

Programmable Spectrometry — Per-pixel Classification of Materials using Learned Spectral Filters

Vishwanath Saragadam, and Aswin C. Sankaranarayanan
Department of ECE, Carnegie Mellon University, USA

vishwanathsrv@cmu.edu

Abstract

Many materials have distinct spectral profiles. This facilitates estimation of the material composition of a scene at each pixel by first acquiring its hyperspectral image, and subsequently filtering it using a bank of spectral profiles. This process is inherently wasteful since only a set of linear projections of the acquired measurements contribute to the classification task. We propose a novel programmable camera that is capable of producing images of a scene with an arbitrary spectral filter. We use this camera to optically implement the spectral filtering of the scene’s hyperspectral image with the bank of spectral profiles needed to perform per-pixel material classification. This provides gains both in terms of acquisition speed — since only the relevant measurements are acquired — and in signal-to-noise ratio — since we invariably avoid narrowband filters that are light inefficient. Given training data, we use a range of classical and modern techniques including SVMs and neural networks to identify the bank of spectral profiles that facilitate material classification. We verify the method in simulations on standard datasets as well as real data using a lab prototype of the camera.

1. Introduction

Material composition of a scene can often be identified by analyzing variations of light intensity as a function of spectrum or wavelengths. Since materials tend to have unique spectral profiles, spectrum-based material classification has found widespread use in numerous scientific disciplines including molecular identification using Raman spectroscopy [5], tagging of key cellular components in fluorescence microscopy [17], land coverage and weather monitoring [4, 11], and even the study of chemical composition of stars and astronomical objects using line spectroscopy. It would not be a stretch to suggest that spectroscopy or its imaging variant, hyperspectral imaging (HSI), is an important scientific tool for material identification.

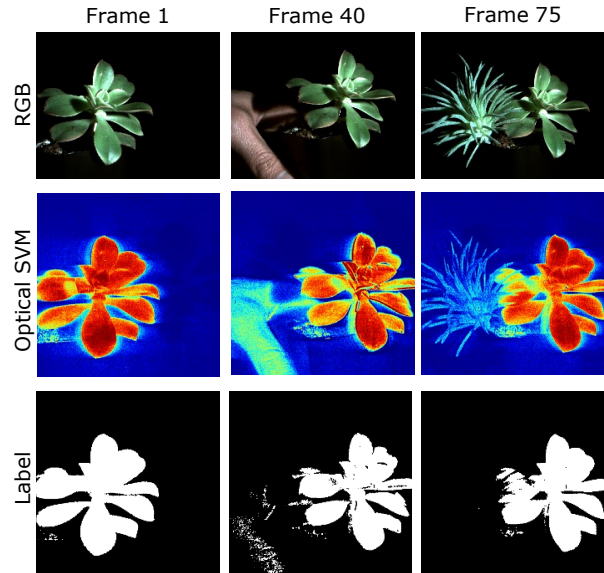


Figure 1: **Video-rate material classifier.** We propose an optical setup that is capable of classifying material on a per-pixel basis. This is achieved by building a programmable spectral filter that can image at high spatial resolution. The images here show a video sequence of a identifying real plants from plastic plants. We captured the data at 4fps, and performed only a per-pixel thresholding to get the video result.

While hyperspectral imaging has also found application in computer vision tasks [14, 25, 31], its widespread adoption has been hindered due to inherent challenges in acquisition them. Capturing a HSI requires sampling of a very high dimensional signal; for example, mega-pixel images at hundreds of spectral bands, a process that is daunting to do at video rate. This problem is further aggravated by the fact that hyperspectral measurements have to combat low signal to noise ratios, as a fixed amount of light is divided in to several spectral bands — leading to long exposure times that can even span several minutes per HSI.

This paper proposes a novel approach for enabling spectrometry-based per-pixel material classification by

overcoming the limitations posed by HSI acquisition. To understand our proposed approach, we first need to delve deeper into the process of classification itself. Classification techniques involve comparing the spectral profile at each pixel with known or *learned* spectra by taking a linear projection. Intuitively, given K material classes, we would compute $\mathcal{O}(K)$ such linear projections. For example, a support vector machine (SVM) classifies by finding distance of features from the separating hyperplane; in the context of spectral classification, this translates to spectrally filtering the scene with the hyperplane coefficients. Hence, spectral classification can be made practical if we can capture the linear projections directly without having to acquire the complete HSI. Such an operation translates to optically filtering the scene’s HSI using known spectral filters, which can be achieved if the camera’s spectral response can be arbitrarily programmed.

To enable per-pixel material classification, we propose a new imaging architecture with a programmable spectral response that can be changed on-the-fly at video rate. Given a training dataset of spectral profiles, we use off-the-shelf classification techniques like SVMs and deep neural networks to identify linear projections that facilitate material classification. For a novel scene, the camera captures multiple images, each with a different spectral response; the captured measurements are used with the classifier to perform per-pixel material classification.

The proposed pipeline has numerous benefits. *Optical computing* of the linear projections allows us to circumvent the measurement of sampling the full HSI. This has the dual benefit of reducing the acquisition time (from minutes to hundreds of milliseconds) as well as increasing light efficiency of each captured image since the linear projections often correspond to broadband spectral profiles. For binary classification problem, our lab prototype provides a classification result every second frame thereby providing material labels at 4 frames per second. We also show results on multi-class labeling problems using a classifier that can differentiate between five distinct material types.

2. Prior Work

We discuss prior work in the areas of material classification using HSIs as well as optical computing and design of programmable spectral filters.

Hyperspectral classification. Consider the HSI of a scene, $H(x, y, \lambda)$, where each pixel (x, y) is assumed to belong to one of K material classes. Specifically, the spectra at each pixel can be written as,

$$H(x, y, \lambda) = \alpha(x, y)S_{L(x, y)}(\lambda), \quad (1)$$

where $L(x, y)$ is label of the material contributing to spectrum at (x, y) , and $\alpha(x, y)$ is scaling parameter. Note that

the model above assumes all spatial pixels are pure, i.e., every pixel gets contribution from only one material. We use this model for the sake of exposition and later discuss about how to relax it later to handle mixed pixel.

The goal of classification is to estimate the label at each pixel, $L(x, y)$, which forms a label map. There are broadly two approaches to spectral classification — generative and discriminative. Generative techniques rely on decomposing the HSI as a linear combination of basic materials that are called end-members [6]. Specifically, the HSI of the scene is decomposed as,

$$H(x, y, \lambda) = \sum_{k=1}^K s_k(\lambda)a_k(x, y), \quad (2)$$

where $s_k(\lambda)$ is the spectra of k^{th} material, and $a(x, y)$ is the relative contribution of material k at (x, y) . The abundances at each pixel along with the end-member spectra provide a feature vector that can be used to spatially cluster the materials and subsequently identify them.

Discriminative techniques rely on directly learning discerning features from the HSI without the intermittent stage of low-dimensional decomposition. Here, we identify a set of spectral filters, $\{(d_k(\lambda), \beta_k)\}_{k=1}^M$ that generate per-pixel feature vector via spectral-domain filtering:

$$F_k(x, y) = \int_{\lambda} H(x, y, \lambda)d_k(\lambda)d\lambda + \beta_k. \quad (3)$$

Hence, each image $F_k(x, y)$ is a spectrally-filtered version of the HSI with an added offset. In case of SVMs, the learned spectral filters form separating hyperplanes; this has been a *de facto* way of HSI classification [7, 22]. More sophisticated learning techniques based on neural networks use spectral features [13] or spatio-spectral features [30, 18, 10, 15, 3, 16, 21, 12] for classification. Invariably, the number of spectral features used, i.e., the dimensionality of the projection, tends to be smaller than the number of spectral channels in the HSI. Hence, we seek to measure the features directly, by computing (3) optically. As is to be expected, such a paradigm of *optical classification* requires the design of cameras that can be programmed with arbitrary spectral filters.

Optical computing. Instead of relying on both spatial and spectral information, we consider a simpler approach which relies only on the spectral profiles for classification. Such a strategy is less accurate than spatial and spectral versions [30, 18, 10, 15, 3, 16, 21, 12], but significantly reduces the complexity of the imaging system. This approach is similar, in spirit, to using BRDFs to perform per-pixel classification by varying the incident illumination [19, 9], or using first layer of a neural network to capture light fields [2]. Such a setup offers two-fold advantage:

1. *Fewer measurements.* Since the number of material classes is far fewer than number of spectral bands, we need to measure far fewer measurements. For example, we show in our experiments that 3 – 5 images suffice for a 5-class classification task.
2. *Increased SNR.* The discriminating filters tend to be spectrally broadband, and hence each image is measured at higher light levels than any individual narrow spectral band. Hence, the images can be captured at higher SNR or at faster acquisition rates.

Optical computing has found use in various computer vision tasks such as capturing light transport matrices [24], low-rank approximation of hyperspectral images [29], and spectral classification using programmable light sources [8, 26]. We adopt the paradigm of optical computing to make discriminative filter measurements by building a camera whose spectral response can be arbitrarily programmed.

Dynamic spectral filters. Spectral filtering can be achieved by modified the response of the camera; a canonical and static example being the Bayer pattern or more interestingly, the case of fluorescence filters in microscopy. It is however more useful to have a camera whose response can be altered arbitrarily in a fast manner. Numerous techniques to achieve spectral filtering have been proposed in the past. Agile spectral imager [23] rely on the coding the so-called “rainbow plane” to achieve arbitrary spectral filtering. This was further developed by [20] where they placed a digital micromirror device (DMD) on the rainbow plane to achieve dynamic spectral filtering.

However, such architectures come with a debilitating problem — usage of simple pupil codes such as open aperture or a slit directly tradeoff spatial resolution for spectral resolution. This was first identified in [29] in the context of hyperspectral imaging. They showed that a slit, a common choice for spectrometry, leads to large spatial blur. Similarly an open aperture, a common choice for high-resolution imaging, leads to large spectral blur. Hence, such apertures are not capable of spectral classification with high accuracy.

We instead rely on the optical setup in [29] to overcome the spatial-spectral tradeoff. The key idea is to use a coded aperture that introduces an invertible blur in both spatial and spectral domains. An important difference is that the setup in [29] is designed for HSI image acquisition; this paper adapts the underlying ideas for performing material classification in the scene.

3. Programmable Spectral Filter

Our optical setup is a modification of the optical setup proposed in [29]. We briefly explain the relevant parts of the optical setup here. The interested reader is referred to [29] as well as appendix for a detailed derivation.

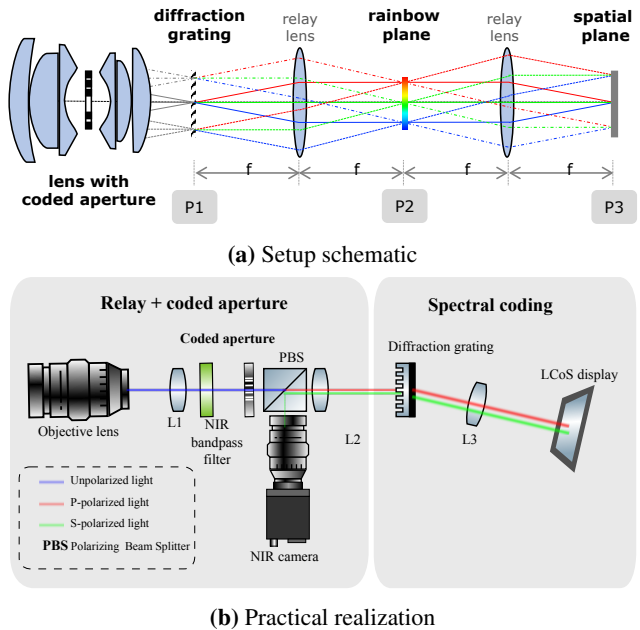


Figure 2: **Schematic for programmable spectral filter.** The optical architecture in (a) consists of a lens assembly with coded aperture which introduces spatial and spectral blurs. By placing an SLM in P2, the HSI of the scene can be spectrally filtered and sensed by a camera sensor on P3. (b) shows a compact realization of the optical setup.

4f system for spectral programming. We borrow the optical schematic for spectral programming from [29], shown in Fig. 2(a). Given the HSI, $H(x, y, \lambda)$, that is focused on the grating at P1, we seek to derive the intensity on planes P2 and P3. The intensity on rainbow plane P2,

$$I_4(x, y) = a^2(-x, -y) * \left(S \left(\frac{x}{f\nu_0} \right) \tilde{c} \left(\frac{x}{f\nu_0} \right) \right), \quad (4)$$

where $S(\lambda) = \int_{(x,y)} H(x, y, \lambda)$ is spectrum of the scene, $\tilde{c}(\lambda)$ is response of the optical system, and ν_0 is the density of groves in mm^{-1} . The intensity on image plane P3,

$$I_5(x, y) = \int_{\lambda} \left(H(x, y, \lambda) * \left| \frac{1}{\lambda^2 f^2} A \left(-\frac{x}{\lambda f}, -\frac{y}{\lambda f} \right) \right|^2 \right) d\lambda, \quad (5)$$

where $A(u, v)$ is the 2D Fourier transform of $a(x, y)$. The key observation from (4), (5) is that a coded aperture placed on plane P2 causes a spectral blur given by $a(x, y)$ and a spatial blur given by $\left| A \left(-\frac{x}{\lambda f}, -\frac{y}{\lambda f} \right) \right|^2$. As shown in Fig. 3, a slit causes a severe spatial blur, whereas an open aperture causes large spectral blur. The solution is to introduce an invertible blur in both domains, which can be achieved using a coded aperture, shown in the last column. We use

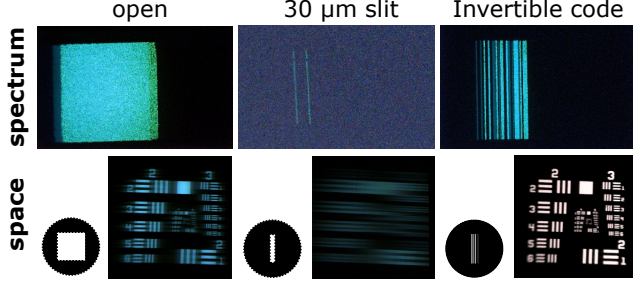


Figure 3: **Spatio-spectral resolution tradeoff.** A slit is capable of high spectral resolution whereas an open aperture is capable of high spatial resolution but both are inappropriate for high spatio-spectral HSI imaging. In contrast, a coded aperture introduces an invertible spatial and spectral blurs which can then be deconvolved. Figure reproduced with permission from [29].

the same coded aperture that was used in [29], as it is designed to promote invertibility in both domains.

Optical setup. Our optical setup is in principle similar to Fig. 2(a). We place a spatial light modulator on the rainbow plane (P2) and sensor on spatial plane (P3) to achieve spectral filtering. The optimized binary code [29] is placed in the lens assembly Figure 2(b) shows a schematic of a practical implementation of the same optical setup. We use a Liquid Crystal on Silicon (LCoS) display as a spatial light modulator for spectral filtering.

Effect of coded aperture. Given the HSI of the scene, $H(x, y, \lambda)$, the coded aperture introduces spatial and spectral blurs in the following way,

$$\hat{H}(x, y, \lambda) = \left(A \left(\frac{x}{\lambda f}, \frac{y}{\lambda f} \right) * H(x, y, \lambda) \right) * a(\lambda \nu_0 f, y), \quad (6)$$

i.e., all operations are now performed on a modified version of the HSI of the scene. Given a spectral profile $s_k(\lambda)$, the proposed setup directly computes filtered image,:

$$\hat{f}_k(x, y) = \int_{\lambda} \hat{H}(x, y, \lambda) s_k(\lambda) c(\lambda) d\lambda, \quad (7)$$

by loading $s_k(\lambda)$ on the spatial light modulator. With the optical setup in place, we will next see how to use the programmable spectral filter to perform optical classification.

4. Learning Discriminant Filters

With camera that is capable of capturing images with arbitrary spectral profiles, we pursue two questions; one, how many filters are required for classifying K classes, and two, what spectral filters maximize classification accuracy. The

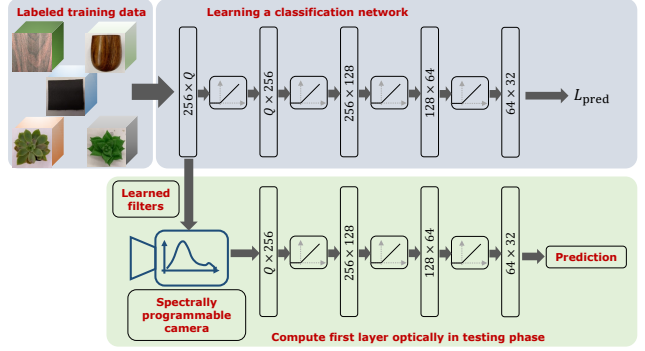


Figure 4: **Proposed optical classifier.** The proposed optical classifier broadly consists of two stages. In the first stage, we learn the weights of a neural network with spectrum as input and class label as output. The training process outputs the set of discerning filters, marked "learned filters" in the image. In testing stage, we filter the HSI of the scene with the learned filters, thereby replacing the first layer of the classifier with an optical implementation. This results in a high accuracy, per-pixel classifier while requiring far fewer measurements than the size of the HSI.

questions above are closely tied to the type of classifier under consideration. We detail the two classifier architectures we explore in this paper which help answer the questions above. Note that any classifier which relies on the linear projection can be used. For the sake of exposition, we only evaluate SVM and neural networks.

4.1. Support Vector Machine

SVMs provide a binary, linear classifier by learning a separating hyperplane on the training dataset. Given a set of data points $\{\mathbf{x}_k, y_k\}_{k=1}^N$, where $y_k \in \{0, 1\}$ is the label of \mathbf{x}_k , SVM seeks to solve the following optimization problem,

$$\min_{\mathbf{w}, c} \frac{1}{N} \sum_{k=1}^k \max(0, 1 - y_k(\mathbf{w}^T \mathbf{x}_k + c)) + \lambda \|\mathbf{w}\|^2, \quad (8)$$

where λ is a tuning parameter. The output of solving the optimization problem is the vector \mathbf{w} and intercept c . In the context of optical classification, \mathbf{w} is the filter that maximizes accuracy for binary decision. For K -class decision, we choose a *one-vs-all* classification strategy, which uses K hyperplanes, and hence K spectral filters.

4.2. Deep Neural Networks

Deep neural networks (DNNs) provide a richer alternative to SVMs. We model the first linear unit of the DNN to be the programmable spectral filter and train a model whose input is the spectral profile at a pixel and whose output is the material class label as a one-hot vector. While there are many possible architectures, we choose a simple, five-layer

Method	Classifier	Coding strategy	#Measurements	Accuracy
Santara et al.	DNN	Non-linear, spatial and spectral	220	96.7% (reported)
Hu et al.	DNN	Convolutional, spectrum-only	220	90.16% (reported)
Lee et al.	DNN	Convolutional, spatial and spectral	220	93.6% (reported)
Melgani et al.	SVM	Linear, spectrum-only	16	84% (computed)
This paper	DNN	Linear, spectrum-only	16	90% (Computed)

Figure 5: **Simulations on the Indian Pines dataset.** We compare state-of-the-art classifiers against the classifiers proposed in this paper. By *reported* we report the accuracy figures listed in the respective papers, while *computed* results were generated by us. A key feature of our optical setup is that it can only compute linear projections of spectra. While this leads to reduction in accuracy, the number of captured images are far fewer.

neural network as an example with all layers being fully connected. Figure 4 gives a brief overview of the proposed training and testing methodology. The weights of first fully connected layer, A_1 are the learned discriminating filters, and hence the first layer can be evaluated optically, thereby circumventing the need to measure the full spectrum at each pixel. The number of filters, Q depends on the number of materials and how easily they can be separated. In our experiments, we classified a total of 5 objects. We then varied the number of filters and computed mean classification accuracy. Based on this, we picked the optimal number of filters. We note that the idea of optically computing the first layer has been explored before in the context of designing color filter arrays [1] and processing light fields [2].

4.3. Simulations.

We compare SVM and the 5-layer DNN classifier to some of the state-of-the-art techniques in spectral-classification on the NASA Indian Pine dataset which consists of 220 spectral bands with 16 object classes. Figure 5 tabulates the accuracies with classifiers used in this paper in bold. We observe that the accuracy is lower than state-of-the-art, which is expected as we only use spectral information, while the other techniques use both spatial and spectral information. However, relying on a spectrum-only classifier lets us capture far fewer images than the number of spectral bands.

5. Experiments

We demonstrate capabilities of our setup for video-rate binary classification with binary SVM as well as matched filtering, and multi-class classification with multi-class SVM and DNNs.

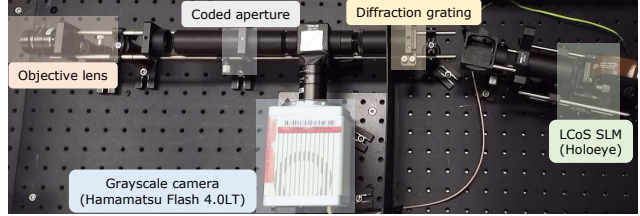


Figure 6: **Lab prototype.** The picture shows the lab prototype we built with only the major components marked. We used an objective lens of 8mm focal length, while all other lenses were 100mm.

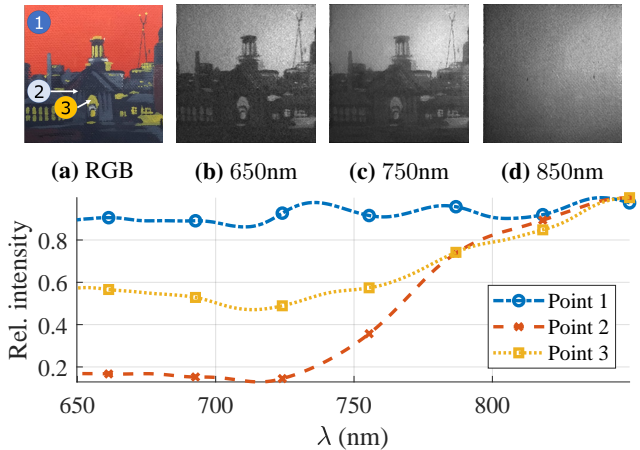


Figure 7: **Example HSI.** Our prototype is designed to capture images from 600nm to 900nm. (a) was captured using a cellphone while (b)-(d) are images captured by our setup. Bottom row shows spectral profiles at three marked points. Note how all the pigments disappear at $\lambda = 850\text{nm}$ in (d).

Optical setup. Figure 6 shows a photograph of the lab prototype we built along with labels for relevant components. A detailed optical layout along with the list of components is in appendix. Our SLM is a Holoeye LCoS SLM with a frame rate of 60 Hz that works as a secondary monitor. We used an NIR-sensitive sCMOS camera (Hamamatsu ORCA Flash 4.0 LT). In order to classify materials accurately, we designed our system to image from 600nm to 900nm, which is the near infrared (NIR) regime. Our setup is capable of coding spectrum at a resolution of 3.3nm, giving us 100 spectral bands. Finally, the SLM acts as a dynamic spectrally-selective camera and hence can be directly used for measuring the complete HSI. To do so, we display permuted Hadamard patterns on the SLM to capture a $512 \times 512 \times 256$ dimensional HSI. Figure 7 shows an example of captured HSI of an acrylic painting.

Calibration. Our optical setup broadly requires calibration of the code resulting in spectral blur, calibration of wavelengths and finally, spatial PSF. We use narrow-band



(a) Visible and NIR images.

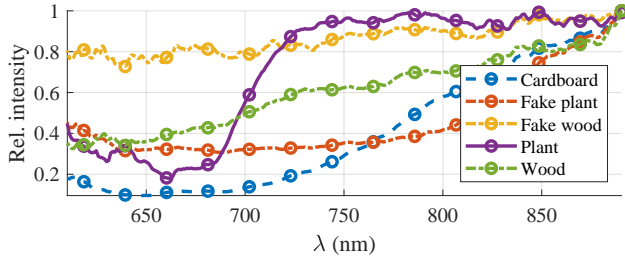
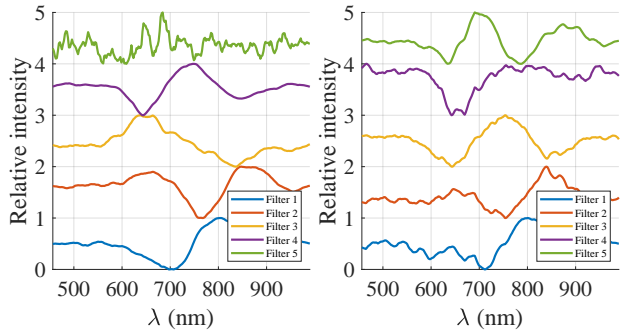


Figure 8: **Material dataset.** The figure shows false colored images of the 5 different materials we collected for classification purpose. (b) shows average spectra of the materials as measured by our lab prototype.



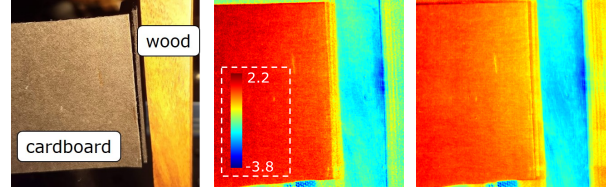
(a) SVM filters

(b) DNN filters

Figure 9: **Learned filters.** The output of a multi-class SVM is K separating hyperplanes, which results in K filters, shown in (a). Similarly, the DNN architecture consists of several layers, of which the first layer is linear. Hence the training process results in weights shown in (b) that can be used as spectral filters.

lasers for calibrating both code and wavelengths, and use a $10\mu\text{m}$ pinhole for calibrating spatial PSF. Details are available in the appendix.

Dataset. We show classification results with a total of five types of subjects: 1) black cardboard, 2) varnished wood, 3) wood-textured paper, 4) real plants, and 5) plastic plants. The choice of objects stems from similarity of these materials (plants vs plastic plants) in visible wavelengths, while



(a) RGB image (b) Full HSI scan + projection (256 meas.) (c) optical projection (2 meas.)

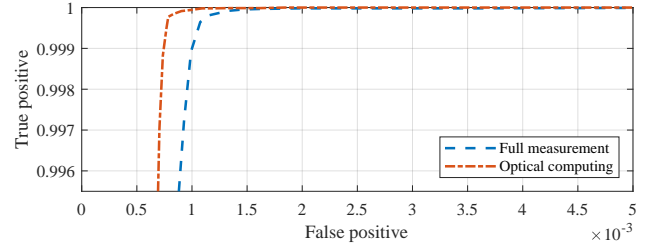


Figure 10: **Advantage of optical computing.** We show an example of binary classification between cardboard and wood (a) using per-pixel SVM. Optical computing achieves higher accuracy with far fewer measurements.

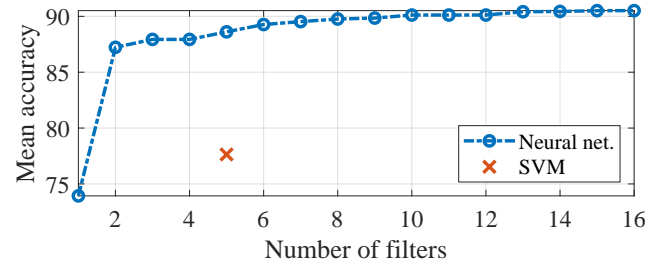
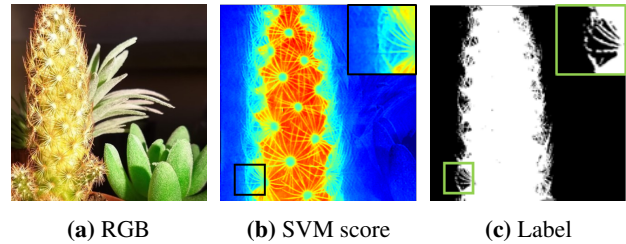


Figure 11: **Accuracy vs. number of filters.** The plot shows accuracy as a function of number of filters. The accuracy increases initially and then saturates. We hence use the knee point of the curve as the optimal number of filters.



(a) RGB (b) SVM score (c) Label

Figure 12: **Per-pixel classification.** Due to per-pixel operation with high spatial resolution, our imager can clearly identify the micro-structures such as the cactus thorns by capturing only two images instead of the complete HSI.

having distinctly different spectra in NIR domain. We collect one HSI for each of the materials and manually label them, giving a total of 5 HSI for training. Figure 8 shows

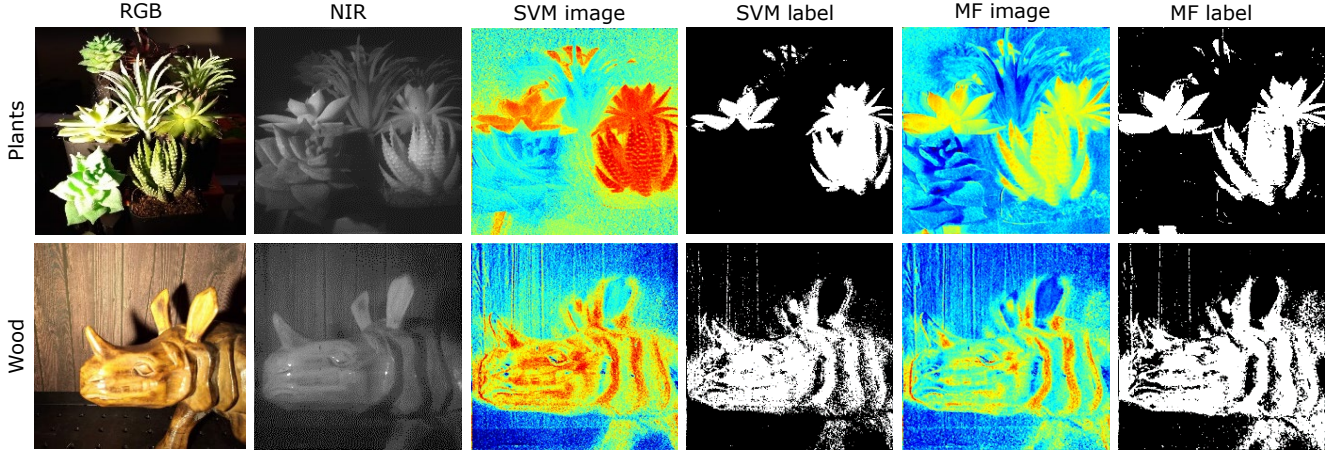


Figure 13: **Various binary classifiers.** We compare binary classification using SVM and matched filtering (MF). First row is a comparison of real wood (rhino) and fake wood (background, printed paper), while the second row is real and fake plants. Due to dynamic programming capability, we can classify with arbitrary filters and hence utilize any classifier that relies on linear projection of spectrum at each pixel.

visible and false-NIR images, as well as the average spectrum for each material. We note that none of the objects used for training the classifiers were reused in testing phase.

Training classifiers. We trained two classifiers – multi-class SVM and DNN with varying number of filters. For SVM, we used Scikit-Learn [28] in a one-vs-all configuration which learned a total of 5 spectral filters. The learned spectral filters are shown in Fig. 9 (a)

DNNs were trained with the network architecture shown in Fig. 4 with loss function set to cross entropy. The number of spectral filters were varied from 1 to 20 to compare performance. We learned the network using the PyTorch framework [27] with learning rate set to 10^{-3} for a total of 50 epochs. We then extracted weights of first layer and used them as spectral filters. The learned filters are shown in Fig. 9 (b). Further details about the learning process are included in appendix.

Figure 11 shows a plot of accuracy as a function of number of filters, Q . Accuracy of the classifier increases sharply initially and then saturates which implies that more spectral filters do increase accuracy but there is diminishing returns after a point. Based on this, we used 3, 5, 10 filters for comparisons in our real experiments.

Handling scale of features. A key requirement of any classifier is that the scale of features be same during training and testing. A common practice is to set the norm of feature at (x_0, y_0) , $\|H(x_0, y_0, \lambda)\|$ to unity, or the maximum value to unity. In our case, this requires having knowledge of the complete spectral profile, which defeats the purpose of optical computing. instead, we normalize our measurements

with sum of the spectrum, $\int_{\lambda} H(x_0, y_0, \lambda)$, which can be measured by displaying a spectral profile with all ones. The measured featured vectors are then,

$$I_{\text{sum}}(x_0, y_0) = \int_{\lambda} H(x_0, y_0, \lambda) d\lambda \quad (9)$$

$$\tilde{I}_k(x_0, y_0) = \int_{\lambda} H(x_0, y_0, \lambda) s_k(\lambda) d\lambda \quad (10)$$

$$I_k(x_0, y_0) = \frac{\tilde{I}_k(x_0, y_0)}{I_{\text{sum}}(x_0, y_0)} \quad (11)$$

We scale the spectra the same way even while training, which makes the scaling consistent. Hence any set of measurements with spectral profiles requires one extra image.

Binary classification. The simplest task possible with our optical setup is a binary classification, where the label at each pixel belongs to one of the two possible classes. In such a situation, one may either use a linear SVM where the spectral filter is the learned supporting hyperplane, \mathbf{w} , or use a matched filter, where the spectral filter is difference of spectra of the two classes, $s_1(\lambda) - s_2(\lambda)$. Figure 10 evaluates the advantages of optical classification. (b) visualizes the SVM score at each pixel obtained by scanning the complete HSI and then computing the projection to the SVM hyperplane, which requires a total of 256 measurements. In contrast, optical projection, shown in (c) requires only two images. Bottom row shows the Receiver operating Characteristic (RoC) of classification for both cases. The SNR advantage is evident; the area under the curve for optical projection (0.7194) is higher than full measurement and then projection (0.7912). Figure 12 shows classification of a real cactus surrounded by several plastic plants. The SVM score

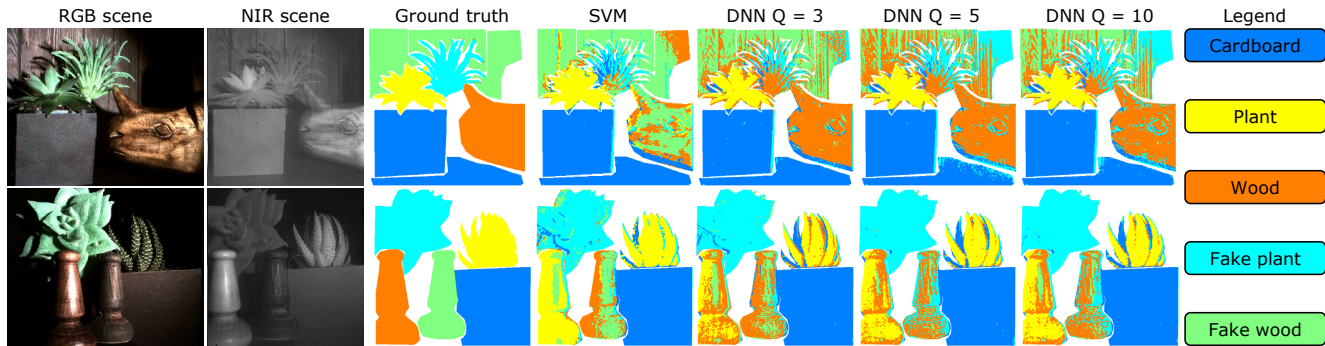


Figure 14: **Optical classification.** We show two examples of classification where the linear operations are directly computed in the optical domain. The ground truth labels were obtained by hand annotation. SVM required a total of 11 measurements for five filters, whereas DNN with 3, 5, 10 filters required 7, 11, 21 images respectively. The RGB images shown how the objects are not easily discernable in the visible domain, while they are accurately identified in the NIR domain along with optical classification.

		Accuracy: 72.55%						
		Cardboard	Fake plant	Fake wood	Plant	Wood		
Output Class	Cardboard	89.2% 70456	9.4% 1393	0.0% 0	0.0% 0	0.0% 0		
	Fake plant	3.9% 3099	56.4% 8329	8.6% 5971	3.6% 804	10.0% 2419		
	Fake wood	0.4% 278	12.3% 1820	63.4% 43822	10.3% 2334	34.8% 8399		
	Plant	0.8% 597	3.6% 533	0.3% 174	72.3% 16321	0.7% 174		
	Wood	5.8% 4551	18.2% 2681	27.8% 19204	13.7% 3100	54.5% 13150		
		Cardboard	Fake plant	Fake wood	Plant	Wood		

(a) SVM

		Accuracy: 77.33%						
		Cardboard	Fake plant	Fake wood	Plant	Wood		
Output Class	Cardboard	91.1% 69997	11.6% 1753	0.1% 21	0.0% 0	0.1% 78		
	Fake plant	1.8% 1413	47.2% 7112	5.7% 2226	0.3% 47	15.2% 9824		
	Fake wood	0.9% 654	26.1% 3929	91.1% 35541	2.5% 351	25.1% 16178		
	Plant	1.2% 887	0.9% 132	0.7% 276	97.2% 13753	4.3% 2751		
	Wood	5.1% 3900	14.2% 2140	2.5% 960	0.0% 3	55.3% 35683		
		Cardboard	Fake plant	Fake wood	Plant	Wood		

(b) Neural net.

Figure 15: **Confusion matrices for classifiers.** Neural networks typically outperform SVM. Among object classes, “wood” and “fake wood” get confused the most, as their spectra are similar. In contrast, “cardboard” is most different from all other spectra and hence has high accuracy.

in (b) as well as the labels show that our setup is capable of resolving very thin structures such as the cactus thorns. Figure 13 shows classification results for real vs. plastic plants and real vs. fake wood with binary SVM as well as matched filtering. Note that the objects are not easily discernable in RGB domain, while they are easily isolated after spectral filtering. Figure 1 shows a video rate classification of a real plant and a fake plant. The video was captured at 4 frames per second with alternating spectral profile and all ones pattern. Note how the real plant is tracked across all frames, while the fake plant is ignored.

Multi-class classification. We test the SVM and DNN filters learned on training data to classify a scene made of various materials from the set of five materials. Figure 14 shows classification results for two scenes for various techniques.

The ground truth annotation was obtained by manually annotating the objects, and then this was used for measuring the accuracy of classification. Visually, DNNs outperform SVM, as is visible from the accurate classification of the wooden rhinoceros head. Figure 15 shows a confusion matrix for SVM and DNN with 5 filters. The accuracies are not very high as we depend on spectral features alone. Accuracy can be significantly increased if spatial information is used along with spectral profiles. This is done by first capturing the Q images and then using the spatial information to classify.

6. Discussions and Conclusion

We propose a per-pixel material classifier that relies on a high resolution programmable spectral filter. We achieve this by learning spectral filters that can achieve high classification accuracy and then measure images of the scene with the learned filters. Owing to a simple, per-pixel decoding strategy, we can achieve classification at video rates. We showed several compelling real world examples with emphasis on binary video-rate and multi-class classification.

Limitations. A key limitation of our setup is the assumption that the pixels come from a single material class. Some real world examples are made of a mixture of materials at each class, an example being land cover. In such a case, outputting just a class label may not suffice but relative probabilities of each class is desired. This can be achieved by modifying the classifiers to output a score for each material at each pixel instead of most probable class.

7. Acknowledgment

The authors acknowledge support from the National Science Foundation under the CAREER award CCF-1652569,

the Expeditions award IIS-1730147, and the National Geospatial-Intelligence Agency's Academic Research Program (Award No. HM0476-17-1-2000).

A. Theoretical background

A.1. Image formation model.

We specified a simplified version of image formation model where we said that the HSI of the scene can be represented as $H(x, y, \lambda)$. We discuss a more precise model here.

Consider the scene's spectral reflectance function, $H_R(x, y, \lambda)$, where we assume that each point in 3D space is well modeled by Lambertian reflectance. Let $L(\lambda)$ be the spectral distribution of a spatially uniform light source. The HSI of the scene under this illumination is then given by,

$$H_o(x, y, \lambda) = H_R(x, y, \lambda)L(\lambda), \quad (12)$$

which was the signal model we used in the main paper. Then, given a camera with spectral response $C(\lambda)$, the measured image is,

$$\begin{aligned} I(x, y) &= \int_{\lambda} H_o(x, y, \lambda)C(\lambda)d\lambda \\ &= \int_{\lambda} H_R(x, y, \lambda)L(\lambda)C(\lambda)d\lambda \end{aligned} \quad (13)$$

From the above equation, we see that the camera measures spectral albedo of the scene's HSI and not the spectral reflectance of models. However, this is not a problem, as long as the light's spectral distribution is known *a priori*.

B. Learning details

We provide details about our training process with emphasis on choice of parameters and hyperparameters. We captured a total of 1,000,000 spectral profiles over 5 material types. For each classifier, we used 20% for training, 5% for validation and 80% for testing. We found that the testing accuracy did not improve even if we used more than 20% data.

Support Vector Machine. We used the function `LinearSVC` from Scikit-Learn [28] for training a *one-vs-all* SVM. The only hyperparameter of tuning was penalty for the hyperplanes, C , which was tuned by performing a grid search over the log space from 10^{-5} to 1. Hyperparameter search was done through a 3-fold cross-validation.

Neural Networks. We used PyTorch [27] for training our neural network (DNN) classifiers. The architecture used for learning is shown in 4 and the details of each layer is provided in Table 1. Q is the number of filters and was varied

Layer	Components
1 (Filters)	Linear (256xQ), ReLU, Dropout (0.1)
2	Linear (Qx256), ReLU, Dropout (0.1)
3	Linear (256x128), ReLU, Dropout (0.1)
4	Linear (128x64), ReLU, Dropout (0.1)
5	Linear (64x32), ReLU, Dropout (0.1)
6 (Output)	Linear (32x5)

Table 1: **Components of our DNN classifier.** All the layers are formed of fully connected layers with a ReLU and dropout added after each linear layer. Here, Q is the number of spectral filters and was variable in our experiments to compare performance. The output was a single linear layer. During training process, we used cross-entropy as loss function.

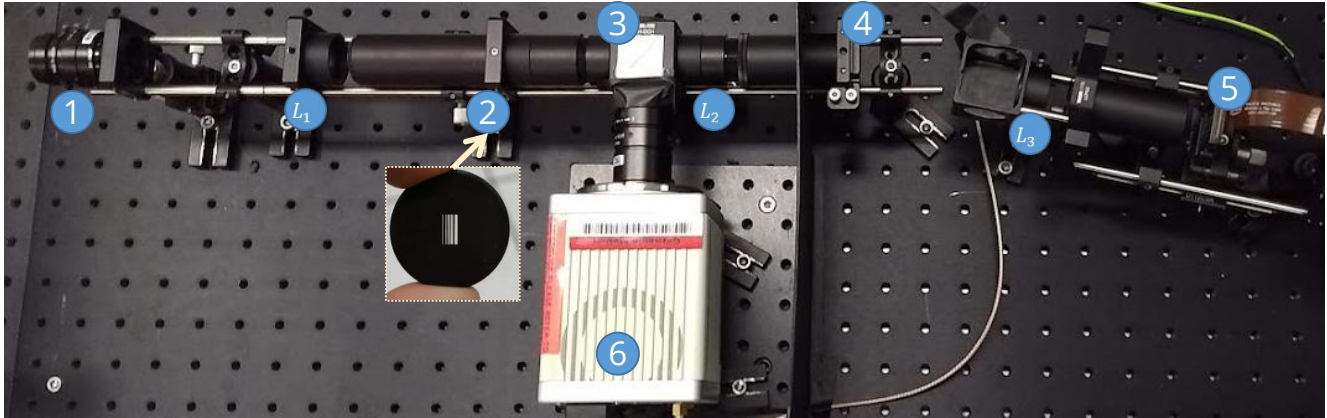
from 1 to 20 to evaluate performance as a function of measurements. We trained the network with an initial learning rate of 10^{-3} and trained for a total of 60 epochs. The filters were initialized with a principal component analysis (PCA) decomposition of training data. This led to smoother filters and higher accuracy. For each Q , we picked the model with best accuracy on validation dataset.

C. Hardware details

Hardware prototype. Figure 16 shows a picture of our lab prototype with names of major components. The last lens in the setup was replaced by a 50mm objective lens focused at infinity. This led to a better spatial resolution than an achromat.

Calibration. As described in the main paper, our setup required calibration of coded aperture, wavelengths and spatial PSF. We detail the calibration procedure here.

1. *Coded aperture calibration:* This is required to capture the code that blurs the spectrum. We measure the coded aperture by illuminating a spectrally flat object (such as spectralon) with a laser of known wavelength and scanning the complete HSI. We then average all spatial pixels to get the spectrum of the scene. Since a laser can be treated as a discrete delta, the measured spectrum will be the coded aperture. We threshold the measured spectrum appropriately to get the binary coded aperture, as shown in Fig. 17 (a).
2. *Wavelength calibration:* To find the correspondence between band index (1 - 256) and the corresponding wavelengths, we capture two scenes, each one comprised of a spectrally flat object illuminated by a narrowband laser light source. The averaged spectrum of the HSI is a blurred version of the laser spectrum. By deconvolving with the previously estimated coded aperture, we get location of the laser in terms of band index. We use this



- | | | |
|--|---|---|
| <p>1 8mm Objective Lens
Thorlabs MVL8M23</p> <p>2 Coded Aperture
Printed by Photoplot store with a feature size of 10μm</p> | <p>3 Polarizing Beam Splitter
Moxtex Wiregrid Polarizing Beamsplitter</p> <p>4 300 groves/mm Diffraction Grating
Dynasil G300TN31.7CC</p> | <p>5 NIR LCoS Display
Holoeye HED6001 with NIR-enhancement</p> <p>6 sCMOS camera
Hamamatsu ORCA Flash 4.0 LT</p> <p>L_1 L_2 L_3 100mm Achromat
Thorlabs AC254-100-B</p> |
|--|---|---|

Figure 16: **Lab prototype.** A picture of the lab prototype along with major components marked with details. We skipped details about opto-mechanical components such as cage plates and posts to avoid clutter. The inset image shows the printed mask we used as coded aperture.

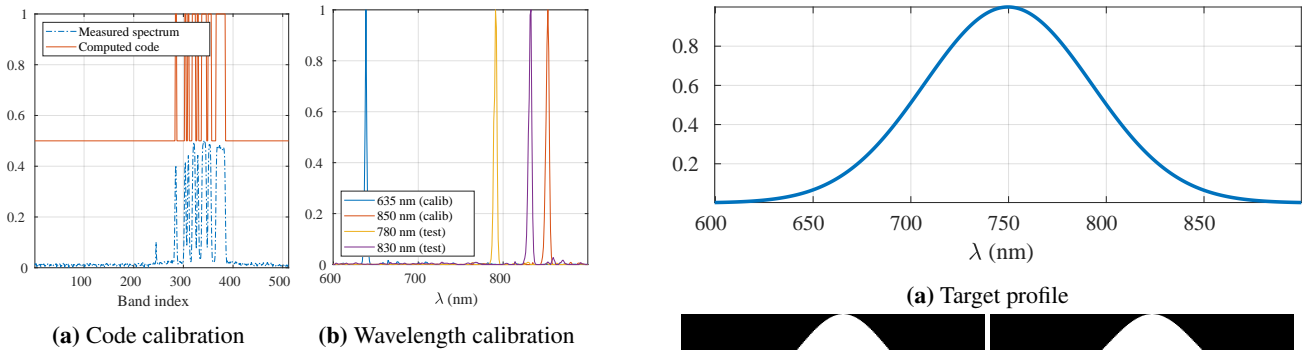


Figure 17: **Wavelength calibration.** We first estimate the blur due to coded aperture by capturing a scene illuminated by a narrowband light source (635nm laser), giving us the code in (a). We then calibrate the correspondence between band index and wavelengths by capturing two separate scenes illuminated by known laser light sources (635nm, 850nm). The results of the two calibration are show in (b), where we capture two more scenes with 780nm and 830nm laser.

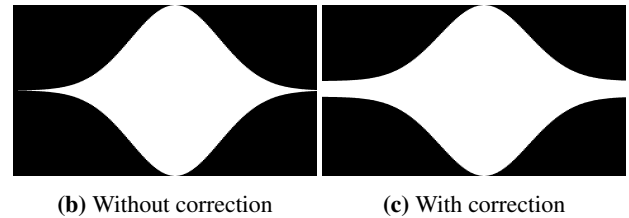


Figure 18: **Displaying desired spectral profile.** Given a target profile (a), we display a binary image on the SLM, as shown in (b), which ensures grayscale modulation despite wavelength-dependent gamma curve. However, since the SLM is $2f$ away from the camera sensor, there will be effects of diffraction. We counter this by adding a small DC offset, as shown in (c).

information along with laser wavelength to calibrate the correspondence.

3. *Spatial PSF*: To find the spatial blur kernel, we capture a single image of a 10 μ m pinhole. Since the PSF is well conditioned, deblurring the spatial images is well conditioned.

4. *Radiometric calibration of SLM*: The LCoS SLM in our optical setup is based on twisted-nematic design, and hence has different gamma curves for different wave-

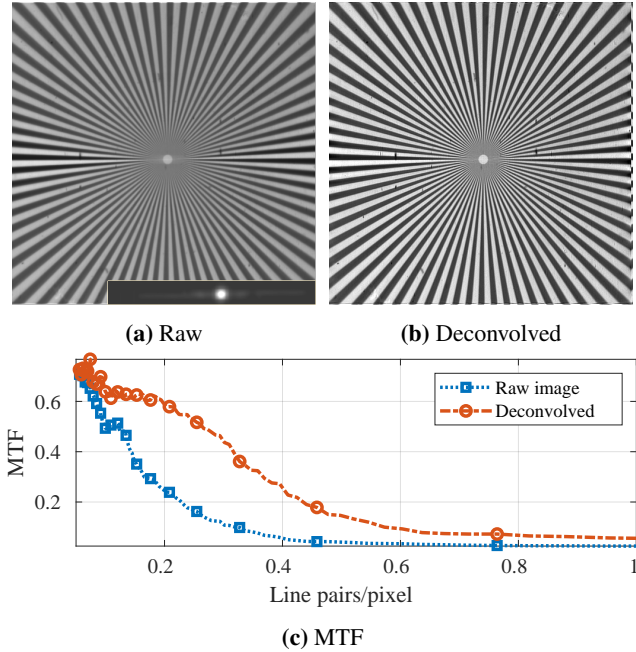


Figure 19: **Spatial deconvolution.** Due to design of an invertible spatial blur, the optical setup is capable of high resolution after deconvolution. (a) shows a raw image, with enlarged PSF in inset, (b) shows result of wiener deconvolution, and (c) shows a comparison of modulation transfer function (MTF). There is a marked increase in resolution both quantitatively and qualitatively.

lengths. Since the spectrum on the SLM is a blurred version of the true spectrum, we cannot perform a column-wise gamma correction. Instead, we use the SLM only as a binary modulator and achieve grayscale modulation by varying height of each column as shown in Fig 18 (b). This way, the SLM has a linear gamma curve for all wavelengths.

Figures of merit. Our setup is capable of achieving spectral resolution of up to 3.3nm over the wavelength range of 600 – 900nm, which is the designed resolution (see KRISM.pdf for further details). Due to invertible spatial blur, our setup is capable of high resolution after deconvolution. Figure 19 visualizes the captured image in (a) and deconvolved image in (b) of a sector star target. (c) shows plot of Modulation Transfer Function (MTF) as a function of line pairs per pixel. Image was deconvolved using simple Wiener deconvolution. The MTF30 after deconvolution was 0.45 linepairs/pixel.

Handling diffraction due to SLM. Since the SLM is placed $2f$ away from the image plane, any pattern displayed on SLM will lead to a diffraction blur. To counter this effect,

we always display ones in the middle of the pattern to be displayed on the SLM. This reduces the effect of diffraction while adding a simple offset to the data, which can be removed by capturing image with only the central part open.

References

- [1] A. Chakrabarti. Learning sensor multiplexing design through back-propagation. In *Advances in Neural Information Processing Systems*, pages 3081–3089, 2016. 5
- [2] H. G. Chen, S. Jayasuriya, J. Yang, J. Stephen, S. Sivaramkrishnan, A. Veeraraghavan, and A. Molnar. Asp vision: Optically computing the first layer of convolutional neural networks using angle sensitive pixels. In *IEEE Conf. Computer Vision and Pattern Recognition*, 2016. 2, 5
- [3] Y. Chen, H. Jiang, C. Li, X. Jia, and P. Ghamisi. Deep feature extraction and classification of hyperspectral images based on convolutional neural networks. *Trans. Geoscience and Remote Sensing*, 54(10):6232–6251, 2016. 2
- [4] E. Cloutis. Review article hyperspectral geological remote sensing: evaluation of analytical techniques. *International J. Remote Sensing*, 17(12):2215–2242, 1996. 1
- [5] N. Colthup. *Introduction to infrared and Raman spectroscopy*. Elsevier, 2012. 1
- [6] N. Dobigeon, J.-Y. Tourneret, C. Richard, J. C. M. Bermudez, S. McLaughlin, and A. O. Hero. Nonlinear unmixing of hyperspectral images: Models and algorithms. *IEEE Signal Processing Magazine*, 31(1):82–94, 2014. 2
- [7] M. Fauvel, J. Chanussot, J. A. Benediktsson, and J. R. Sveinsson. Spectral and spatial classification of hyperspectral data using svms and morphological profiles. In *IEEE Int. Geoscience and Remote Sensing Symposium*, pages 4834–4837, 2007. 2
- [8] M. Goel, E. Whitmire, A. Mariakakis, T. S. Saponas, N. Joshi, D. Morris, B. Guenter, M. Gavrilu, G. Borriello, and S. N. Patel. Hypercam: hyperspectral imaging for ubiquitous computing applications. In *ACM Intl. Joint Conf. Pervasive and Ubiquitous Computing*, pages 145–156, 2015. 3
- [9] R. Gross, I. Matthews, and S. Baker. Fisher light-fields for face recognition across pose and illumination. In *Joint Pattern Recognition Symposium*, pages 481–489, 2002. 2
- [10] A. B. Hamida, A. Benoit, P. Lambert, and C. B. Amar. 3-d deep learning approach for remote sensing image classification. *Trans. Geoscience and Remote Sensing*, 56(8):4420–4434, 2018. 2
- [11] J. C. Harsanyi and C.-I. Chang. Hyperspectral image classification and dimensionality reduction: An orthogonal subspace projection approach. *IEEE Trans. Geoscience and Remote Sensing*, 32(4):779–785, 1994. 1
- [12] M. He, B. Li, and H. Chen. Multi-scale 3d deep convolutional neural network for hyperspectral image classification. In *Intl. Conf. Image Processing*, 2017. 2
- [13] W. Hu, Y. Huang, L. Wei, F. Zhang, and H. Li. Deep convolutional neural networks for hyperspectral image classification. *Journal of Sensors*, 2015, 2015. 2
- [14] M. H. Kim, T. A. Harvey, D. S. Kittle, H. Rushmeier, J. Dorsey, R. O. Prum, and D. J. Brady. 3d imaging spec-

- troscopy for measuring hyperspectral patterns on solid objects. *ACM Trans. Graphics (TOG)*, 31(4):38, 2012. 1
- [15] H. Lee and H. Kwon. Contextual deep cnn based hyperspectral classification. In *Intl. Geoscience and Remote Sensing Symposium*, 2016. 2
- [16] Y. Li, H. Zhang, and Q. Shen. Spectral-spatial classification of hyperspectral imagery with 3d convolutional neural network. *Remote Sensing*, 9(1):67, 2017. 2
- [17] J. W. Lichtman and J.-A. Conchello. Fluorescence microscopy. *Nature methods*, 2(12):910, 2005. 1
- [18] B. Liu, X. Yu, P. Zhang, X. Tan, A. Yu, and Z. Xue. A semi-supervised convolutional neural network for hyperspectral image classification. *Remote Sensing Letters*, 8(9):839–848, 2017. 2
- [19] C. Liu and J. Gu. Discriminative illumination: Per-pixel classification of raw materials based on optimal projections of spectral brdf. *IEEE trans. pattern analysis and machine intelligence*, 36(1):86–98, 2014. 2
- [20] S. P. Love and D. L. Graff. Full-frame programmable spectral filters based on micromirror arrays. *J. Micro/Nanolithography, MEMS, and MOEMS*, 13(1):011108, 2014. 3
- [21] Y. Luo, J. Zou, C. Yao, X. Zhao, T. Li, and G. Bai. Hsi-cnn: A novel convolution neural network for hyperspectral image. In *Intl. Conf. Audio, Language and Image Processing*, 2018. 2
- [22] F. Melgani and L. Bruzzone. Classification of hyperspectral remote sensing images with support vector machines. *IEEE Trans. geoscience and remote sensing*, 42(8):1778–1790, 2004. 2
- [23] A. Mohan, R. Raskar, and J. Tumblin. Agile spectrum imaging: Programmable wavelength modulation for cameras and projectors. In *Computer Graphics Forum*, 2008. 3
- [24] M. O’Toole and K. N. Kutulakos. Optical computing for fast light transport analysis. *ACM Trans. Graph.*, 29(6):164, 2010. 3
- [25] Z. Pan, G. Healey, M. Prasad, and B. Tromberg. Face recognition in hyperspectral images. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 25(12):1552–1560, 2003. 1
- [26] J.-I. Park, M.-H. Lee, M. D. Grossberg, and S. K. Nayar. Multispectral imaging using multiplexed illumination. In *IEEE Intl. Conf. Computer Vision*, 2007. 3
- [27] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer. Automatic differentiation in pytorch. In *NIPS-W*, 2017. 7, 9
- [28] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011. 7, 9
- [29] V. Saragadam and A. Sarankaranarayanan. KRISM—krylov subspace-based optical computing of hyperspectral images. *arXiv:1801.09343*, 2018. 3, 4
- [30] V. Sharma, A. Diba, T. Tuytelaars, and L. Van Gool. Hyperspectral cnn for image classification & band selection, with application to face recognition. *Technical report KUL/ESAT/PSI/1604, KU Leuven, ESAT, Leuven, Belgium*, 2016. 2
- [31] Y. Tarabalka, J. Chanussot, and J. A. Benediktsson. Segmentation and classification of hyperspectral images using watershed transformation. *Pattern Recognition*, 43(7):2367–2379, 2010. 1