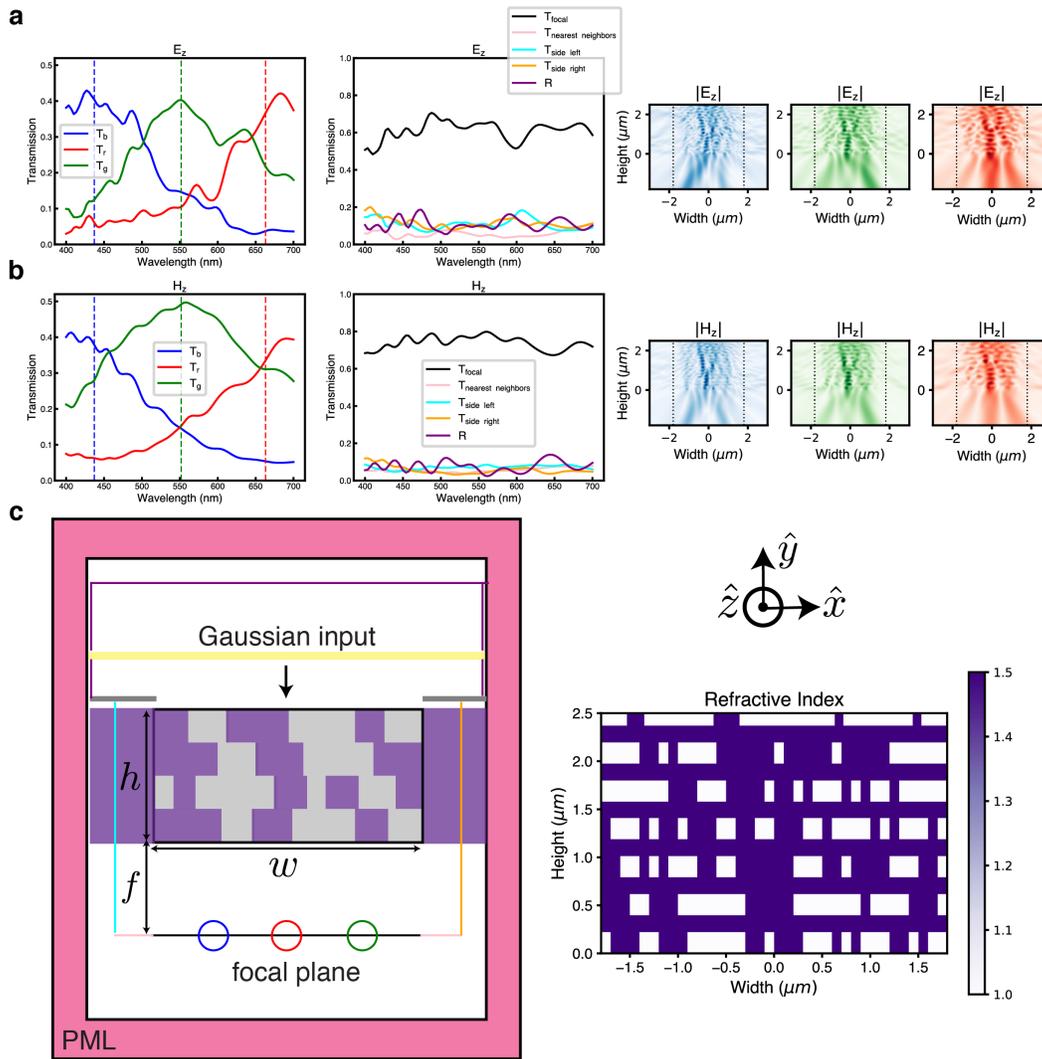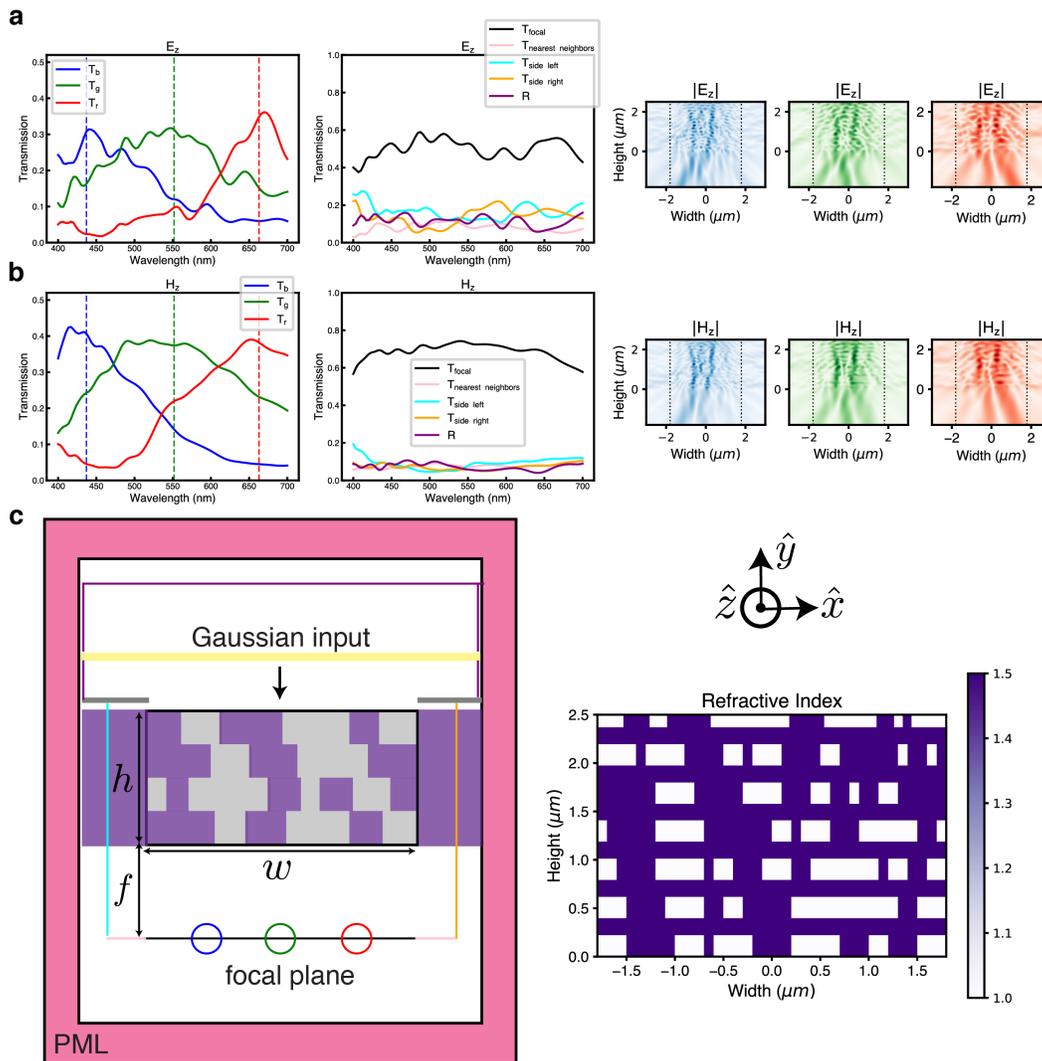
**Global Optimization**

Electromagnetic inverse design is a nonlinear optimization problem, typically with many local optima that may be settled on in gradient-driven or local optimization techniques. The variance in quality of the local optima is usually unknown and non-trivial, physics-based performance upper bounds to compare against optimization results are not straightforwardly derived for all problem structures. Further, while the continuous phase of an optimization may reach high performance when able to utilize a fully grayscale refractive index, pushing towards a nearby binary solution may significantly affect the overall merit function [29]. Many optimizations also contain hyperparameters that need to be selected like choice of design materials, focal length, width of device, and thickness of individual layers. These parameter choices do not always have associated gradients and even when they do, it is not always fruitful to choose their values with only local information.

Global approaches search a larger portion of the design space, usually at the cost of overall optimization time, to find higher performing solutions. Some of these approaches build the global search on top of gradient information and local optimizers. Examples of this include using a variety of initial device seeds, optimizing hyperparameters via particle swarm with particle performance defined by the result of a local optimization, and backpropagating gradient information for a set of devices produced by a generative neural network to the network model parameters so that it can learn to create a family of designs from which to choose [30, 31]. Other global optimization techniques may not even rely directly on gradients and instead on heuristic and statistical approaches to searching the design space [32, 33].

A drawback of global optimization is the additional computational cost since they often require more device evaluations or amount to running large numbers of local optimizations. For simpler design spaces with fewer degrees of freedom or problems that can take advantage of fast solvers, these methods can converge to high performing solutions in reasonable optimization times. For multilayer 3D design spaces with currently available solvers, many global optimization approaches exceed computational budgets. However, when possible, it is useful to assess if higher performing solutions exist compared to a locally optimized design especially when in a challenging design space. Here, given the results of two of the optimizations presented above, we utilize a direct binary search (DBS) method that has been proposed and established as useful tool for inverse design [34]. In this technique, the design is assumed at all times to be binary, so we take as seed the binarized

devices at the end of the optimizations. Then, for each pass through the device, a random path steps through each of the design voxels. The voxel is first flipped to the other material domain and the change is kept only if the figure of merit improves. Otherwise, the optimization proceeds to the next location. This requires running one simulation for every design voxel, which does not scale well to larger problems. In the case of these optimizations, though, with faster simulations and less design voxels given the 2D simplification, we were able to make 5 passes through the whole design to see the improvement in overall performance.

The DBS method was applied to the two devices optimized for focal spot configuration A. Instead of using the intensity figure of merit to evaluate device performance, the average transmission to each focal location across the band and the desired polarizations was used. This is the more desirable quantity to maximize, but it is associated with a more complex adjoint source compared to the focusing figure of merit. However, in the global search, we can easily evaluate the transmission performance of the device. In Figure 4.11 and Figure 4.12 are the results for the $E_z$ polarization and dual polarization devices, respectively. Transmission curves in panels a,b of both figures plot the device performance before the DBS method is employed for comparison. The algorithm found, in both cases, several changes that increased transmission as seen in the bottom of 4.11d and 4.12d. Likely due to a lack of uniformity enforcement, which is a benefit of the weighting function in the local optimization, some instances of improved transmission in one band came at the cost of lower transmission and sometimes loss of contrast for other bands. For example, 4.12b shows that in the dual-polarization case, for $H_z$ input, the 'red' band transmission increased significantly, but continued to cut off slowly into the 'green' band thus wiping out a lot of contrast for the middle part of the spectrum. Using a figure of merit that further accounts for balance between each band would help in this global search. This can be achieved, for example, with a product figure of merit instead of a sum. Instead of summing together transmission performance for each band, the product figure of merit multiplies each individual merit function together, the result being that the largest improvements are made by increasing the worst performing individual. Global search techniques also allow the ability to craft more complex figures of merit to fully define desired performance. In local optimizers, some figures of merit may have difficult gradients to compute or may not be readily improved with gradient information. The results here demonstrate an ability to also combine local and global techniques where a large part of the performance comes from the more efficient local optimizer while the global optimizer is used as a final

step to enhance performance or satisfy additional constraints and metrics.

Figure 4.11: Optimization results for $E_z$ polarization, configuration A. (a) Transmission and field intensity plots for $E_z$ polarization normalized to transmission through an empty aperture with the same Gaussian source. The left plot shows the transmission to each focal quadrant compared to the transmission before applying DBS and the middle plot shows the transmission to various scattering channels. Panel (c) shows the location of transmission monitors corresponding to curves in these plots. On the right is the field magnitude corresponding to the dotted vertical lines in the transmission plots. Each plot is individually normalized by its maximum intensity. (b) Same plots as in (a), but for $H_z$ polarization. (c) Simulation geometry and serves as a key for various scattering channels corresponding to the transmission plots. (d) Optimized device refractive index profile on the top and the change in refractive index after 5 DBS passes compared to the index at the end of the local optimization.
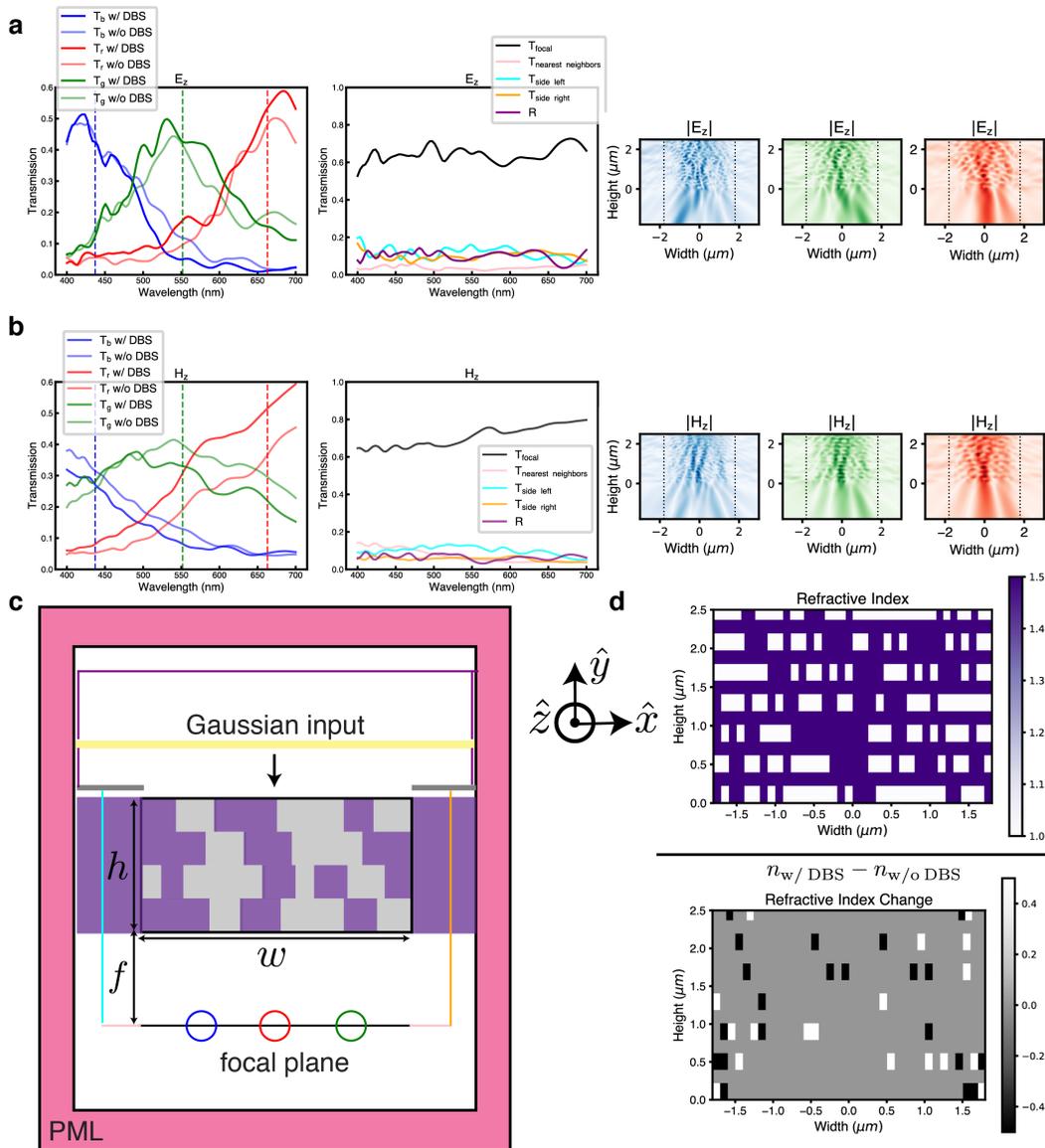
Figure 4.12: Optimization results for $(E_z, H_z)$ polarization, configuration A. (a) Transmission and field intensity plots for $E_z$ polarization normalized to transmission through an empty aperture with the same Gaussian source. he left plot shows the transmission to each focal quadrant compared to the transmission before applying DBS and the middle plot shows the transmission to various scattering channels. Panel (c) shows the location of transmission monitors corresponding to curves in these plots. On the right is the field magnitude corresponding to the dotted vertical lines in the transmission plots. Each plot is individually normalized by its maximum intensity. (b) Same plots as in (a), but for $H_z$ polarization. (c) Simulation geometry and serves as a key for various scattering channels corresponding to the transmission plots. (d) Optimized device refractive index profile on the top and the change in refractive index after 5 DBS passes compared to the index at the end of the local optimization.
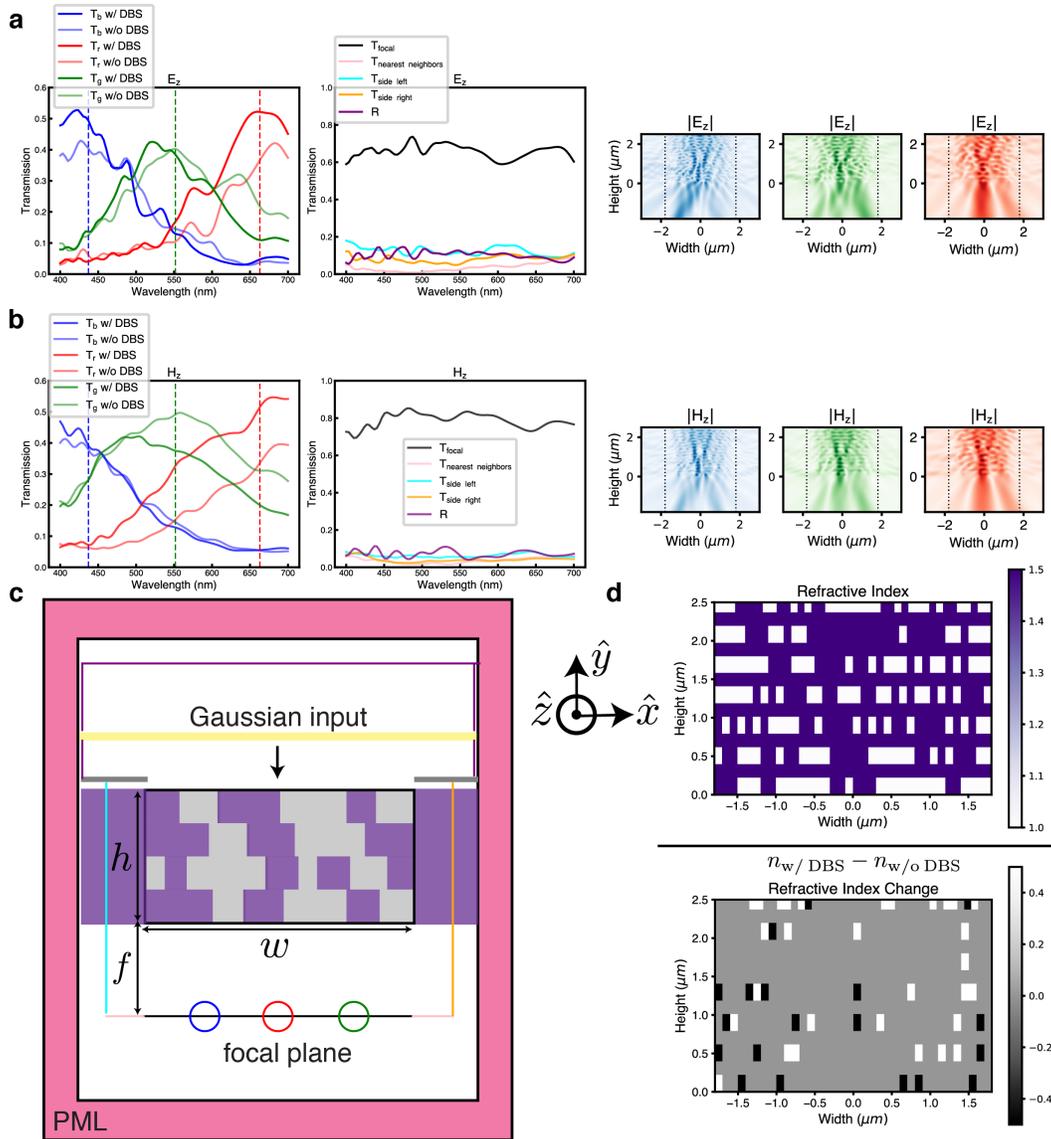
## 4.6 Conclusion

The visible light design space for volumetric metaoptic devices is challenging due to fabrication requirements. It also contains vast applications and promise for practical technologies, making it a worthwhile pursuit. Beyond color imaging on CMOS image sensors, augmented and virtual reality displays, multispectral remote sensing, active optical devices for beam steering and wavefront shaping, fiber optic and non-traditional imaging geometries, high color purity structural color, and spectrum sorters and angular filters for solar energy efficiency all will benefit from efficient, compact, visible light, volumetric metaoptics [35–42]. Fabrication remains challenging to prototype new devices in academic and research settings. Foundries to create these devices for industrial applications will drive forward progress in this area much as it did for integrated circuits. Access to foundry processes via MPW organization can bring those advanced fabrication techniques to the research community to create the next generation of metaoptics.

## References

[1] N. Antipa, G. Kuo, R. Heckel, B. Mildenhall, E. Bostan, R. Ng, and L. Waller, "DiffuserCam: lensless single-exposure 3D imaging", Optica **5**, 1–9 (2018).

[2] H. Taguchi, "Technology of color filter materials for image sensor", IISW, 2011, 34–37 (2011).

[3] J. Adams, K. Parulski, and K. Spaulding, "Color processing in digital cameras", IEEE micro **18**, 20–30 (1998).

[4] S. Yokogawa, S. P. Burgos, and H. A. Atwater, "Plasmonic color filters for CMOS image sensor applications", Nano letters **12**, 4349–4354 (2012).

[5] T. Xu, Y.-K. Wu, X. Luo, and L. J. Guo, "Plasmonic nanoresonators for high-resolution colour filtering and spectral imaging", Nature communications **1**, 59 (2010).

[6] P. B. Catrysse and B. A. Wandell, "Integrated color pixels in 0.18-µm complementary metal oxide semiconductor technology", JOSA A **20**, 2293–2306 (2003).

[7] T. W. Ebbesen, H. J. Lezec, H. F. Ghaemi, T. Thio, and P. A. Wolff, "Extraordinary optical transmission through sub-wavelength hole arrays", nature **391**, 667–669 (1998).

[8] P. J. W. Noble, "Self-scanned silicon image detector arrays", IEEE Transactions on electron Devices **15**, 202–209 (1968).

[9] E. R. Fossum, S. Mendis, and S. E. Kemeny, "Active pixel sensor with intra-pixel charge transfer", (2003).

[10]E. R. Fossum, "CMOS image sensors: Electronic camera-on-a-chip", IEEE transactions on electron devices **44**, 1689–1698 (1997).

[11]R. J. Gove, "CMOS image sensor technology advances for mobile devices", in *High performance silicon imaging* (Elsevier, 2020), pp. 185–240.

[12]B. Pain, T. Cunningham, S. Nikzad, M. Hoenk, T. Jones, C. Wrigley, and B. Hancock, "A back-illuminated megapixel CMOS image sensor", (2005).

[13]B. Pain, "Fabrication and initial results for a back-illuminated monolithic APS in a mixed SOI/bulk CMOS technology", in Ieee workshop on ccd & ais, 2005 (2005), pp. 102–104.

[14]Y. Oike, "Evolution of image sensor architectures with stacked device technologies", IEEE Transactions on Electron Devices **69**, 2757–2765 (2021).

[15]B. E. Bayer, "Color imaging array", United States Patent 3,971,065 (1976).

[16]Q. Chen, X. Nan, M. Chen, D. Pan, X. Yang, and L. Wen, "Nanophotonic color routing", Advanced Materials **33**, 2103815 (2021).

[17]S. Nishiwaki, T. Nakamura, M. Hiramoto, T. Fujii, and M.-a. Suzuki, "Efficient colour splitters for high-pixel-density image sensors", Nature Photonics **7**, 240–246 (2013).

[18]B. H. Chen, P. C. Wu, V.-C. Su, Y.-C. Lai, C. H. Chu, I. C. Lee, J.-W. Chen, Y. H. Chen, Y.-C. Lan, and C.-H. Kuan, "GaN metalens for pixel-level full-color routing at visible light", Nano letters **17**, 6345–6352 (2017).

[19]M. Miyata, N. Nemoto, K. Shikama, F. Kobayashi, and T. Hashimoto, "Full-color-sorting metalenses for high-sensitivity image sensors", Optica **8**, 1596–1604 (2021).

[20]X. Zou, Y. Zhang, R. Lin, G. Gong, S. Wang, S. Zhu, and Z. Wang, "Pixel-level Bayer-type colour router based on metasurfaces", Nature Communications **13**, 3288 (2022).

[21]S. Yun, S. Roh, S. Lee, H. Park, M. Lim, S. Ahn, and H. Choo, "Highly efficient color separation and focusing in the sub-micron CMOS image sensor", in 2021 ieee international electron devices meeting (iedm) (2021), pp. 30–31.

[22]P. B. Catrysse, N. Zhao, W. Jin, and S. Fan, "Subwavelength Bayer RGB color routers with perfect optical efficiency", Nanophotonics **11**, 2381–2387 (2022).

[23]P. Camayd-Muñoz, C. Ballew, G. Roberts, and A. Faraon, "Multifunctional volumetric meta-optics for color and polarization image sensors", Optica **7**, 280–283 (2020).

[24]X. Zou, G. Gong, Y. Lin, B. Fu, S. Wang, S. Zhu, and Z. Wang, "Metasurface-based polarization color routers", Optics and Lasers in Engineering **163**, 107472 (2023).

[25]J. D. Jackson, *Classical electrodynamics*, 1999.

[26]L. Hong, S. McManus, H. Yang, and K. Sengupta, "A fully integrated CMOS fluorescence biosensor with on-chip nanophotonic filter", in 2015 symposium on vlsi circuits (vlsi circuits) (2015), pp. C206–C207.

[27]R. Fatemi, C. Ives, A. Khachaturian, and A. Hajimiri, "Subtractive photonics", Optics Express **29**, 877–893 (2021).

[28]R. Hubner, *Advanced Ta-based diffusion barriers for Cu interconnects* (Nova Science Publishers: Lamont, 2008).

[29]O. Teytaud, P. Bennet, and A. Moreau, "Discrete global optimization algorithms for the inverse design of silicon photonics devices", Photonics and Nanostructures-Fundamentals and Applications **52**, 101072 (2022).

[30]J. Jiang and J. A. Fan, "Global optimization of dielectric metasurfaces using a physics-driven neural network", Nano letters **19**, 5366–5372 (2019).

[31]H. Chung and O. D. Miller, "Tunable metasurface inverse design for 80% switching efficiencies and 144 angular deflection", Acs Photonics **7**, 2236–2243 (2020).

[32]S. Jafar-Zanjani, S. Inampudi, and H. Mosallaei, "Adaptive genetic algorithm for optical metasurfaces design", Scientific reports **8**, 11040 (2018).

[33]P.-I. Schneider, X. Garcia Santiago, V. Soltwisch, M. Hammerschmidt, S. Burger, and C. Rockstuhl, "Benchmarking five global optimization approaches for nano-optical shape optimization and parameter reconstruction", ACS Photonics **6**, 2726–2733 (2019).

[34]B. Shen, P. Wang, R. Polson, and R. Menon, "An integrated-nanophotonics polarization beamsplitter with 2.4× 2.4 $\mu$m2 footprint", Nature Photonics **9**, 378–382 (2015).

[35]G.-Y. Lee, J.-Y. Hong, S. Hwang, S. Moon, H. Kang, S. Jeon, H. Kim, J.-H. Jeong, and B. Lee, "Metasurface eyepiece for augmented reality", Nature communications **9**, 4562 (2018).

[36]X. Zhang, F. Zhang, Y. Qi, L. Deng, X. Wang, and S. Yang, "New research methods for vegetation information extraction based on visible light remote sensing images from an unmanned aerial vehicle (UAV)", International Journal of Applied Earth Observation and Geoinformation **78**, 215–226 (2019).

[37]P. L. Hatfield and P. J. Pinter Jr, "Remote sensing for crop protection", Crop protection **12**, 403–413 (1993).

[38]S.-Q. Li, X. Xu, R. Maruthiyodan Veetil, V. Valuckas, R. Paniagua-Domínguez, and A. I. Kuznetsov, "Phase-only transmissive spatial light modulator based on tunable dielectric metasurface", Science **364**, 1087–1090 (2019).

[39]J. E. Fröch, L. Huang, Q. A. A. Tanguy, S. Colburn, A. Zhan, A. Ravagli, E. J. Seibel, K. F. Böhringer, and A. Majumdar, "Real time full-color imaging in a meta-optical fiber endoscope", eLight **3**, 1–8 (2023).

[40]Y. Zhao, Y. Zhao, S. Hu, J. Lv, Y. Ying, G. Gervinskas, and G. Si, "Artificial structural color pixels: A review", Materials **10**, 944 (2017).

[41]A. Mojiri, R. Taylor, E. Thomsen, and G. Rosengarten, "Spectral beam splitting for efficient conversion of solar energy—A review", Renewable and Sustainable Energy Reviews **28**, 654–663 (2013).

[42]O. Höhn, T. Kraus, G. Bauhuis, U. T. Schwarz, and B. Bläsi, "Maximal power output by solar cells with angular confinement", Optics Express **22**, A715–A722 (2014).

*C h a p t e r   5*

# OPTIMIZATION LANDSCAPE VISUALIZATION

## 5.1    Introduction

To this point, the main focus has been on optimizing a permittivity distribution that maximizes the value of a figure of merit given a fixed geometry. In some cases, this geometry has been fixed, at least in part, by the choice of application or the available fabrication method. For example, in visible light color splitting targeted for CMOS image sensors, the lateral dimension is fixed by the assumed underlying pixel size. However, even in this case, there are many choices left for the designer including focal length, material choices, number of patterned layers, and device height. Parameter choices that influence the final optimization result but are not the direct optimization parameters are referred to as hyperparameters. For some hyperparameters, there are logical ranges from which to choose dictated by physical intuition. When optimizing for a focal spot in a given pixel, for example, the device width and focal length should be chosen such that the diffraction limited spot of an equivalent numerical aperture lens is at least contained within the target pixel size. However, other choices that influence the final device performance are not always immediately apparent and running large parameter sweeps to directly search for their optimal values can be prohibitively expensive computationally.

There are certain physical intuitions that typically guide choices of device parameters when setting up optimizations. Devices with large index of refraction contrast have been at the core of dielectric metasurface research [1]. High index contrast allows for full phase coverage and confinement of light to individual nanoposts such that coupling is minimal between elements. The wavelength of light in a medium scales inversely with the index of refraction, so for a given thickness, a higher index feature appears thicker. Larger optical thickness is typically helpful when making increasingly advanced devices. Thicker devices allow more propagation length and interaction with patterned material through which different degrees of freedom can be sorted and processed. Empirically, starting from a single or few layer design space, adding layers to inverse design optimizations tends to have a positive impact on performance [2]. Especially in the case of fully binarized designs, it is important how finely the index can be patterned in both the lateral and axial direction. Lateral

feature sizes are usually associated with lithographic patterning resolution while axial feature sizes can be controlled via deposition thickness. They are usually chosen independently in designs and thus far we have been utilizing relatively thicker layer sizes compared to minimum lateral feature size. These fabrication considerations are mostly independent except that in the limit of increasingly thick individual layers with fine feature sizes where there may be aspect ratio dependent etching limitations on pattern fidelity. Patterning increasingly finer feature sizes across a fixed overall lateral and axial dimensions offers the optimization larger numbers of degrees of freedom and the ability to reach more effectively grayscale index values by virtue of being deeply subwavelength.

Despite physical intuition behind which direction to push geometric and optical parameters to improve performance, we must also contend with the nature of local optimization and the possibility that even given a design space that can theoretically efficiently couple to and thus control a greater number of modes, the optimizer may settle on a lower performing solution. The severity of this problem is dependent on many factors including the form and regularization of the objective function, the underlying optimizer, properties of the governing physics, and geometrical device constraints. Index of refraction contrast and average index, minimum feature size, layer thickness, and overall device thickness are examples of physical properties that affect the optimization process. In a continuous optimization space, for the same average index of refraction, a greater available index contrast contains all the designs using a lower contrast and thus it is expected that the optimization solution is no worse by increasing contrast. However, in practical scenarios, the final device needs to be binary and in this case, the lower index contrast solutions are not part of the higher index binary design space. Choosing hyperparameter values therefore requires physical insight such that high performing solutions are possible as well as construction of a design space where these solutions are readily discovered via local optimization.

The following sections use techniques from machine learning to explore where design spaces with seemingly more physical degrees of freedom lead to lower performing optimization results. By visualizing the figure of merit for sample permittivity distributions and the behavior of devices near local optima, we seek to understand why certain optimizations are more difficult than others. In machine learning, the figure of merit is often called the loss function and the values of this function for every configuration of the learnable parameters describes what is

referred to as the loss landscape. Here, we adopt this term for photonic optimizations, often shortening to just 'landscape.'

## 5.2 Device Configurations

To explore cases that show both sides of this optimization reality, we consider a variety of device configurations. In Configuration A, we consider a baseline set of parameters given in Table 5.1 for $\lambda_0 = 512.5$nm.

| Property | Value |
|---|---|
| Device width, $w$ | 1.872µm (3.65$\lambda_0$) |
| Device thickness, $t$ | 1.008µm (1.97$\lambda_0$) |
| Focal length, $f$ | 1.404µm (2.74$\lambda_0$) |
| Minimum feature size | 36nm (0.070$\lambda_0$) |
| Layer thickness, $t_{\text{layer}}$ | 36nm (0.070$\lambda_0$) |

Table 5.1: Device Configuration A Properties

Then, we check the effect of reducing vertical patterning resolution in Configuration B. This means the optimization has half the number of parameters with respect to permittivity design voxels and less vertical resolution, with the same thickness of device. The optimization parameters for Configuration B are given in Table 5.2. Then, keeping the number of degrees of freedom the same, we probe the effect of recovering the vertical patterning resolution at the cost of overall device thickness. Parameters for Configuration C are listed in Table 5.3.

| Property | Value |
|---|---|
| Device width, $w$ | 1.872µm (3.65$\lambda_0$) |
| Device thickness, $t$ | 1.008µm (1.97$\lambda_0$) |
| Focal length, $f$ | 1.404µm (2.74$\lambda_0$) |
| Minimum feature size | 36nm (0.070$\lambda_0$) |
| Layer thickness | 72nm (0.14$\lambda_0$) |

Table 5.2: Device Configuration B Properties

In each design space, depicted in Figure 5.1, we consider different index of refraction contrasts, where each device is binarized between $n = 1$ and $n = \{1.5, 2.5, 3, 3.5\}$. Devices are pushed to binary via a sigmoid filtering technique with increasingly large strength, in line with how devices have been binarized in previous chapters. Since index contrasts and device heights are different for and within each configuation, the number of iterations per binarization epoch is scaled by the overall index contrast and device thickness. This way, devices that have to binarize an overall larger volume of

| Property | Value |
|---|---|
| Device width, $w$ | 1.872µm ($3.65\lambda_0$) |
| Device thickness, $t$ | 0.504µm ($0.98\lambda_0$) |
| Focal length, $f$ | 1.404µm ($2.74\lambda_0$) |
| Minimum feature size | 36nm ($0.070\lambda_0$) |
| Layer thickness | 36nm ($0.070\lambda_0$) |

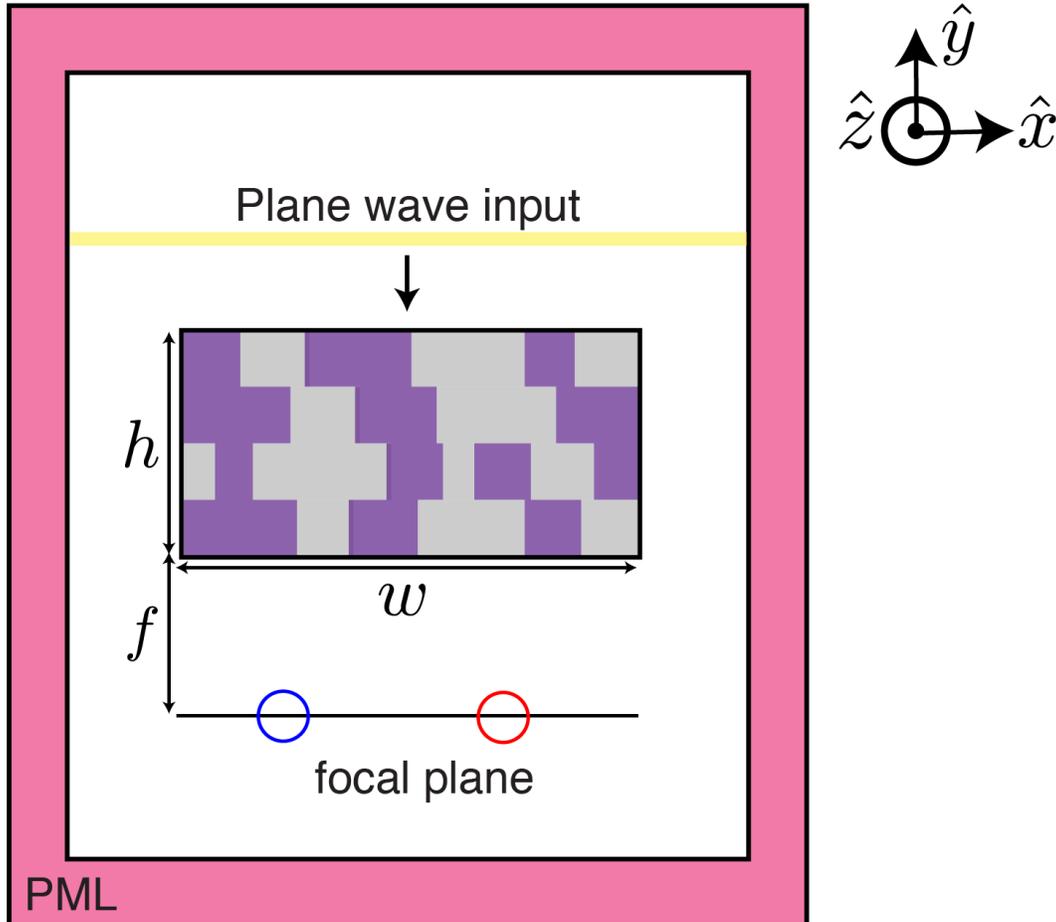Table 5.3: Device Configuration C Properties



Figure 5.1: 2D design space for probing optimization landscape. Intensity in the $z$-polarization under plane wave illumination is optimized for one of two spots depending on which spectral band the input wavelength belongs to.

permittivity a further amount have a longer optimization time to do so. For each final device, the total amount of achieved binarization, $B$, is indicated. Given a material density $\boldsymbol{\rho}$, $B = \frac{2}{N} \sum_{k=1}^{N} |\rho_k - 0.5|$. The optimization is setup to focus two bands of wavelengths, $\lambda \in [450\text{nm}, 500\text{nm}]$ and $\lambda \in [525\text{nm}, 575\text{nm}]$, to two discrete locations on the focal plane. Similar to other exploratory studies in this body of

work (see Chapters 2, 4), given the extensive simulation time in 3D optimizations, we optimize in the extruded 2D design space with a single polarization, $E_z$, and a normally incident plane wave illumination. In each band, 4 intensities from evenly spaced wavelengths are averaged together with the figure of merit formed as the product of the two individual band performances. The product helps re-balance the emphasis on each band performance depending on which is performing higher. For $g(\boldsymbol{\epsilon}) = g_{\text{low}}(\boldsymbol{\epsilon}) * g_{\text{high}}(\boldsymbol{\epsilon})$, the gradient is given by $\frac{\partial g}{\partial \boldsymbol{\epsilon}} = g_{\text{high}}(\boldsymbol{\epsilon})\frac{\partial g_{\text{low}}}{\partial \boldsymbol{\epsilon}} + g_{\text{low}}(\boldsymbol{\epsilon})\frac{\partial g_{\text{high}}}{\partial \boldsymbol{\epsilon}}$. The weighting for each individual gradient component component scales as the figure of merit of the other, so the higher performing figures of merit will be emphasized less [3]. The optimization problem is:

$$
\begin{aligned}
\max_{\boldsymbol{\epsilon} \in \{\epsilon_{\text{min}}, \epsilon_{\text{max}}\}^N} g(\mathbf{E}_z) &= g_{\text{low}}(\mathbf{E}_z) g_{\text{high}}(\mathbf{E}_z) \\
g_{\text{low}}(\mathbf{E}_z) &= \sum_{\lambda \in [450\text{nm}, 500\text{nm}]} \frac{I_z(\mathbf{r}_{\text{low}}, \lambda)}{I_{\text{max}}(\lambda)} \\
g_{\text{high}}(\mathbf{E}_z) &= \sum_{\lambda \in [525\text{nm}, 575\text{nm}]} \frac{I_z(\mathbf{r}_{\text{high}}, \lambda)}{I_{\text{max}}(\lambda)} \\
I_z(\mathbf{r}, \lambda) &= ||\mathbf{E}_z(\mathbf{r}, \lambda)||^2 \\
I_{\text{max}}(\lambda) &= \frac{\lambda_0^2}{\lambda^2} \\
\mathbf{r}_{\text{low}} &= (\frac{w_{\text{device}}}{4}, f) \\
\mathbf{r}_{\text{high}} &= (\frac{3w_{\text{device}}}{4}, f)
\end{aligned}
\tag{5.1}
$$

In Figure 5.2a,b, the optimization results are plotted for Configuration A and B, respectively. Cuts across the focal plane are shown for each wavelength and scaled by the maximum intensity value across all plots in each configuration for each wavelength. This is intended for qualitative comparison purposes. The figure of merit value is written for reference at the top of each plot. In Figure 5.2a, we note that an increase in index contrast is helpful moving from 1.5/1 to 2.5/1 and 3.0/1. The larger figure of merit comes from higher color contrast in the sorting. Beyond this, the result does not improve by simply increasing the index contrast and the 3.5/1 optimization result dips back down again.

Figure 5.3 shows the results, for each index contrast, of moving between Configurations B and C, from a thicker to a thinner device with the same overall number of design degrees of freedom. In the low index contrast case of 1.5/1 in Figure 5.3a, the reduction in thickness comes with a reduction in overall performance. However, for
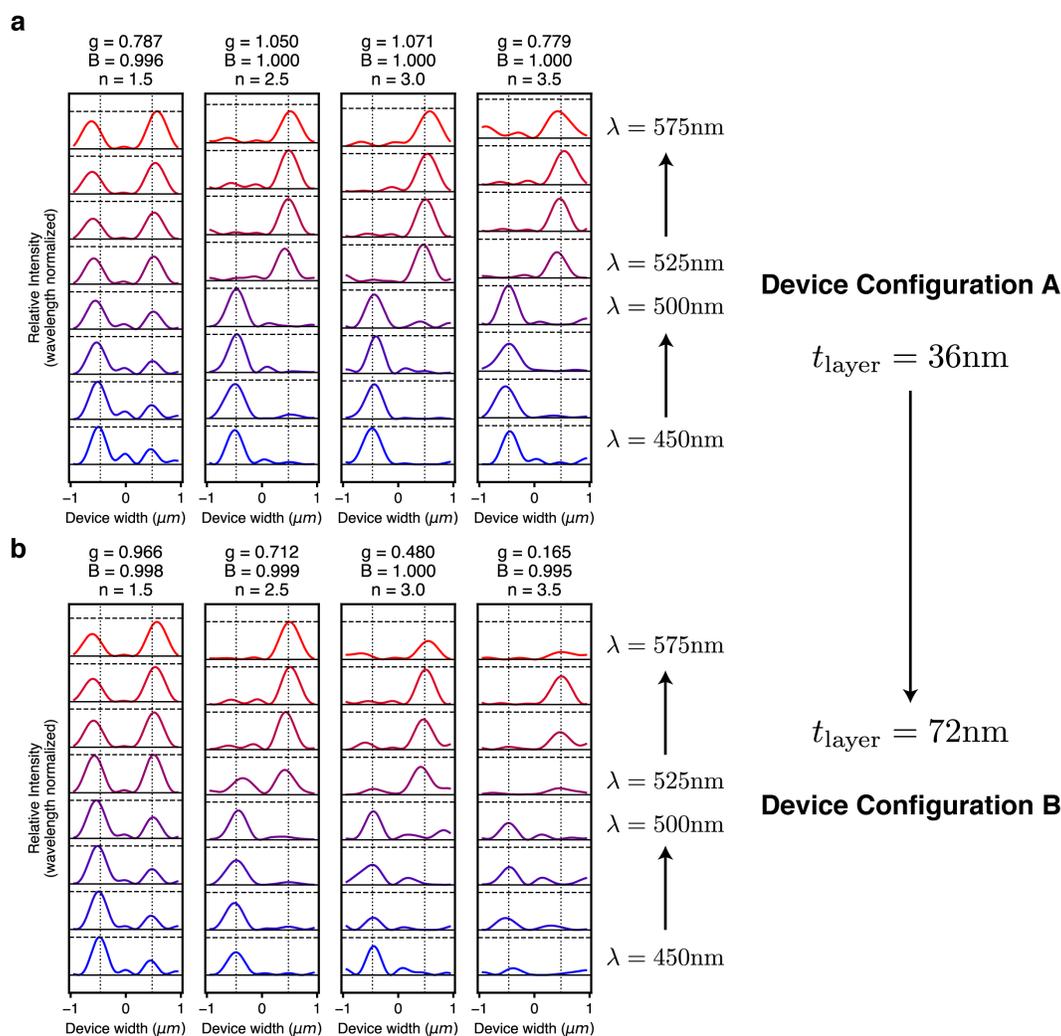
Figure 5.2: Optimization results for Configurations A and B. (a) Cut across focal plane intensity $\mathbf{I}_z$ for each wavelength in the optimization. The lower band of wavelengths is optimized to focus to the left focal point and the upper band to the right focal point. For each wavelength, the intensities are normalized by the maximum intensity for that wavelength in all four plots. At the top of the plot is the final figure of merit, $g$, the binarization, $B$, and the upper bound index of refraction in the optimization, $n$. The vertically dotted lines show the lateral location of the focal points and the horizontal solid and dashed lines indicate the bounds [0,1], respectively for each wavelength plot. (b) Same plots as (a) but for Configuration B.

the higher index contrast cases, the optimizer finds a better solution when moving to the thinner devices with more axial patterning resolution. To understand these observations better, in the next section we start by relating photonic inverse design problems to modern machine learning training and adopt visualization techniques from that field.
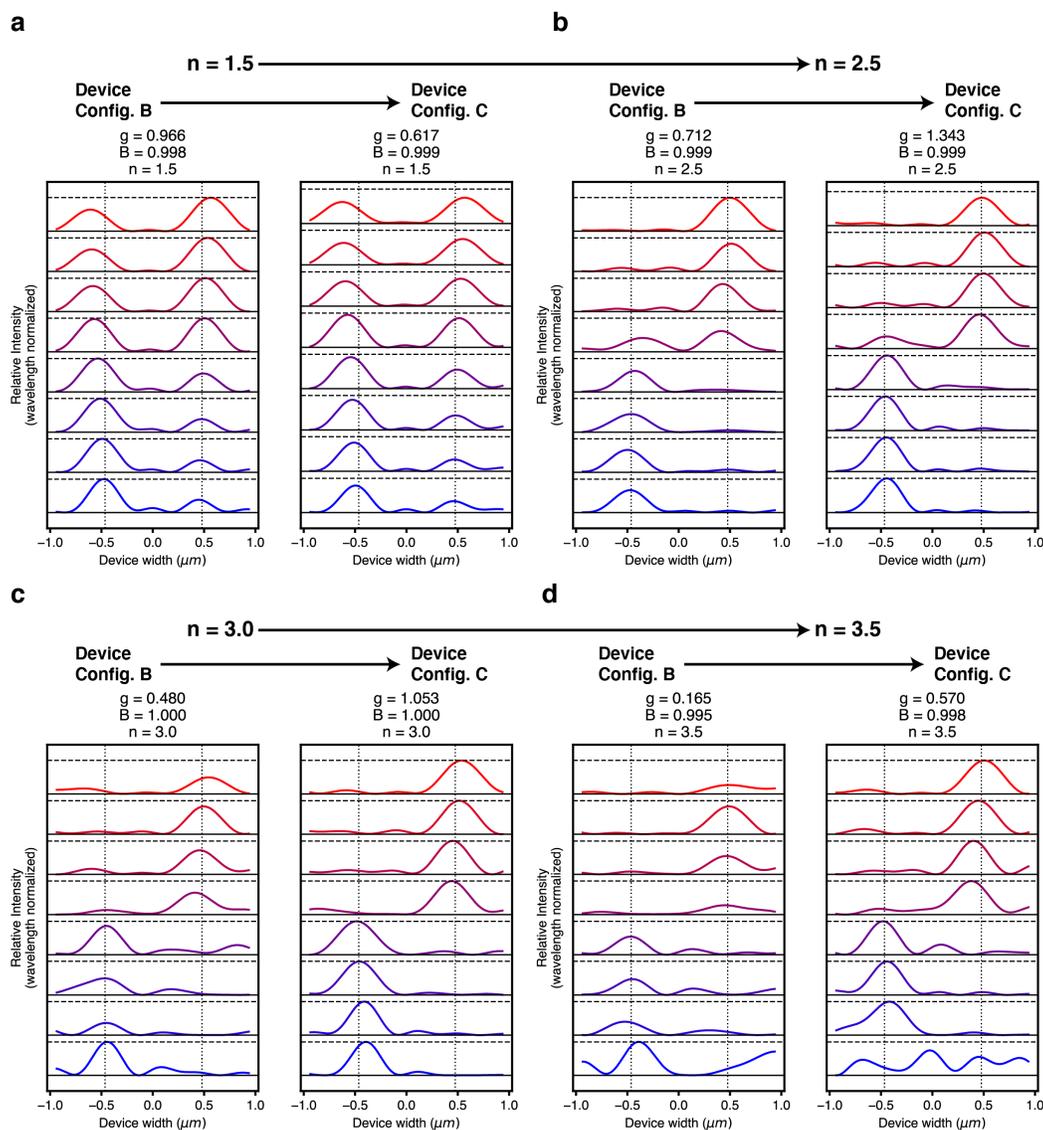
Figure 5.3: Comparison of Configurations B and C. (a) Cut across focal plane intensity $\mathbf{I}_z$ for each wavelength with the same structure as in Figure 5.2 where intensities are normalized by maximum intensity between both plots for each wavelength. These plots are for an upper index bound of $n = 1.5$. (b) Same plots as (a) but for $n = 2.5$. (C) Same plots as (a) but for $n = 3$. (d) Same plots as (a) but for $n = 3.5$.

## 5.3   Line Cuts

The dimensionality of a typical inverse design problem is large and optimizations search over a space well beyond typically low-dimensional visualization methods. This is the case in many practical optimization problems, especially deep learning which deals with an exceptional number of trainable parameters. For example, the

VGG network architecture for image classification problems that was published in 2014 contained over 100 million parameters [4]. With each dimension, the volume of the search space grows exponentially and thus we only ever search, compute, and visualize slices from this data [5]. One technique for building intuition about the landscape the optimizer is traversing is to look at projections along a small number of directions as a way to generate a low-dimensional view of its curvature. The simplest approach is to monitor the figure of merit along a line cut [6, 7]. Given two points in the design space, $\epsilon_A$ and $\epsilon_B$, we can plot the figure of merit, $g(\epsilon)$ on a line between the two. Mathematically, we probe:

$$\alpha \in [0, 1]$$
$$g((1 - \alpha)\epsilon_A + \alpha\epsilon_B)$$
$$(5.2)$$

For Configuration A, we see the performance initially improve with additional index contrast and then fall off for an index contrast of 3.5/1. In Figures 5.4, 5.5, 5.6, the line cuts are plotted when moving from the final permittivity for the 1.5/1, 2.5/1, and 3.0/1 optimizations to the final permittivity for the 3.5/1 optimization, respectively. For example, in Figure 5.4, according to Equation (5.2), $\epsilon_A$ is the final device for the 1.5/1 optimization and $\epsilon_B$ is the final device for the 3.5/1 optimization. Panel a of each figure shows the product figure of merit and panel b shows the component figures of merit for the lower and upper bands. Note in all plots the y-axis scaling is different as the cuts are meant to show landscape roughness even for lower overall figure of merit values. Since each line cut spans between different permittivity distributions, the x-axis for each plot moves over a different amount of total permittivity. This metric is indicated for each plot and computed as the amount of permittivity movement per step (for $N_\alpha$ points on the line) normalized by the areal fraction of the voxels with respect to a square with dimension equal to the midpoint wavelength, $\lambda_0$:

$$\frac{|\Delta\epsilon|}{\alpha}(\epsilon_A, \epsilon_B) = \frac{\sqrt{\sum_k (\epsilon_{A,k} - \epsilon_{B,k})^2}}{(N_\alpha - 1)} (\frac{\delta}{\lambda_0})^2 \qquad (5.3)$$

The general trend in each plot is that approaching the 3.5/1 optimization side of the plot, the landscape looks rougher with more local maxima along the way to the design point. This is especially noticeable in Configuration B where these

maxima are at nearly the same height as the final solution. We note again that these are not necessarily local maxima of the whole design space since the line cut is searching along just one of the many dimensions. This means the optimization will not necessarily settle in these places, especially since they are not binary solutions, but they still show noisy topology of the design landscape. There is not a one-to-one correspondence between this rough nature and the final optimization figure of merit. For example, Figure 5.4a shows a smooth approach towards the 1.5/1 optimization result but an overall lower figure of merit compared to those of the 2.5/1 and 3.0/1 Configuration A results. However, on the contrary, we may often find that rough optimization landscapes correlate to lower performing binary designs. If we consider the binarization process typically employed of slowly pushing a high performing continuous optimization towards a binarized one, movement through a highly varying landscape can break the assumption that a high quality local maximum to the binary problem is nearby the local maxima explored for the increasingly binarized grayscale problems. Once again, we highlight the difference here between the existence of a good solution in a design space and ability of the local optimizer to find such a solution. A summary of the line cuts, plotted with an x-axis adjusted for the total permittivity spacing between endpoints is shown in Figure 5.7.

In Figure 5.8, we look at another set of line cuts for the 2.5/1 and 3.0/1 index contrast optimizations in panels a and b, respectively. Here, instead of moving between devices with different index contrasts, we look at the effect of moving between Configurations A and B where the layer thickness doubles between the two with the overall thickness kept constant. We note that for lower index contrast 2.5/1, the maxima are broader along the direction of the line cut. This width is more easily comparable in the summary plot in Figure 5.9 where the x-axes have been scaled by the permittivity difference between the endpoints. For the thicker layers, the lower index contrast 2.5/1 optimization finds a smoother, wider, higher performing local maximum compared to the 3.0/1 optimization.
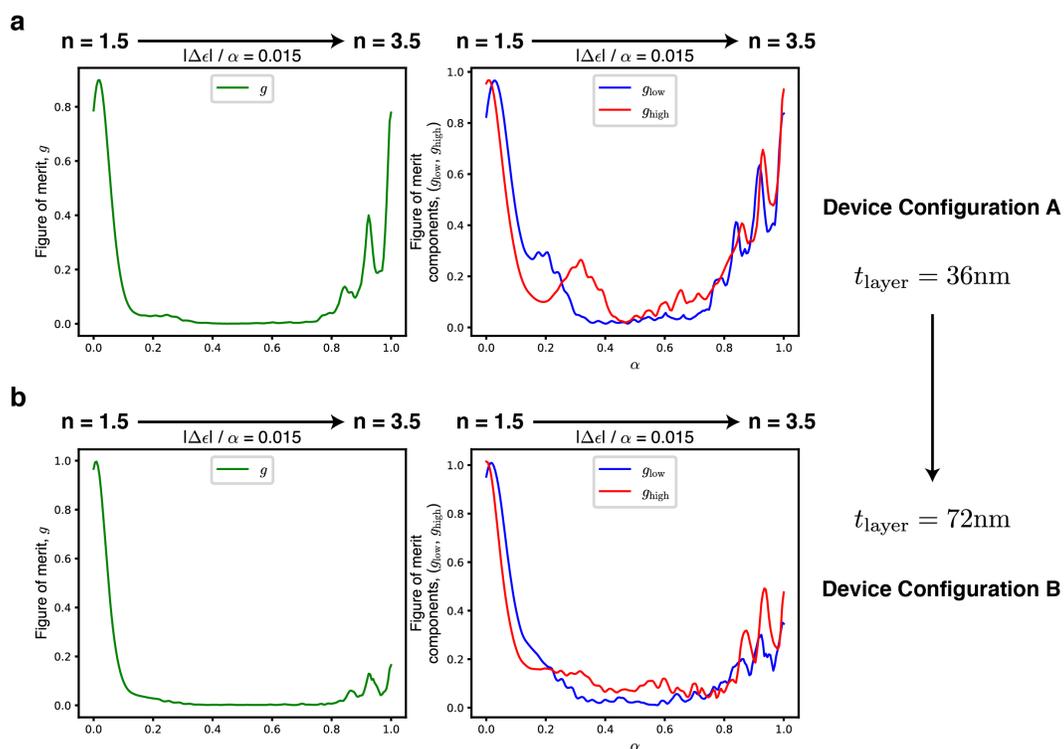
Figure 5.4: Line cuts between the $n = 1.5$ and $n = 3.5$ optimization results for Configurations A and B. (a) Configuration A line cuts for (left) full figure of merit, $g$, and (right) component figures of merit $g_{\mathrm{low}}$ and $g_{\mathrm{high}}$. (b) Same plots as (a) but for Configuration B.

Figure 5.5: Line cuts between the $n = 2.5$ and $n = 3.5$ optimization results for Configurations A and B. (a) Configuration A line cuts for (left) full figure of merit, $g$, and (right) component figures of merit $g_{low}$ and $g_{high}$. (b) Same plots as (a) but for Configuration B.
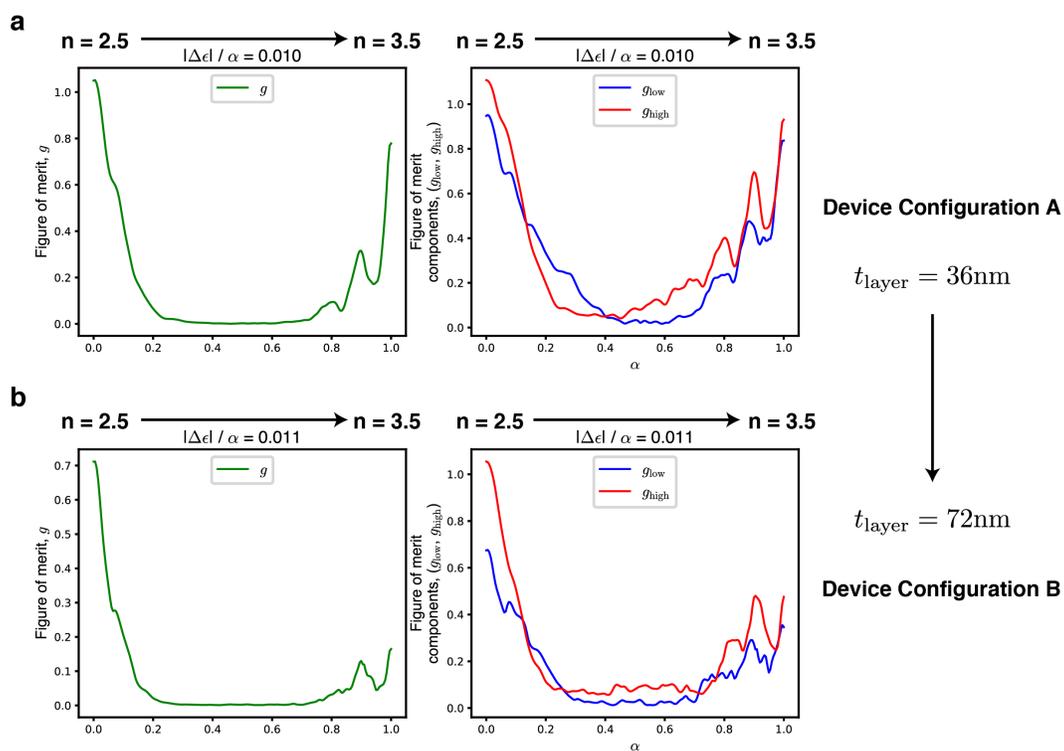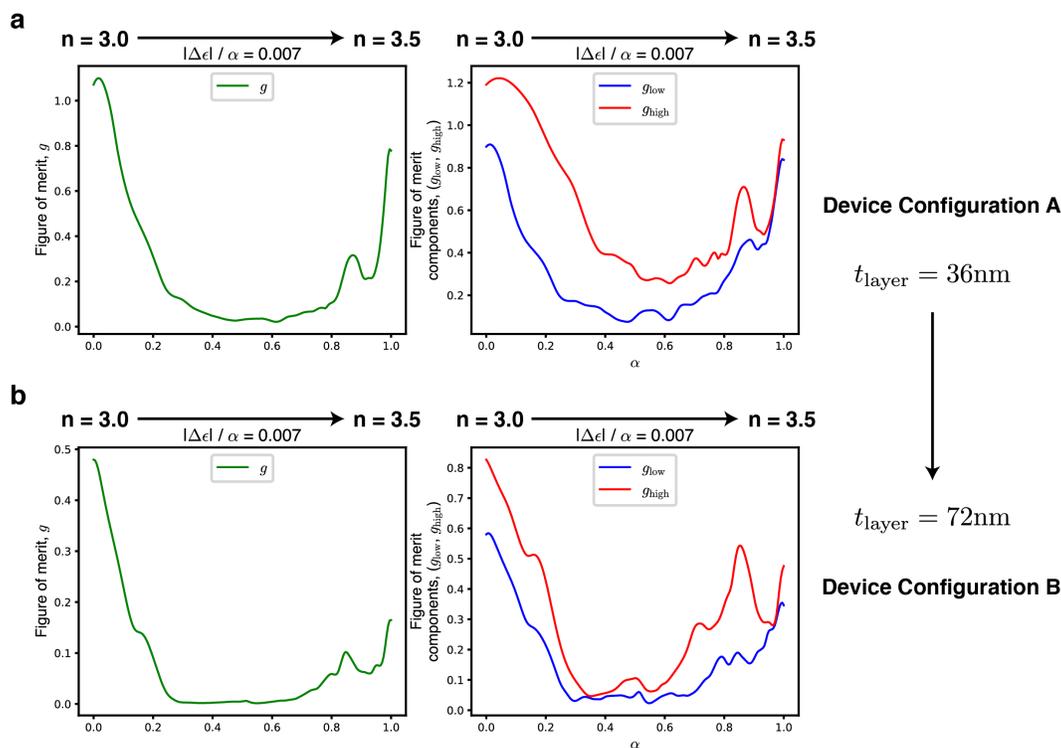
Figure 5.6: Line cuts between the $n = 3.0$ and $n = 3.5$ optimization results for Configurations A and B. (a) Configuration A line cuts for (left) full figure of merit, $g$, and (right) component figures of merit $g_{\text{low}}$ and $g_{\text{high}}$. (b) Same plots as (a) but for Configuration B.
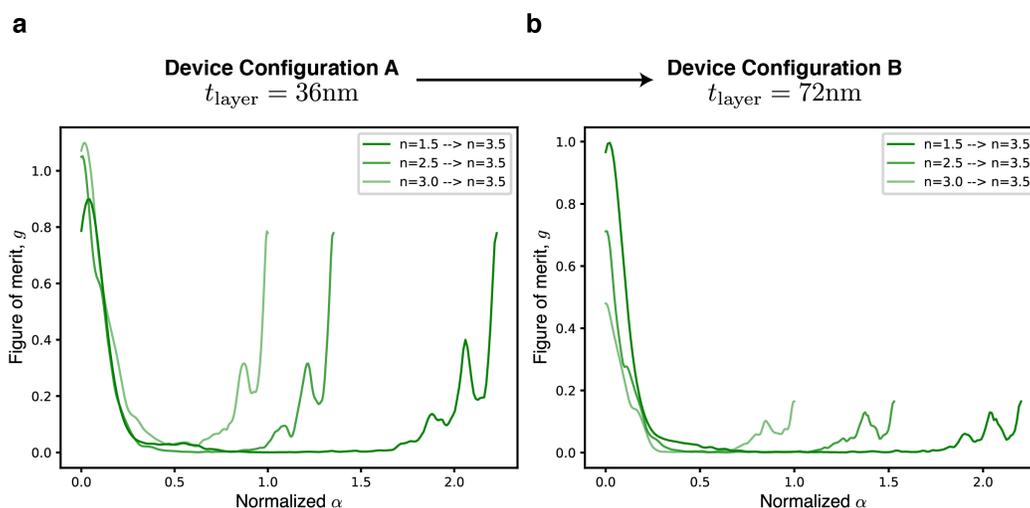


Figure 5.7: Line cuts for all index contrasts with comparable x-axis for Configurations A and B. (a) Configuration A line cuts for full figure of merit, $g$ with x-axis proportional to $\frac{|\Delta\boldsymbol{\epsilon}|}{\alpha}$. (b) Same plots as (a) but for Configuration B.
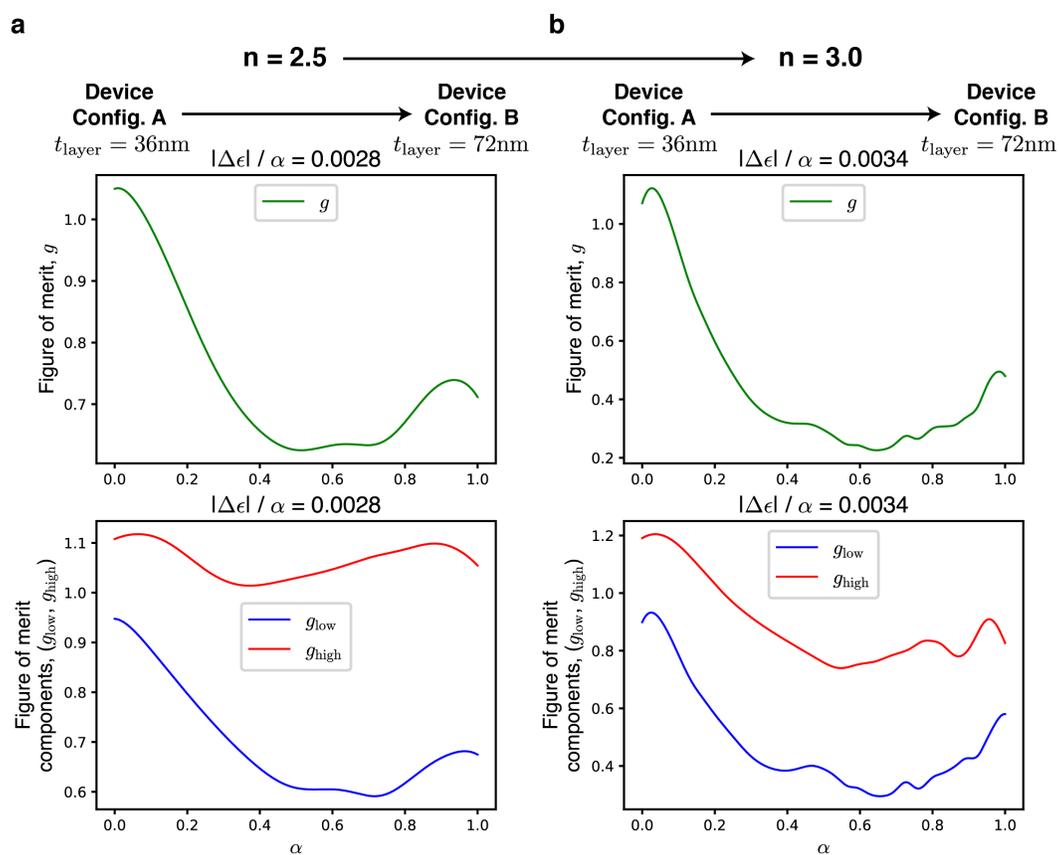
Figure 5.8: Line cuts between Configurations A and B for $n = 2.5$ and $n = 3.0$. (a) Configuration A to B line cuts for (left) full figure of merit, $g$, and (right) component figures of merit $g_{\text{low}}$ and $g_{\text{high}}$ for index of refraction contrast 2.5/1. (b) Same plots as (a) but for index of refraction contrast 3.0/1.

**Device Configuration A**
$t_{\text{layer}} = 36\text{nm}$

**Device Configuration B**
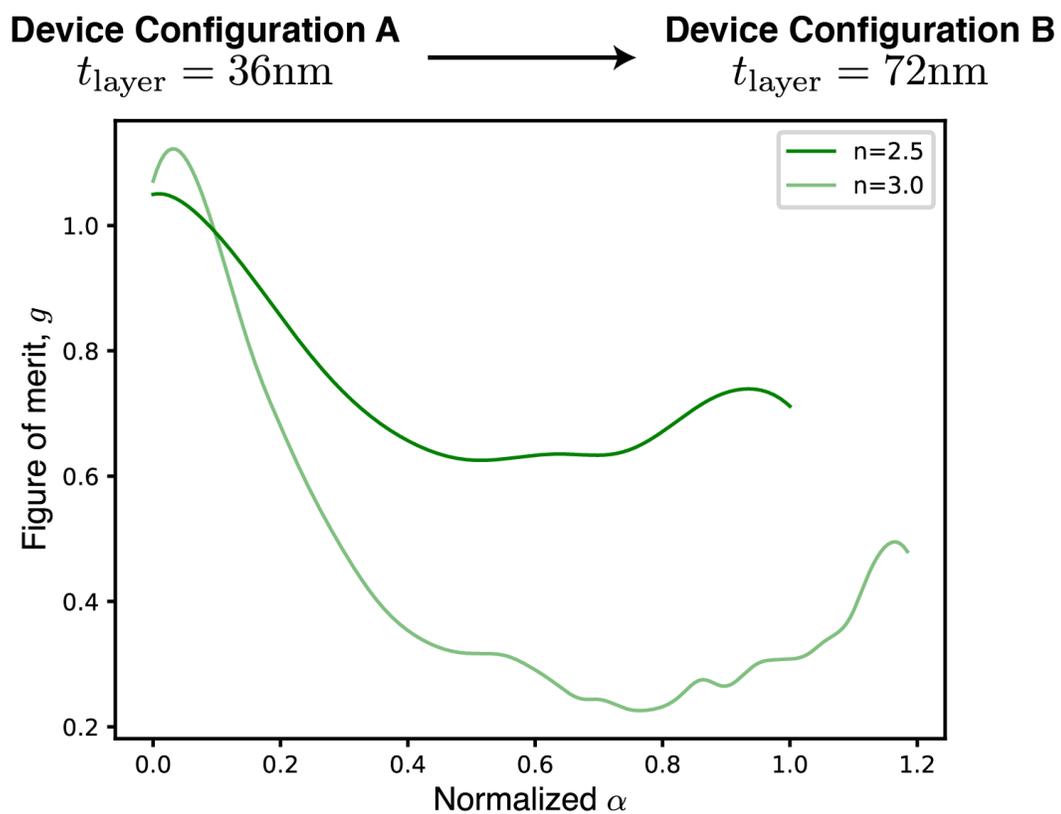$t_{\text{layer}} = 72\text{nm}$



Figure 5.9: Line cuts for each configuration sweep with x-axes proportional to $\frac{|\Delta\epsilon|}{\alpha}$ for $n = 2.5$ and $n = 3.0$. (a) Configuration A to B line cuts for full figure of merit, $g$.

## 5.4  2D Sweeps

Line cuts visualize approaches to optima along a single direction. While moving to probing the full dimensionality of the design space around the optima is outside the grasp of visualization and compute budget capabilities, at least one dimension can be added to the current technique. Similar to work done in this area for analyzing machine learning models [6], the landscape can be probed as described in Equation (5.4).

$$\beta, \gamma \in [0, 1]$$

$$\mathbf{p}_1 \cdot \mathbf{p}_2 = \sum_k p_{1,k} p_{2,k} = 0 \tag{5.4}$$

$$g(\boldsymbol{\epsilon}^* + \beta \mathbf{p}_1 + \gamma \mathbf{p}_2)$$

In Equation (5.4), $\boldsymbol{\epsilon}^*$ is the result of an optimization. Perturbations $\mathbf{p}_1$, $\mathbf{p}_2$ are chosen such that the evaluated permittivity is kept within the optimization bounds. The permittivity, in general, will be grayscale as it is perturbed around the binary optimum, but it will always remain between the lower and upper permittivity bounds for the optimization. For a given design with binary profile, $\boldsymbol{\rho}$, we fill in two perturbation vectors $\mathbf{p}_1$ and $\mathbf{p}_2$. At every point, based on a fair coin flip, we assign a random value to one of the two perturbation vectors and $0$ to the other. The nonzero value is randomly sampled from a uniform distribution between $[0, 1]$. If $\boldsymbol{\rho}$ at that point is $1$, then the perturbation value is multiplied by $-1$ such that it searches values below the upper permittivity bound. Otherwise, it is left positive for the opposite effect. Afterwards, $\mathbf{p}_{1,2}$ are rescaled to have the same length. For each optimization, the number of design voxels, maximum permittivity bound, and overall device thickness may differ. For comparison purposes, $\mathbf{p}_{1,2}$ are scaled such that the maximum perturbation size is the same (i.e., $|\mathbf{p}_{1,2}|$ is kept constant across plots). The total perturbation normalized to the voxel size in units of wavelength in every plot below is set to be $|\Delta \mathbf{p}_k| = (\frac{\delta}{\lambda_0})^2 \sqrt{\sum_k p_k^2} = 0.74$.

Figures 5.10-5.11 each show, for a fixed index contrast (3.0/1 and 3.5/1), the comparison between the 3 optimization configurations. Given the construction of the sweep, the device output from the optimization occurs in the bottom left corner for $\beta = \gamma = 0$. Between Configuration A and B, where the layer thickness increases by a factor of 2, we see a drop in the optimized device figure of merit along with a narrowing of the local maxima along these two axes. Qualitatively, the landscape also appears more irregular for the thicker layers and even contains

an additional nearby maxima for index contrast 3.5/1 as seen in Figure 5.11b. For higher index contrasts, Figure 5.3c,d showed device performance improved when moving between Configurations B, C. These configurations have the same number of degrees of freedom with respect to design voxels, but Configuration C has half the total thickness. This represents a trade-off between overall device thickness and axial patterning resolution. In this case, the optimization of the thinner device with better resolution finds a better maxima. We see in Figure 5.10 the recovery of the performance lost between Configurations A and B by optimizing with the same layer thickness as in Configuration A but reducing the number of layers by a factor of 2. This recovery in performance comes with the area around the local optimum becoming smoother and wider again. A similar effect with less, but still significant, recovery of performance happens in Figure 5.11.
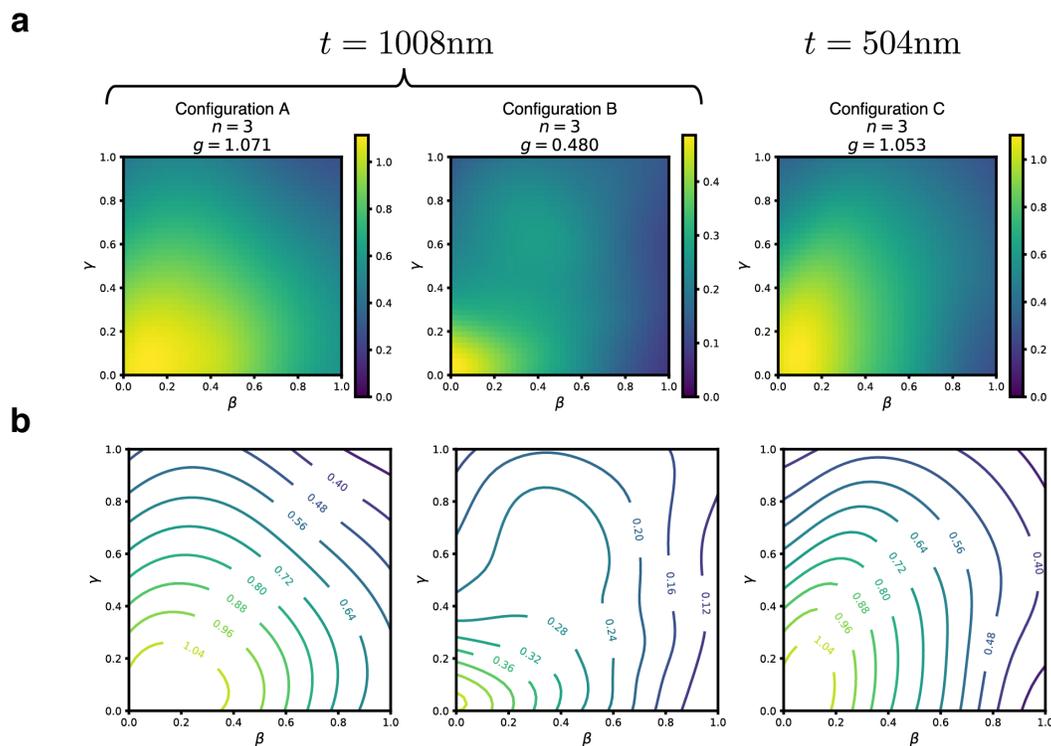


Figure 5.10: 2D sweep around optimization solution for Configurations A-C for index contrast 3.0/1. (a) Landscape around randomly chosen perturbation vectors $\mathbf{p}_{1,2}$ for each configuration. The colorbars are scaled between the maximum value for each plot and 0 to highlight relative changes in figure of merit. The optimized device figure of merit, $g$, are given for each device. (b) Contour lines for the respective plots in (a).

A comparison of Configuration B devices is shown in a single plot in Figure 5.12.

Figure 5.11: 2D sweep around optimization solution for Configurations A-C for index contrast 3.5/1. (a) Landscape around randomly chosen perturbation vectors $\mathbf{p}_{1,2}$ for each configuration. The colorbars are scaled between the maximum value for each plot and 0 to highlight relative changes in figure of merit. The optimized device figure of merit, $g$, are given for each device. (b) Contour lines for the respective plots in (a).

Here, as the index contrast increases, the figure of merit steadily drops and the landscape becomes increasingly irregular nearby the local optima.
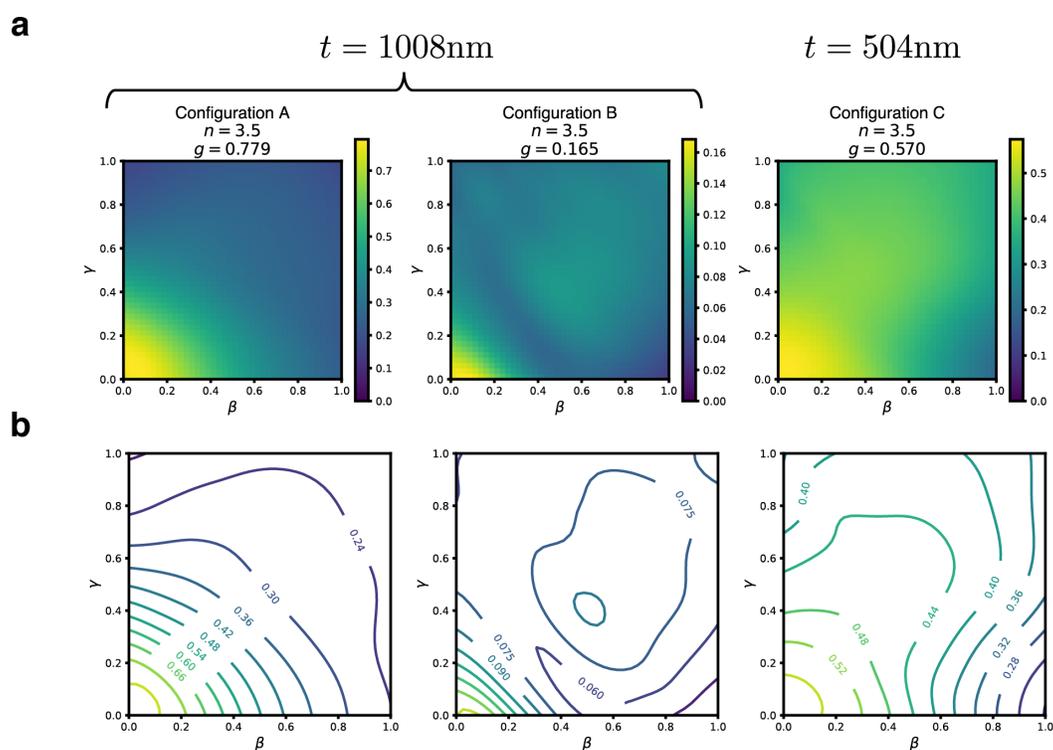
Figure 5.12: 2D sweep around optimization solutions for Configuration B and index contrast 2.5/1, 3.0/1, 3.5/1. (a) Landscape around randomly chosen perturbation vectors $\mathbf{p}_{1,2}$ for each configuration. The colorbars are scaled between the maximum value for each plot and 0 to highlight relative changes in figure of merit. The optimized device figure of merit, $g$, are given for each device. (b) Contour lines for the respective plots in (a).

## 5.5   Conclusion

While these analyses are qualitative and based on visualization, high-dimensional non-convex optimization spaces come with few guarantees about their shape and properties making them difficult to analyze exhaustively. These types of visualizations should be carried out cautiously and similar to interpreting the single optimization results, their outcomes should not be taken as firm rules. As is present even in this limited study, the various hyperparameters controlling optimization results are interwoven and cannot be treated as independent entities. In practice, however, there are design spaces that seem easier to optimize within than others even when they physically contain similar amounts of degrees of freedom. Considering many of the devices we are interested in are fully 3D and thus costly to optimize, setting up local optimizers for success is critical.

The visualizations in this study were based on the result of a simple, gradient descent optimizer. A different optimizer, possibly one that includes a momentum contribution, or stochastic element that can jump the design into the realm of wider local optima may improve the quality of solutions in the tougher design spaces [8]. A formidable challenge in photonic inverse design is the restriction to a binary space of solutions. Using continuous optimization techniques that slowly push towards binary solutions via filtering has had a lot of success in ultimately converging to high performing binary designs [9]. However, this becomes increasingly challenging in large index contrast design spaces when feature resolution is not unlimited [10]. Via regularization methods or explicit constraints, geometries of binary optimized solutions could be restricted in these types of optimizations if it can be reasoned what properties tend to be associated with smooth, wide optima. Visualization methods presented in this chapter can help profile how landscapes change when restricted to certain geometries without needing to run full optimizations.

In some cases, we may have a choice to move to a more promising design space and these visualization methods can help justify an otherwise less intuitive geometric change. For example, we saw that when restricted to a specific number of device layers, with index contrasts of 3.5/1 or 3.0/1, our optimization space was easier to navigate for half of the overall device thickness (see Figures 5.10-5.11). Reducing thickness or limiting degrees of freedom to smooth design spaces will not always be a good solution. For optimization targets that are increasingly complex, larger degrees of freedom will be required to accomplish them. This may come in the form of increased design volume or index contrast where simply reducing these factors

will not achieve sufficient performance. Modern deep learning has been on a trend of increased degrees of freedom for many years, where the number of weights in models are getting larger to accommodate the accomplishment of complex tasks. Models do not, however, simply increase in training and generalization performance by the addition of more network layers and parameters. The architecture of a network needs to be carefully considered for the performance to increase with network complexity. A good example of this is the residual net architecture, or ResNet, which uses shortcut connections with blocks of convolutions to allow learning of a residual mapping as opposed to the original underlying mapping [11]. Using these residual blocks, ResNet architectures can continue to gain accuracy while adding more layers to the networks. In photonic inverse design, geometries that steadily add performance with additional degrees of freedom are critical to continuing to miniatiurize and improve efficiency of metaoptics.

## References

[1] H.-T. Chen, A. J. Taylor, and N. Yu, "A review of metasurfaces: physics and applications", Reports on progress in physics **79**, 076401 (2016).

[2] P. Camayd-Muñoz, C. Ballew, G. Roberts, and A. Faraon, "Multifunctional volumetric meta-optics for color and polarization image sensors", Optica **7**, 280–283 (2020).

[3] M. Mansouree, H. Kwon, E. Arbabi, A. McClung, A. Faraon, and A. Arbabi, "Multifunctional 2.5 D metastructures enabled by adjoint optimization", Optica **7**, 77–84 (2020).

[4] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition", arXiv preprint arXiv:1409.1556 (2014).

[5] R. Bellman, "Dynamic programming", Science **153**, 34–37 (1966).

[6] H. Li, Z. Xu, G. Taylor, C. Studer, and T. Goldstein, "Visualizing the loss landscape of neural nets", Advances in neural information processing systems **31** (2018).

[7] I. J. Goodfellow, O. Vinyals, and A. M. Saxe, "Qualitatively characterizing neural network optimization problems", arXiv preprint arXiv:1412.6544 (2014).

[8] S. Ruder, "An overview of gradient descent optimization algorithms", arXiv preprint arXiv:1609.04747 (2016).

[9] F. Wang, B. S. Lazarov, and O. Sigmund, "On projection methods, convergence and robust formulations in topology optimization", Structural and multidisciplinary optimization **43**, 767–784 (2011).

[10]C. Ballew, G. Roberts, T. Zheng, and A. Faraon, "Constraining continuous topology optimizations to discrete solutions for photonic applications", ACS photonics **10**, 836–844 (2023).

[11]K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition", in Proceedings of the ieee conference on computer vision and pattern recognition (2016), pp. 770–778.

*C h a p t e r   6*

# CONCLUSION

Volumetric metaoptics is in its infancy, ripe with interesting research questions to explore as well as an ability to make an instant impact on optical technologies. In the realm of research, there are still open questions on how much device complexity is needed for different types of functionalities and how to get the most out of a certain number of degrees of freedom. While modern nanofabrication technology may be capable of creating intricate structures, there are cost and robustness considerations that come with using increasingly complex manufacturing processes. This question of required complexity can be addressed by squeezing performance out of as few degrees of freedom as possible and also approached through rethinking the architectures that host inverse designed devices. For example, considering both the device and reconstruction algorithm simultaneously, entire optical systems can be optimized holistically [1]. Further, expanding the optimizable size of devices by using more efficient solvers can open up new domains for inverse design [2]. In Chapter 4, we observed that with a coarse feature size constraint, it was difficult to find high performing solutions even with several patterned layers. Understanding the nature of and how to work in these types of design spaces is important for easing nanofabrication burden especially in academic settings. Coupling local optimizations with efficient global and machine learning techniques may be part of the answer to this challenge [3]. Finally, optimizing for device tunability and changes in functionality inside dynamic environments can open up new paths for compact sensing and active optical systems. Devices can be made to change functionality with modulation of a nearby refractive index using phase change material or liquid crystals, for example [4].

From a practical standpoint, volumetric metaoptics have tangible near-term promise. It is interesting already from the perspective of improving color imaging efficiency by completely or partially replacing absorptive filters with color sorting devices. Spectral routing structures have been a research topic for several years and some are being actively pursued in industry [5]. In Chapter 3, we showed that the color routing problem could be thought of more generally and devices occupying the same footprint could sort on combinations of spectral, polarization, and spatial properties of incoming light. This puts within practical grasp the ability to change

the function of a camera by just modifying the properties of scattering structures built into the sensor array. By working with experts in various imaging and sensing domains, we can propose new optical technologies that may have been considered too challenging in the past. This could have been for a variety of reasons including required miniaturization or form factor, or a need to integrate a device in a complex environment. For example, there is increasing interest in using metaoptics for fiber-based imaging systems where the system geometry and fiber transmission function create new challenges [6, 7].

The devices proposed in this work require dedicated design and fabrication effort with the promise that they can achieve unprecedented levels of performance in compact form factors. An analogy that often comes to mind is that of hardware accelerators in computing. General computing hardware trades off flexibility for power efficiency and speed. In some processors there are known and oft-used bottleneck functionalities like, for example, an image signal processing pipeline connected to one or many camera sources. This functionality is sometimes offloaded into a specialized hardware accelerator. The hardware accelerator is a dedicated set of transistors and memory access specifically designed to perform a task at high speed and low power. If we map this to optics, the standard bulk optical components are our general purpose computing; they offer high performance for a variety of tasks via reconfiguration of their relative positions. However, if we have a task we know needs to be done in an extremely small volume or at high efficiency, we can think of the inverse designed devices as accelerators, able to be designed specifically for those tasks. While it requires an investment, the accelerator can be an enabling technology in resource constrained settings. It also can go beyond what the general purpose compute can do and its function is limited only by the imagination and creativity of the designer. As this field of volumetric metaoptics grows and scientists continue to find ways to prototype this novel class of devices experimentally, we will see new types of optics and the ability to explore photonics in fully 3D settings. We will continue to learn and understand the limits of optical components and push the boundaries of sensing and imaging.

## References

[1] Z. Lin, C. Roques-Carmes, R. Pestourie, M. Soljačić, A. Majumdar, and S. G. Johnson, "End-to-end nanophotonic inverse design for imaging and polarimetry", Nanophotonics **10**, 1177–1187 (2021).

[2]J. Skarda, R. Trivedi, L. Su, D. Ahmad-Stein, H. Kwon, S. Han, S. Fan, and J. Vučković, "Low-overhead distribution strategy for simulation and optimization of large-area metasurfaces", npj Computational Materials **8**, 78 (2022).

[3]J. Jiang and J. A. Fan, "Global optimization of dielectric metasurfaces using a physics-driven neural network", Nano letters **19**, 5366–5372 (2019).

[4]H. Chung and O. D. Miller, "Tunable metasurface inverse design for 80% switching efficiencies and 144 angular deflection", ACS Photonics **7**, 2236–2243 (2020).

[5]S. Yun, S. Roh, S. Lee, H. Park, M. Lim, S. Ahn, and H. Choo, "Highly efficient color separation and focusing in the sub-micron CMOS image sensor", in 2021 ieee international electron devices meeting (iedm) (2021), pp. 30–31.

[6]H. Ren, J. Jang, C. Li, A. Aigner, M. Plidschun, J. Kim, J. Rho, M. A. Schmidt, and S. A. Maier, "An achromatic metafiber for focusing and imaging across the entire telecommunication range", nature communications **13**, 4183 (2022).

[7]J. E. Fröch, L. Huang, Q. A. A. Tanguy, S. Colburn, A. Zhan, A. Ravagli, E. J. Seibel, K. F. Böhringer, and A. Majumdar, "Real time full-color imaging in a meta-optical fiber endoscope", eLight **3**, 1–8 (2023).