

3D Shape and Indirect Appearance by Structured Light Transport

Matthew O'Toole, John Mather, and Kiriakos N. Kutulakos, *Member, IEEE*

Abstract—We consider the problem of deliberately manipulating the direct and indirect light flowing through a time-varying, general scene in order to simplify its visual analysis. Our approach rests on a crucial link between stereo geometry and light transport: while direct light always obeys the epipolar geometry of a projector-camera pair, indirect light overwhelmingly does not. We show that it is possible to turn this observation into an imaging method that analyzes light transport in real time in the optical domain, prior to acquisition. This yields three key abilities that we demonstrate in an experimental camera prototype: (1) producing a live indirect-only video stream for any scene, regardless of geometric or photometric complexity; (2) capturing images that make existing structured-light shape recovery algorithms robust to indirect transport; and (3) turning them into one-shot methods for dynamic 3D shape capture.

Index Terms—Light transport, coded exposure, coded illumination, epipolar constraints, dynamic 3D shape capture, structured light 3D scanning, inter-reflections, subsurface scattering, direct/global separation, multi-path interference, primal-dual coding

1 INTRODUCTION

A common assumption in computer vision is that light travels along *direct* paths, i.e., it goes from source to camera by bouncing at most once in the scene. While this assumption works well in many cases, light propagation through natural scenes is actually a much more complex phenomenon: light reflects and refracts, it undergoes specular and diffuse inter-reflections, it scatters volumetrically and creates caustics, and it may do all of the above in the same scene. Analyzing all these phenomena with a conventional camera is a hard, open problem—and is even harder when the scene is dynamic and light transport changes unpredictably.

Despite the problem's intrinsic difficulty, indirect transport is a major component of real-world appearance [1] and an important cue for scene and material understanding [2]. It is also a major factor preventing broader use of structured-light techniques, which largely assume direct or low-frequency light transport (e.g., 3D laser scanning [3], [4], active triangulation [5], [6] and photometric stereo [7]).

As a step toward analyzing scenes that exhibit complex light transport, in this paper we introduce novel computational cameras for imaging such scenes in real time and for recovering their 3D shape. Our focus is on the general case where the scene is unknown; its motion and photometric properties unrestricted; and its illumination comes from a projector in general position.

Working from first principles, we show that two families of transport paths dominate image formation in a projector-camera system: *epipolar paths*, which satisfy the well-known

epipolar constraint [8] and contribute to a scene's direct image, and *non-epipolar paths* which contribute to its indirect image. Crucially, while the contributions of these paths are hard to separate computationally once an image has been captured, the paths themselves can be blocked optically *before* acquisition takes place. This motivates an alternative approach to imaging complex scenes in which the key question is how to optically control the relative contribution of epipolar and non-epipolar paths in the image being captured.

Using this idea as a starting point, we show how to realize this type of control over light transport with the help of a programmable sensor mask and a programmable projector that operate at rates much faster than the camera itself (i.e., masks and projection patterns change tens to hundreds of times during the exposure of a single image). Thus the optical behavior of the imaging system is determined by the precise sequence of masks and projection patterns used during image exposure. Here we apply this general technique, which we call *Structured Light Transport (SLT)*, to the following four imaging problems and derive their associated masks and projection patterns:

- *indirect-only imaging*: capture an image that records only contributions from indirect light;
- *two-shot direct-only imaging*: capture two images whose difference contains only contributions from direct light;
- *indirect-invariant structured light*: given any conventional structured-light pattern used for 3D shape acquisition, capture a view of the scene under that pattern that is guaranteed to be invariant to indirect light; and
- *one-shot multi-pattern imaging*: given any $S \geq 2$ conventional structured-light patterns used for 3D shape acquisition, capture one image that contains S separate views of the scene “packed” into it, each corresponding to a different structured-light pattern.

• The authors are with the Department of Computer Science, University of Toronto, 40 St. George St., Bahen Centre, Toronto, Ontario, Canada M5S 2E4. E-mail: {motoole, jmather, kyros}@cs.toronto.edu.

Manuscript received 20 Nov. 2014; revised 28 Oct. 2015; accepted 17 Nov. 2015. Date of publication 28 Mar. 2016; date of current version 10 June 2016. Recommended for acceptance by A. Martinez, R. Basri, R. Vidal, and C. Fermuller.

For information on obtaining reprints of this article, please send e-mail to: reprints@ieee.org, and reference the Digital Object Identifier below. Digital Object Identifier no. 10.1109/TPAMI.2016.2545662

Little is currently known about how to solve these problems in the general setting we consider. Our solutions, while firmly rooted in computer vision, operate exclusively in the optical domain and require no computational post-processing: our implementation is a physical device that just outputs live video; this video is optionally processed after acquisition by standard 3D reconstruction algorithms [5] which can be oblivious to the complexity of light transport occurring in a scene. The device itself is a novel combination of existing off-the-shelf components—a conventional video camera operating at 28Hz, a pair of synchronized digital micro-mirror devices (DMDs) operating at 2.7 to 24 kHz for sensor masking and pattern projection, and optics for coupling them.

From a practical point of view, our work offers four main contributions over the state of the art. First, it is the first demonstration of an “indirect-only video camera,” i.e., a camera that outputs a live stream of indirect-only video for general scenes—exhibiting arbitrary motion, caustics, specular inter-reflections and numerous other transport effects (Fig. 1). Prior work on indirect imaging was either constrained to static scenes [10], [11], or assumed diffuse/low-frequency transport [2], [12] and accurate 2D motion estimation [12]. Second, we show how to capture—with just one SLT shot—views of a scene that are invariant to indirect light. This is particularly useful for imaging dynamic scenes and represents an advance over direct-only imaging [2], [10], which requires at least two images. Third, we show that *any* ensemble of structured-light patterns can be made robust to indirect light, regardless of the patterns’ frequency content. This involves simply switching from conventional to SLT imaging—without changing the patterns or the algorithm that processes them. As such, our work stands in contrast to prior work on transport-robust structured light, which places the onus on the design of the patterns themselves [6], [13], [14], [15]. Fourth, we show that SLT imaging can turn any multi-pattern 3D structured-light method into a one-shot technique for dynamic shape capture. Thus an entire family of previously-inapplicable techniques can be brought to bear on this much-studied problem [5], [16], [17], [18], [19] in order to improve depth map resolution and robustness to indirect light. As a proof of concept, we demonstrate in Fig. 14 the reconstruction of dense depth and albedo from individual frames of monochrome video, acquired by combining indirect-invariant SLT imaging and conventional six-pattern phase-shifting.

Conceptually, our work has one essential difference from conventional structured-light methods in computer vision [2], [5]: instead of controlling light only at its source by projecting patterns, we control light at its destination as well, with a DMD mask in front of the camera pixels. This simultaneous projection and masking makes it possible to analyze light transport *geometrically* (by blocking 3D light paths), rather than *photometrically* (by blocking certain transport frequencies and assuming constrained scene reflectance [2]). It also enables optical-domain implementations, which can have a significant speed and signal-to-noise ratio advantage over post-capture processing. Restricted forms of such imaging have enjoyed widespread use in confocal microscopy [20], [21],

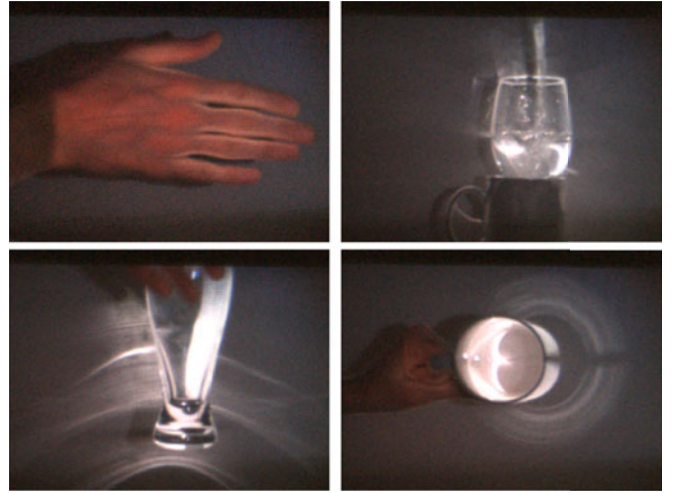


Fig. 1. Snapshots from raw live indirect video. Clockwise from top: (1) A hand; note the vein pattern and the inter-reflections between fingers. (2) Pouring water into a glass. (3) Caustics formed inside a mug from specular inter-reflections; note the secondary reflections to the board behind the mug and from the board onto the mug’s exterior surface. (4) Refractions and caustics from a beer glass. See Figs. 12 and 18 for more indirect-only images and [9] for videos.

[22] and were first used by O’Toole et al. [10] for macroscopic imaging of static scenes with a coaxial projector/camera. In all these cases, however, epipolar geometry is degenerate and stereo is impossible. While SLT imaging builds on that work, its premise, theory, applications, and physical implementation are different.

2 THE STEREO TRANSPORT MATRIX

We begin by relating scene geometry to the light transported from a projector to a camera in general position. Consider a scene whose shape potentially varies with time. If the camera and projector respond linearly to light, the scene’s instantaneous image satisfies the light transport equation [23]:

$$\mathbf{i} = \mathbf{T} \mathbf{p}. \quad (1)$$

where \mathbf{i} is the image represented as a column vector of I pixels; \mathbf{p} is the P -pixel projected pattern, also represented as a column vector; and \mathbf{T} is the scene’s $I \times P$ instantaneous light transport matrix.

Intuitively, element $\mathbf{T}[i, p]$ of the transport matrix specifies the total radiance transported from projector pixel p to image pixel i over all possible paths. As such, \mathbf{T} models image formation in very general settings: the scene may have non-Lambertian reflectance, it may scatter light volumetrically, exhibit specular inter-reflections, etc.

Anatomy of the stereo transport matrix. Since a projector and a camera in general position define a stereo pair, their transport matrix is best understood by taking two-view geometry into account. More specifically, we classify the elements of \mathbf{T} into three categories based on the geometry of their transport paths (Fig. 2):

- *Epipolar elements*, whose projector and camera pixels are on corresponding epipolar lines. These are the only elements of \mathbf{T} whose transport paths begin and end on rays that can intersect in 3D. By performing stereo calibration [8] and vectorizing patterns and

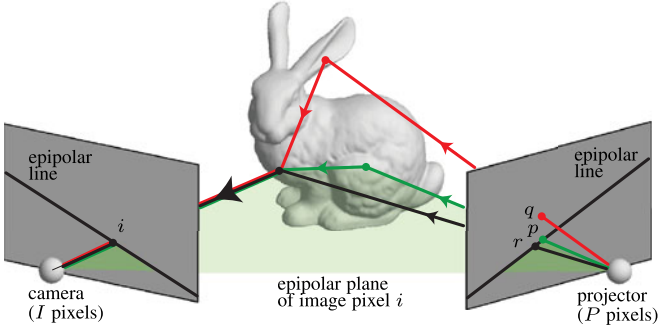


Fig. 2. Light transport in a stereo projector-camera system. Light can reach pixel i on the image in one of three general ways: by indirect transport from an arbitrary pixel p on the corresponding epipolar line (green path); by indirect transport from a pixel q that is *not* on that line (red path); or by direct surface reflection, starting from projector pixel r on the epipolar line (black path).

images according to Fig. 3, these elements can be made to occupy a known, time-invariant, block-diagonal subset of the transport matrix.

- *Non-epipolar elements*, whose projector pixel and camera pixel are not on corresponding epipolar lines. Non-epipolar elements are significant because they vastly outnumber the other elements of \mathbf{T} and *never*-account for direct transport. This is because their transport paths begin and end with rays that do not intersect, so light must bounce at least twice to follow them.
- *Direct elements*, whose camera and projector pixels are in stereo correspondence, i.e., they are the perspective projections of a visible surface point. Direct elements are where direct surface reflection actually occurs in the scene; although they always lie within \mathbf{T} 's epipolar blocks, their precise location is scene dependent and thus unknown. Indeed, locating the direct elements is equivalent to computing the scene's instantaneous stereo disparity map (Fig. 4).

We can therefore express every image of the scene as a sum of three components that arise from distinct “slices” of the transport matrix:

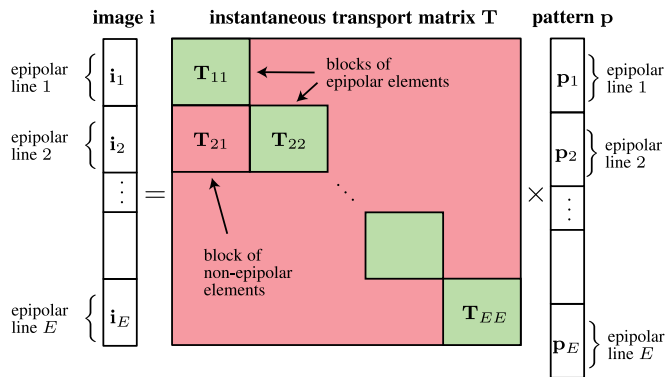


Fig. 3. The light transport equation when patterns and images are vectorized so that consecutive pixels on corresponding epipolar lines form subvectors \mathbf{p}_e and \mathbf{i}_e , respectively. Under this vectorization scheme, block \mathbf{T}_{ee} of the transport matrix describes transport from epipolar line f on the pattern to epipolar line e on the image. Blocks \mathbf{T}_{ee} , shown in green, contain the epipolar elements.

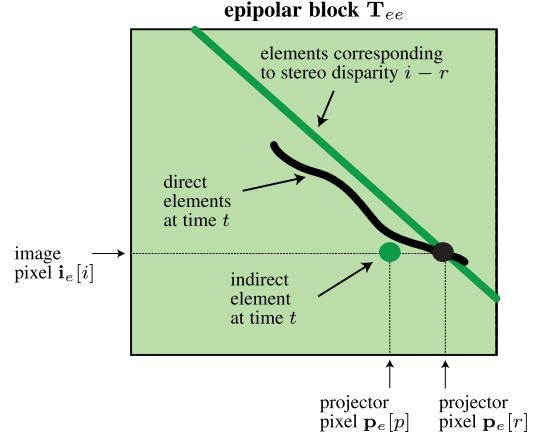


Fig. 4. Structure of an epipolar block \mathbf{T}_{ee} . Element $\mathbf{T}_{ee}[i, r]$ describes transport from projector pixel $\mathbf{p}_e[r]$ to image pixel $\mathbf{i}_e[i]$. This element is direct if and only the scene point projecting to both pixels is the same, i.e., the point's stereo disparity is $i - r$. The set of direct elements therefore represents the scene's instantaneous disparity map. Conventional stereo algorithms attempt to localize this set while assuming that the transport matrix is zero everywhere else—both inside and outside its epipolar blocks.

$$\mathbf{i} = \underbrace{\mathbf{T}^D}_{\text{direct image}} \mathbf{p} + \underbrace{\mathbf{T}^{EI}}_{\text{epipolar indirect image}} \mathbf{p} + \underbrace{\mathbf{T}^{NE}}_{\text{non-epipolar indirect image}} \mathbf{p}, \quad (2)$$

where the $I \times P$ matrices \mathbf{T}^D , \mathbf{T}^{EI} and \mathbf{T}^{NE} hold the direct, epipolar indirect, and non-epipolar elements, respectively, and are zero everywhere else.

3 DOMINANCE OF NON-EPIPOLAR TRANSPORT

Although in theory all three image components in Eq. (2) may contribute to scene appearance, in practice their contributions are not equal. The key observation underlying our work is that the non-epipolar component is very large relative to the epipolar indirect for a broad range of scenes:

$$\mathbf{i} \approx \underbrace{\mathbf{T}^D}_{\text{direct image}} \mathbf{p} + \underbrace{\mathbf{T}^{NE}}_{\text{non-epipolar indirect image}} \mathbf{p}. \quad (3)$$

We call this the *non-epipolar dominance assumption*. The transport matrix is much simpler when this assumption holds because we can treat it as having a time-invariant structure with two easily-identifiable parts: the epipolar blocks, which contribute only to the direct image, and the non-epipolar blocks, which contribute only to the indirect.

To motivate this assumption on theoretical grounds, we prove that it holds for two very general scene classes: (1) scenes whose transport function for indirect elements is square-integrable everywhere and (2) generic scenes containing pure specular reflectors and transmitters. These two cases can be thought of as representing opposite extremes, with the former covering low-frequency transport phenomena such as diffuse inter-reflection and diffuse isotropic subsurface scattering [24] and the latter covering transport whose frequency content is not band limited. In particular, we prove the following:

Proposition 1. *If \mathbf{T}^{EI} and \mathbf{T}^{NE} are discretized forms of transport functions that are square-integrable and positive over the rectified projector and image planes, then*

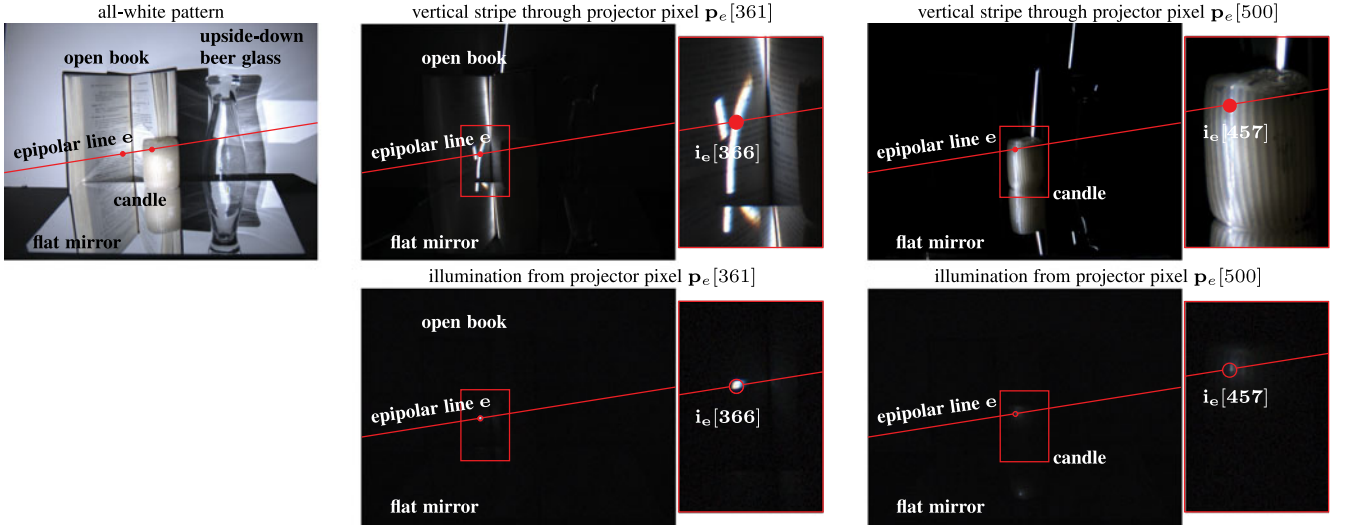


Fig. 5. Experimental validation of non-epipolar dominance for a scene containing diffuse, translucent, refractive and mirror-like objects. From left-to-right: (1) View under an all-white projection pattern. (2) View when just one white vertical stripe is projected onto the scene. The many bright regions in this image occur because the stripe illuminates the book’s pages in three different ways: directly from the projector, by diffuse inter-reflection from the opposite page, and by specular reflection via the mirror. Their existence makes the scene hard to reconstruct with conventional techniques such as laser-stripe 3D scanning [4]. A magnified view of these regions is shown in the inset on the right. (3) View for another vertical stripe, part of which falls on the candle. The stripe appears very broad and poorly localized there, because of strong sub-surface scattering. (4) View when just one projector pixel illuminates the scene. Camera pixels along the epipolar line receive light travelling along both direct and epipolar indirect paths; note that, unlike (2), these camera pixels receive no light travelling along non-epipolar indirect paths. (5) View when a single pixel illuminates a point on the candle.

$$\lim_{\epsilon \rightarrow 0} \frac{\mathbf{T}^{\text{EI}} \mathbf{p}}{\mathbf{T}^{\text{NE}} \mathbf{p}} = \mathbf{0} \quad (4)$$

where division is entrywise, $\mathbf{0}$ is a vector of zeros, and ϵ is the pixel size for discretization.

Proposition 2. Two generic n -bounce specular transport paths that originate from corresponding epipolar lines do not intersect for $n > 1$.

See Appendix A for proofs. Intuitively, both propositions are consequences of a “dimensionality gap”: the set of transport paths contributing to the epipolar indirect image has lower dimension than the set of paths contributing to the non-epipolar image (Fig. 2). Thus contributions

accumulated in one image are negligible relative to the other in generic settings.

On the practical side, we have found non-epipolar dominance to be applicable quite broadly; see Figs. 5 and 6 for a detailed analysis of non-epipolar dominance in a complex scene, Figs. 12 to 18 for more examples, and [9] for videos confirming the assumption’s validity in a variety of settings.

4 IMAGING BY STRUCTURED LIGHT TRANSPORT

The rich structure of the stereo transport matrix cannot be exploited by simply projecting a pattern onto the scene. This is because projection gives no control over how light flows through the scene: all elements of \mathbf{T} —regardless of

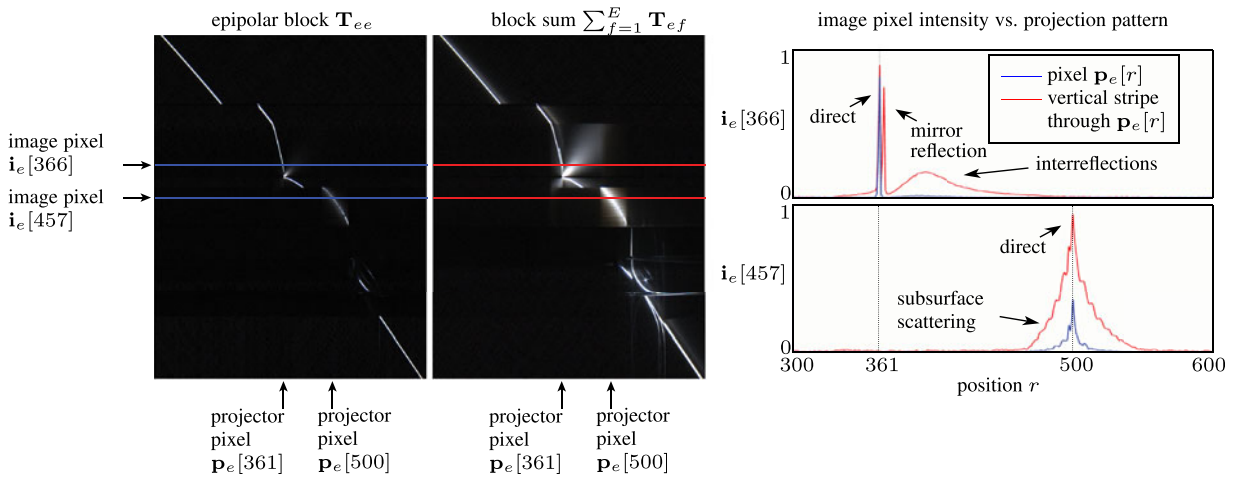


Fig. 6. *Left:* The epipolar block \mathbf{T}_{ee} for epipolar line e . We show \mathbf{T}_{ee} using the conventions of Fig. 4, i.e., its r th column comes from an image of the scene acquired with only projector pixel $\mathbf{p}_e[r]$ turned on. *Middle:* To assess the image contribution of non-epipolar transport, we acquire the block sum $\sum_{f=1}^E \mathbf{T}_{ef}$ and compare it to block \mathbf{T}_{ee} —observe that non-epipolar contributions indeed far surpass the epipolar indirect ones. To acquire the block sum, we capture images of the scene while sweeping a vertical stripe on the projector plane (see [9] for a video of the captured image sequence). The r th column of the block sum is given by the pixels on epipolar line e when the stripe is at $\mathbf{p}_e[r]$. *Right:* Horizontal cross-section of \mathbf{T}_{ee} and $\sum_{f=1}^E \mathbf{T}_{ef}$ for two image pixels. Observe that \mathbf{T}_{ee} ’s cross-section (blue) is sharp and unimodal whereas the block sum’s (red) is trimodal for one pixel and very broad for the other.

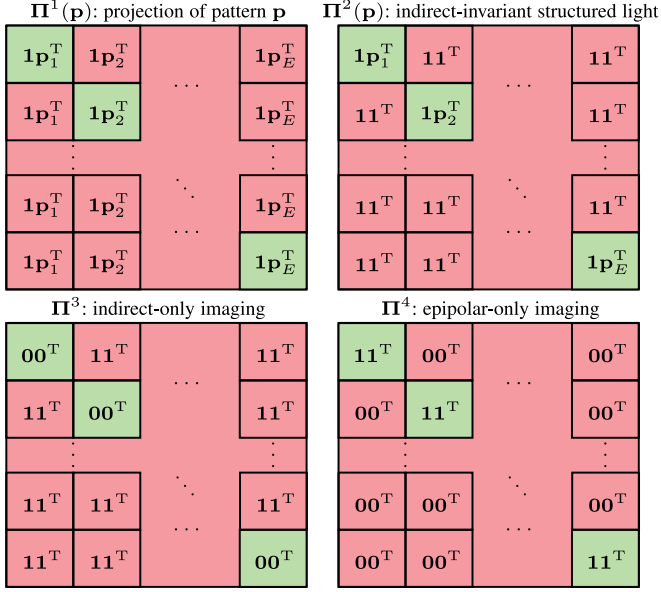


Fig. 7. The four basic probing matrices used in this paper. Their block structure mirrors the structure of \mathbf{T} in Fig. 3.

position—will participate in image formation. To make full use of \mathbf{T} 's structure, *we structure the flow of light itself*.

Our starting point is an imaging procedure first proposed by O'Toole et al. [10]. Its main advantage is that the contribution of individual elements of \mathbf{T} can be weighted according to a user-defined “probing matrix” Π :

$$\mathbf{i} = [\Pi \circ \mathbf{T}] \mathbf{1}, \quad (5)$$

where \circ denotes entrywise (*a.k.a.* Hadamard) product and $\mathbf{1}$ is a column vector of all ones. Images captured this way are said to be the result of *probing* the scene's transport matrix with matrix Π . Conceptually, they correspond to a scene that is illuminated by an all-white pattern and whose transport matrix is $\Pi \circ \mathbf{T}$.

Two basic questions arise when considering Eq. (5) for image acquisition and shape recovery: (1) what should Π be, and (2) how to design an imaging system that implements the equation? The answers in [10] were restricted to static scenes and projector/camera arrangements that share a single viewpoint, none of which apply here. Below we focus on the first question—designing Π —and discuss live imaging of dynamic scenes in Section 5.

Conventional structured-light imaging. To gain some insight, let us re-cast as a probing operation the act of projecting a fixed pattern \mathbf{p} and capturing an image \mathbf{i} . Applying the vectorization scheme of Fig. 3 to the light transport equation and re-arranging terms we get for epipolar line e :

$$\mathbf{i}_e = \sum_{f=1}^E \mathbf{T}_{ef} \mathbf{p}_f = \left[\sum_{f=1}^E \underbrace{(\mathbf{1p}_f^T)}_{\text{block of probing matrix}} \circ \underbrace{\mathbf{T}_{ef}}_{\text{block of } \mathbf{T}} \right] \mathbf{1}, \quad (6)$$

where E is the number of epipolar lines. Equation (6) implies that projecting \mathbf{p} is equivalent to probing with the matrix $\Pi^1(\mathbf{p})$ shown in Fig. 7. Observe that if we capture images for a whole sequence of projection patterns—as is often the case in structured-light systems—the non-epipolar

blocks of the probing matrix will be different for each pattern. Indirect transport will therefore contribute to each captured image differently in a way that strongly depends on the particular pattern. This makes structured-light 3D scanning difficult when indirect transport is present because its contributions cannot be easily identified and removed.

Indirect-invariant structured light. The contribution of indirect transport becomes much easier to handle if we ensure it is *the same for every pattern*. Since this contribution is dominated by the non-epipolar blocks of the transport matrix, we can achieve (almost) complete invariance to indirect transport by probing with a matrix whose non-epipolar blocks are independent of \mathbf{p} . In particular, probing with the matrix $\Pi^2(\mathbf{p})$ in Fig. 7 yields

$$\mathbf{i}_e = \underbrace{\left[(\mathbf{1p}_e^T) \circ \mathbf{T}_{ee} \right] \mathbf{1}}_{\text{direct image (depends on } \mathbf{p})} + \underbrace{\left[\sum_{f=1, f \neq e}^E \mathbf{T}_{ef} \right] \mathbf{1}}_{\text{non-epipolar indirect image (ambient)}}. \quad (7)$$

The image in Eq. (7) has two properties: (1) its direct component is identical to the direct component we would get by projecting \mathbf{p} conventionally onto the scene, and (2) its non-epipolar component is independent of \mathbf{p} . This independence essentially turns indirect contributions into an “ambient light” term that does not originate from the projection pattern.¹ To see the practical significance of this independence, Fig. 15 compares views of a scene under conventional and one-shot indirect-invariant structured light, for the same projection pattern.

An important corollary of Eq. (7) is that indirect-invariant structured light images can be acquired for *any* sequence of patterns—regardless of frequency content or other properties—using the corresponding sequence of probing matrices.

Indirect-only imaging. A notable special case of indirect-invariant structured light is to set \mathbf{p} to zero (matrix Π^3 in Fig. 7). This yields an image guaranteed to have no contributions from direct transport. Moreover, almost all indirect light will be recorded when non-epipolar dominance holds.

Epipolar-only imaging. The exact opposite effect can be achieved with a probing matrix that is zero everywhere except along the epipolar blocks (matrix Π^4 in Fig. 7). When non-epipolar dominance holds, images captured this way can be treated as (almost) purely direct.

One-shot, multi-pattern, indirect-invariant structured light. All four probing matrices in Fig. 7 produce views of the scene under a fixed illumination pattern \mathbf{p} . With probing, however, it is possible to capture—in just one shot—spatially-multiplexed views of the scene for a whole sequence of structured-light patterns, $\mathbf{p}(1), \dots, \mathbf{p}(S)$. The probing matrix to achieve this can be thought of as defining a “projection pattern mosaic,” much like the RGB filter mosaic does for color (Fig. 8). Moreover, we can confer invariance to indirect light by defining the mosaic in terms of *probing matrices* rather than conventional patterns.

1. Other examples of ambient terms with identical behavior include image contributions from the projector's black level and contributions from light sources other than the projector. Because such terms are often unavoidable yet easy to handle, many structured-light algorithms are designed to either recover them explicitly or be robust to their existence [5]. Non-zero ambient terms do, however, reduce contrast and may affect SNR.

color filter mosaic			6-pattern mosaic			6-pattern indirect-invariant mosaic		
R	G	R	p(1)	p(2)	p(3)	$\Pi^2(\mathbf{p}(1))$	$\Pi^2(\mathbf{p}(2))$	$\Pi^2(\mathbf{p}(3))$
G	B	G	p(4)	p(5)	p(6)	$\Pi^2(\mathbf{p}(4))$	$\Pi^2(\mathbf{p}(5))$	$\Pi^2(\mathbf{p}(6))$

Fig. 8. Example layouts for color RGB, monochrome 6-pattern, and monochrome 6-pattern indirect-invariant structured light imaging.

Specifically, suppose we partition the I image pixels into S sets and let $\mathbf{b}(1), \dots, \mathbf{b}(S)$ be binary vectors of size I indicating the pixel membership of each set. The matrix

$$\Pi^5(\mathbf{p}(1), \dots, \mathbf{p}(S)) = \sum_{s=1}^S [\mathbf{b}(s) \mathbf{1}^T] \circ \Pi^2(\mathbf{p}(s)), \quad (8)$$

interleaves the rows of S indirect-invariant probing matrices. Thus, probing with this matrix yields an image containing S sub-images, each of which is a view of the scene under a specific structured-light pattern in the sequence.

5 LIVE STRUCTURED-LIGHT-TRANSPORT IMAGING

The feasibility of probing comes from re-writing Eq. (5) as a bilinear matrix-vector product [10]:

$$\mathbf{i} = \sum_{t=1}^T \mathbf{m}(t) \circ [\mathbf{T} \mathbf{q}(t)], \quad (9)$$

where the transport matrix \mathbf{T} is constant in time and $\Pi = \sum_{t=1}^T \mathbf{m}(t)(\mathbf{q}(t))^T$ is a rank-1 decomposition of the probing matrix. According to Eq. (9), optical probing is possible by (1) opening the camera's shutter, (2) projecting pattern $\mathbf{q}(t)$ onto the scene, (3) using a pixel mask $\mathbf{m}(t)$ to modulate the light arriving at individual camera pixels, (4) changing the pattern and mask synchronously T times, and (5) closing the shutter. This procedure acquires one image; it was implemented in [10] for low-resolution probing matrices using an LCD panel for pixel masking, an SLR camera for image acquisition, and $T \in [100, 1,000]$.

Although results were promising, LCDs are not suitable for video-rate (30 Hz) probing: they refresh at 30–200 Hz, limiting T to an unusable 1–6 masks/projections per frame; and they have low transmittance, requiring long exposure times.

Our approach, on the other hand, is to use a pair of off-the-shelf digital micro-mirror devices (DMDs) for projection and masking (Figs. 10 and 11). These devices are compact, incur no light loss and can operate synchronously at 2.7–24 kHz. To implement Eq. (9), we couple them with a conventional video camera operating at 28 fps. This allows 96–800 masks/projections within the 36 msec exposure of each frame.² To our knowledge, such a coupling has not been proposed before.³

A major difference between LCDs and DMDs is that DMDs are *binary*. This turns the derivation of masks and projection patterns into a combinatorial optimization problem. Formally, given an *integer*⁴ probing matrix Π and an

upper bound on T , we seek a length- T rank-1 decomposition into binary vectors such that the decomposition approximates Π as closely as possible. This problem is difficult and we know of no general solution. Indeed, estimating the length of the shortest *exact* decomposition is itself NP-hard [25].

Our approach, below, is to derive randomized decompositions of Π that approximate Eq. (9) in expectation. Although our experience is that this approach works well in practice, it should not be treated as optimal.

Indirect-only imaging. Matrix Π^3 is a special case where short decompositions are easy. Let $\mathbf{q}(e)$ be a pattern whose pixels are 1 along epipolar line e and 0 everywhere else and let $\mathbf{m}(e)$ be a mask that is 1 everywhere except at epipolar line e . Then it is easy to show that $\Pi^3 = \sum_{e=1}^E \mathbf{m}(e)(\mathbf{q}(e))^T$. This corresponds to a sequence of mask/projection pairs where only one epipolar line is “off” in the mask and only the corresponding epipolar line is “on” in the pattern. Even though this decomposition is exact—and feasible for near-megapixel images—it has poor light efficiency because only one epipolar line is “on” at any time. To improve light efficiency we use random patterns instead, which yield good approximations that are much shorter.

Specifically, consider the random pattern

$$\mathbf{q} = \{\text{each epipolar line is 1 with probability } 0.5\}, \quad (10)$$

let the projection pattern $\mathbf{q}(t)$ be a sample of \mathbf{q} , and let the mask $\mathbf{m}(t)$ be equal to $\overline{\mathbf{q}(t)}$. See Fig. 9, Row 1 for an example of $\mathbf{q}(t)$, $\overline{\mathbf{q}(t)}$, and $\mathbf{m}(t)$. Taking expectations in Eq. (9), the epipolar line e of the expected image is given by

$$\mathcal{E}[\mathbf{i}_e] = \mathcal{E}[\overline{\mathbf{q}_e}] \circ \sum_{\substack{f=1 \\ f \neq e}}^E \mathbf{T}_{ef} \mathcal{E}[\mathbf{q}_f] = 0.25 \sum_{\substack{f=1 \\ f \neq e}}^E \mathbf{T}_{ef} \mathbf{1}, \quad (11)$$

where $\mathcal{E}[\cdot]$ denotes expectation. This is the result of probing with matrix Π^3 , albeit at one quarter of the “ideal” image intensity.⁵ Note that corresponding epipolar lines are never on at the same time in the pattern and mask; thus no epipolar transport path ever contributes to the captured image.

Epipolar-only imaging. Matrix Π^4 is a special case at the other extreme, where *no* short rank-1 decompositions exist. Since $\Pi^4 = \Pi^1(1) - \Pi^3$, we compute the result of probing with Π^4 by subtracting two adjacent video frames—one captured by projecting an all-white pattern and one captured by indirect-only imaging. Naturally, two-frame motion estimation may be necessary to handle fast-moving scenes (but we do not estimate motion in our experiments).

Indirect-invariant structured light. A perhaps counterintuitive result is that even though epipolar-only imaging requires two frames, indirect-invariant structured light requires just one. This is important because probing with matrix $\Pi^2(\cdot)$ is all we need for reconstruction with structured light. Let \mathbf{p} be an arbitrary structured-light pattern scaled to $[0, 1]$. Define mask $\mathbf{m}(t)$ to be a sample of \mathbf{q} from Eq. (10) and the pattern to be

5. Intuitively, since half the epipolar lines are “off” in the pattern and the mask, only 1/4th of the total light is transported from projector to camera.

2. See [10] for an analysis of the SNR advantage conferred by performing T mask/projection operations in a single exposure versus capturing an image for each projection pattern, at $1/T$ th the exposure.

3. The closest design we are aware of comes from confocal microscopy [20]. Its optical path was less challenging to implement, however, because imaging was both coaxial and orthographic.

4. Since any grayscale structured-light pattern \mathbf{p} must be quantized before projection, probing matrices are always integer, including $\Pi^2(\mathbf{p})$.

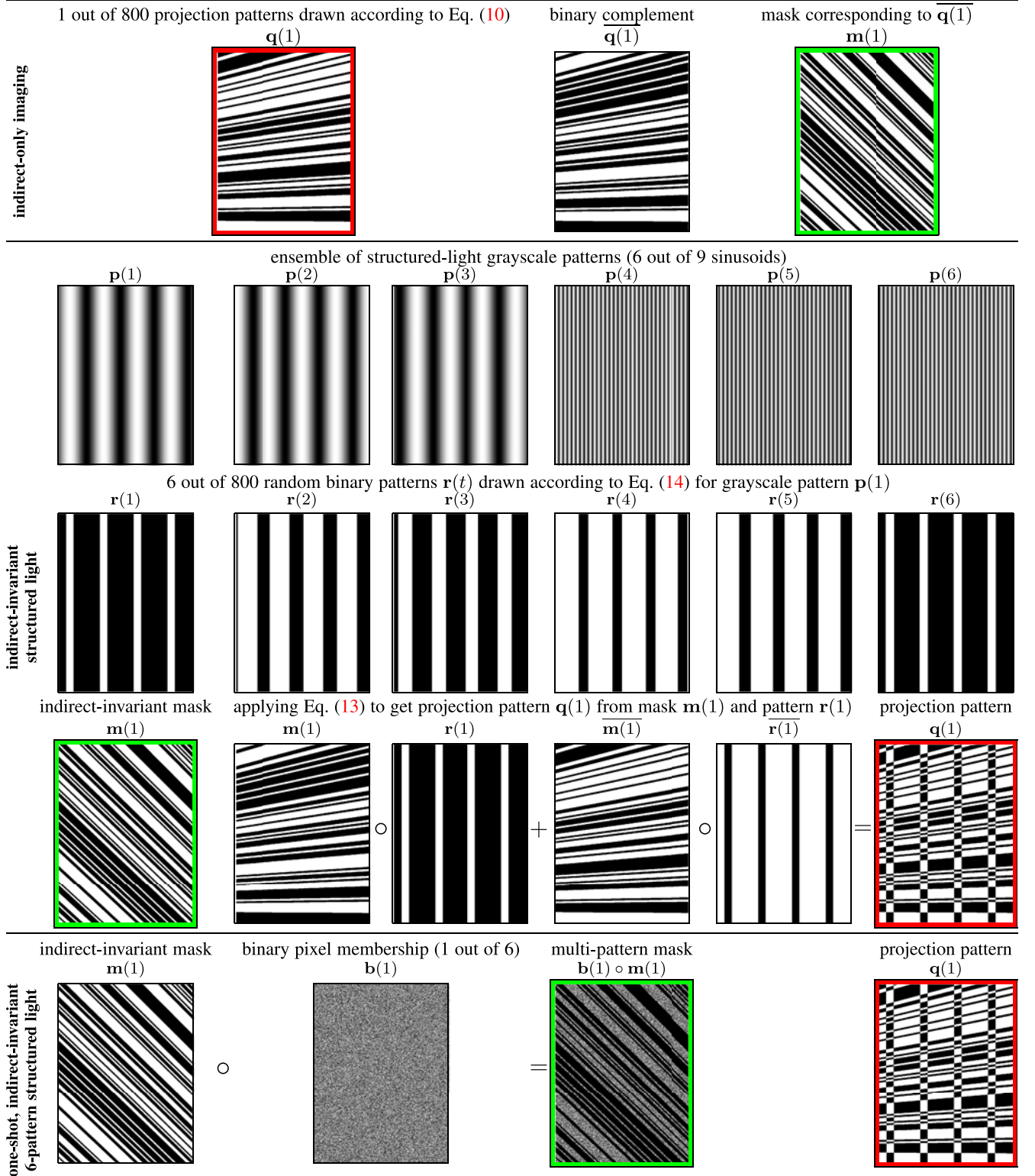


Fig. 9. Deriving random pattern/mask pairs for three cases of SLT imaging. The derived patterns and masks are indicated with red and green borders, respectively. *Row 1:* For indirect-only imaging, the patterns and masks are constant along epipolar lines, with approximately half of them “on.” *Row 2:* Six of the nine structured-light patterns we used. *Rows 3-4:* The masks for indirect-invariant structured light are identical to those for indirect-only imaging but the projection patterns differ. To generate them for a given grayscale structured-light pattern, we first generate a random sequence of binary patterns (Row 3) and then use that sequence, along with the sequence of masks, to compute the projection patterns. Row 4 shows one such example. *Row 5:* We generate pattern/mask pairs for 6-shot imaging as follows: (1) create 6 random binary images representing pixel membership for each pattern; (2) generate a sequence of 132 indirect-invariant binary pattern/mask pairs for each of 6 grayscale structured-light patterns, as outlined in Rows 2-4; (3) use the 792 projection patterns as is, and (4) multiply the masks element-wise with the associated pixel memberships. Row 5 shows one such calculation, for grayscale pattern $p(1)$.

$$\mathbf{q}(t) = \mathbf{m}(t) \text{ XNOR } \mathbf{r}(t) \quad (12)$$

$$\stackrel{\text{def}}{=} \mathbf{m}(t) \circ \mathbf{r}(t) + \overline{\mathbf{m}(t)} \circ \overline{\mathbf{r}(t)}, \quad (13)$$

where XNOR is the exclusive nor operator of binary vectors $\mathbf{m}(t)$ and $\mathbf{r}(t)$, and $\mathbf{r}(t)$ is a sample of yet another random pattern:

$$\tau = \{\text{pixel } p \text{ on epipolar line } e \text{ is } 1 \text{ with probability } \mathbf{p}_e[p]\}. \quad (14)$$

A pictorial illustration of Eq. (13) can be found in Fig. 9 on Row 4, with example random binary patterns $\mathbf{r}(t)$ sampled from τ (Eq. (14)) shown on Row 3. From calculations similar to Eq. (11), the expected image is

$$\begin{aligned} \mathcal{E}[\mathbf{i}_e] &= 0.5\mathbf{T}_{ee}\mathbf{p}_e + 0.25 \sum_{f=1, f \neq e}^E [\mathbf{T}_{ef}\mathbf{p}_f + \mathbf{T}_{ef}(\mathbf{1} - \mathbf{p}_f)] \\ &= \underbrace{0.5\mathbf{T}_{ee}\mathbf{p}_e}_{\substack{\text{direct image} \\ (\text{depends on } \mathbf{p})}} + \underbrace{0.25 \sum_{f=1, f \neq e}^E \mathbf{T}_{ef}\mathbf{1}}_{\substack{\text{indirect image (ambient)}}}, \end{aligned} \quad (15)$$

which is equivalent to the result of probing with $\Pi^2(\cdot)$.

One-shot, multi-pattern, indirect-invariant structured light. Here we use the mask for indirect-invariant structured light and temporally multiplex S random projection patterns—each defined by Eq. (13) and corresponding to a different structured-light pattern—across our “budget” of T total projections per video frame. Fig. 8 and Row 5 of Fig. 9 illustrate the construction of a multi-pattern indirect-invariant mask and corresponding projection pattern for $S = 6$. After the video is recorded, we “demosaic” each frame \mathbf{i} independently to infer S full-resolution images, one for each structured-light pattern. Following work on compressed sensing [26], [27] we do this by solving for S images that reproduce frame \mathbf{i} and are sparse under a chosen basis \mathbf{W} :

$$\text{minimize } \|\mathbf{W}^T[\mathbf{i}(1) \ \dots \ \mathbf{i}(S)]\|_n \quad (16)$$

$$\text{subject to } \left\| \sum_{s=1}^S \mathbf{b}(s) \circ \mathbf{i}(s) - \mathbf{i} \right\|_2 \leq \epsilon, \quad (17)$$

where $\|\cdot\|_n$ is a sparsity-inducing norm⁶ and $\mathbf{b}(s)$ is the binary vector holding pixel memberships for pattern s .

6 IMPLEMENTATION

Experimental prototypes. We created two experimental systems for performing optical probing according to Eq. (9):

- a low-speed, low-cost system for video-rate indirect-only and epipolar-only imaging (Fig. 10) whose components are listed in Table 1; and
- a high-speed system for indirect-invariant structured light shape acquisition and one-shot multi-pattern imaging (Fig. 11), whose components are also listed in Table 1.

6. We use the (1,2)-norm because it promotes group sparsity and thus concentrates non-zero terms to the same pixels across views.

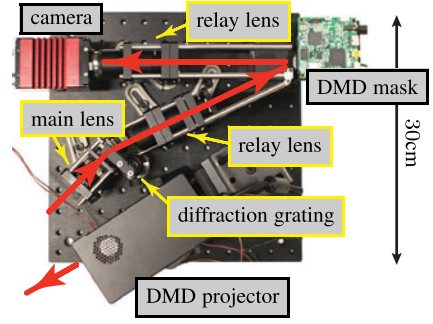


Fig. 10. Photo of our low-speed, low-cost prototype. The projector can be detached to change the stereo baseline. The optical path is shown in red.

Our low-speed, low-cost system included a color AVT GT1920C camera for acquisition, a Texas Instruments LightCrafter for pixel masking and a 100 lumen Keynote Photonics LightCrafter kit for projection. The DMDs were synchronized at 2.7 kHz, permitting $T = 96$ patterns and masks per video frame. The camera and DMD resolutions were quite different— $1,936 \times 1,456$ versus 608×684 —with each DMD pixel mapping to a 2×2 block of camera pixels. System calibration consists of computing the epipolar geometry between the two DMDs. We did this by first computing correspondences between the camera and each DMD separately. Patterns are uploaded to both DMDs once, at the beginning of an imaging session.

For our high-speed system, we used a monochrome AVT GT1920 camera and a pair of high-end DMDs from Texas Instruments (DLi 4130) with a 2,000 lumen light source. These operate at 22.2 kHz, permitting $T = 800$ patterns per video frame. Although the DMD resolution was fairly high at $1,024 \times 768$, its effective resolution was much lower, 484×364 , because of the different physical dimensions and orientation of the camera sensor and DMD.

In one-shot multi-pattern imaging, the effective DMD resolution was even lower, 256×256 , because of the scene’s limited extent within the camera’s field of view.

Indirect-only mask & projection patterns. We use random mask/pattern pairs like those shown in Row 1 of Fig. 9. To reduce the sensation of flicker by users who are physically present during video acquisition, we generate a random sequence of $T/2$ mask/pattern pairs and then generate a second mask/pattern sequence whose projection patterns are the binary complement of the first $T/2$ projection patterns. This ensures a stable perception because the image integrated by the eye (or by a mask-less camera) over the period of one video frame corresponds to a view of the scene under an all-white projection pattern.⁷

For indirect-only imaging, it is also important to ensure that no direct light “leaks” accidentally through the DMD mask. Such leaks can occur because of pixel misalignments between the DMD mask and the camera’s sensor; because of the binary rasterization of epipolar lines; and because of projector/camera defocus. To make

7. We emphasize that flicker is a purely subjective sensation that may be experienced by users who view the scene directly, without the benefit of the DMD mask. In particular, flicker *does not* occur in the videos captured by our prototypes.

TABLE 1
List of Parts for the Low-Speed System Shown in Fig. 10 and the High-Speed System Shown in Fig. 11

Item #	Part Description	Quantity	Model Name	Company	Low-speed Part	High-speed Part
1	color camera	1	GT1920C	Allied Vision Technologies	✓	
2	monochrome camera	1	GT1920	Allied Vision Technologies		✓
3	connector housing	2	WM1722-ND	Digi-Key Corporation	✓	
4	connector housing	2	WM1728-ND	Digi-Key Corporation		✓
5	power supply	1	T1228-Z12P-ND	Digi-Key Corporation	✓	
6	power supply	1	LC3000-Pro Power Supply	Keynote Photonics	✓	
7	low-speed DMD (projector)	1	LC3000-Pro Pico Projector	Keynote Photonics	✓	
8	low-speed DMD (mask)	1	DLP LightCrafter	Texas Instruments	✓	
9	high-speed DMD	2	DLi4130VIS-7XGA	Digital Light Innovations		✓
10	high-power LED light engine	1	High Power S2+ w/ LED	Digital Light Innovations		✓
11	fixed filter holder 40 mm Sq.	1	#54-997	Edmund Optics		✓
12	45 degree mounting adapter	1	#59-001	Edmund Optics		✓
13	crimp	4	WM1142CT-ND	Digi-Key Corporation	✓	✓
14	Hirose contact plug	1	HR1623-ND	Digi-Key Corporation	✓	✓
15	12 mm f/1.4 objective lens	1	Cinegon 1.4/12-0906	Schneider Optics	✓	✓
16	visible achromatic doublet pairs	2	MAP10100100-A	Thorlabs	✓	✓
17	300 grooves/mm transmission grating	1	GT25-03	Thorlabs	✓	✓
18	ring-activated threaded iris diaphragm	2	SM1D12D	Thorlabs	✓	✓
19	C-mount to SM1 adapter	1	SM1A9	Thorlabs	✓	✓
20	SM1 to C-mount adapter	1	SM1A10	Thorlabs	✓	✓
21	SM1 Coupler	1	SM1T10	Thorlabs	✓	✓
22	SM1 Lens Tube, 2 inch Thread Depth	1	SM1L20	Thorlabs	✓	✓
23	SM1 Lens Tube, 3 inch Thread Depth	1	SM1L30	Thorlabs	✓	✓
24	SM1-threaded cage plate	1	CP4S	Thorlabs	✓	✓
25	cage plate with 1.2 inch double bore	5	CP12	Thorlabs	✓	✓
26	cage plate with 35 mm aperture	4	CP03/M	Thorlabs	✓	✓
27	cylindrical lens mount	1	CH1A	Thorlabs	✓	✓
28	rod swivel coupler (set of four)	1	C2A	Thorlabs	✓	✓
29	rod end swivel connector (set of four)	1	C3A	Thorlabs	✓	✓
30	cage assembly rod, 2 inch long	2	ER2	Thorlabs	✓	✓
31	cage assembly rod, 3 inch long	4	ER3	Thorlabs	✓	✓
32	cage assembly rod, 4 inch long	4	ER4	Thorlabs	✓	✓
33	aluminum breadboard	1	MB3030/M	Thorlabs	✓	✓
34	12.7 mm x 40 mm optical post	1	TR40/M	Thorlabs	✓	✓
35	12.7 mm x 100 mm optical post	6	TR100/M	Thorlabs	✓	✓
36	post holder, 40 mm	1	PH40/M	Thorlabs	✓	✓
37	post holder, 100 mm	6	PH100/M	Thorlabs	✓	✓
38	studded pedestal base adapter	7	BE1/M	Thorlabs	✓	✓
39	small clamping fork	7	CF125	Thorlabs	✓	✓
40	30 mm single axis translation stage	1	#66-397	Edmund Optics	✓	✓
41	bottom adapter plate	1	#66-620	Edmund Optics	✓	✓
42	top adapter plate	1	#66-493	Edmund Optics	✓	✓
43	metric base plate	1	#54-975	Edmund Optics	✓	✓
44	thread-to-thread adapter	1	#56-323	Edmund Optics	✓	✓

Both systems use identical optics; the only differences between the two systems are (1) the DMD projector and mask, (2) their mounts, and (3) a color vs. monochrome camera. The camera in both systems outputs live video at a rate of 28 frames per second. For the low-speed system, each video frame requires 96 binary masks/projection patterns, i.e., each mask/pattern is active for 375 μsec of the frame's 36,000 μsec total exposure time. Each video frame of the high-speed system consists of 800 binary mask/projection patterns, i.e., each mask/pattern is active for 45 μsec .

indirect-only acquisition robust to such effects, we slightly dilate the “off” regions on the generated masks. This reduces the occurrence of such leaks at the expense of a slight reduction in light efficiency. We found this approach to be very effective in practice; a similar idea was used in [10].

Epipolar-only patterns. We generate epipolar-only video by operating the camera at 56 fps and configuring the DMD of our low-speed prototype as follows:

- *odd video frames:* display $T/2$ all-on mask/pattern pairs
- *even video frames:* display a sequence of $T/2$ indirect-only mask/pattern pairs.

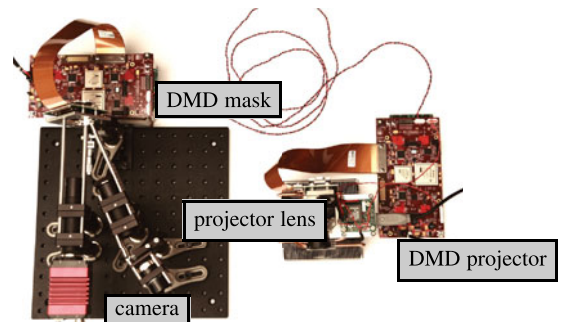


Fig. 11. Our high-speed system. The key differences between this system and that shown in Fig. 10 are a monochrome camera, the DMD mask, and the DMD projector.

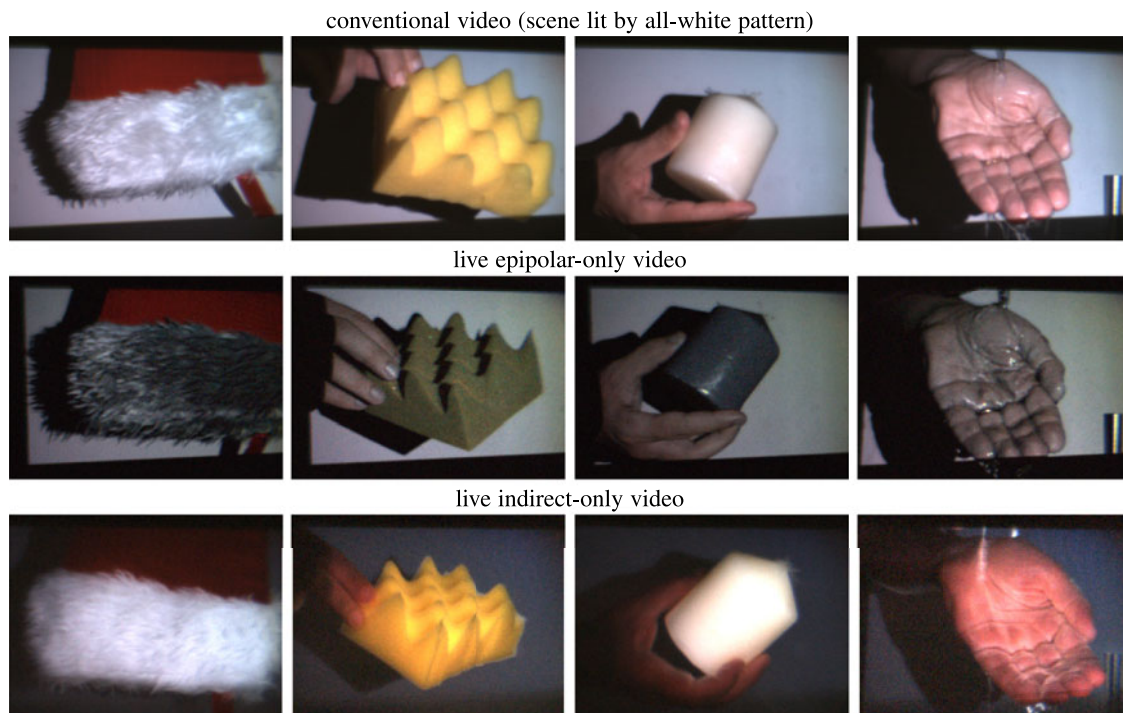


Fig. 12. Frames from conventional, epipolar-only, and indirect-only video. From left-to-right: (1) Faux-fur; note the marked difference between the epipolar-only component, which appears very shiny due to direct near-specular reflection, versus the diffuse appearance of the indirect-only component, caused by sub-surface scattering. (2) A piece of packing foam. (3) A translucent candle. The indirect-only frames of both the packing foam and candle demonstrate that the color of volumetric or translucent materials is often attributed to light traveling through the sub-surface. (4) Water flowing over a hand; the indirect-only frame makes apparent the very dramatic change in a hand's reflectance properties when water flows over it. We hypothesize that these changes are caused by scattering in the thin film of water flowing over the hand.

Epipolar-only video at 28fps is generated by (1) scaling the odd frames by 0.25 to account for the reduced intensity of indirect-only imaging (Eq. (11)) and (2) subtracting in real time the even frames from the scaled odd ones.

Indirect-invariant patterns. We generate a sequence of T mask/pattern pairs for each of S grayscale structured-light

patterns, as illustrated in Rows 2-4 of Fig. 9. We then capture one raw image of the scene for each of the S generated mask/pattern sequences. These S images are supplied, unaltered, to the 3D reconstruction algorithm.

One-shot, multi-pattern, indirect-invariant patterns. We generate a sequence of T mask/pattern pairs, as outlined in Row 5 of Fig. 9, and upload them to the DMDs. We then

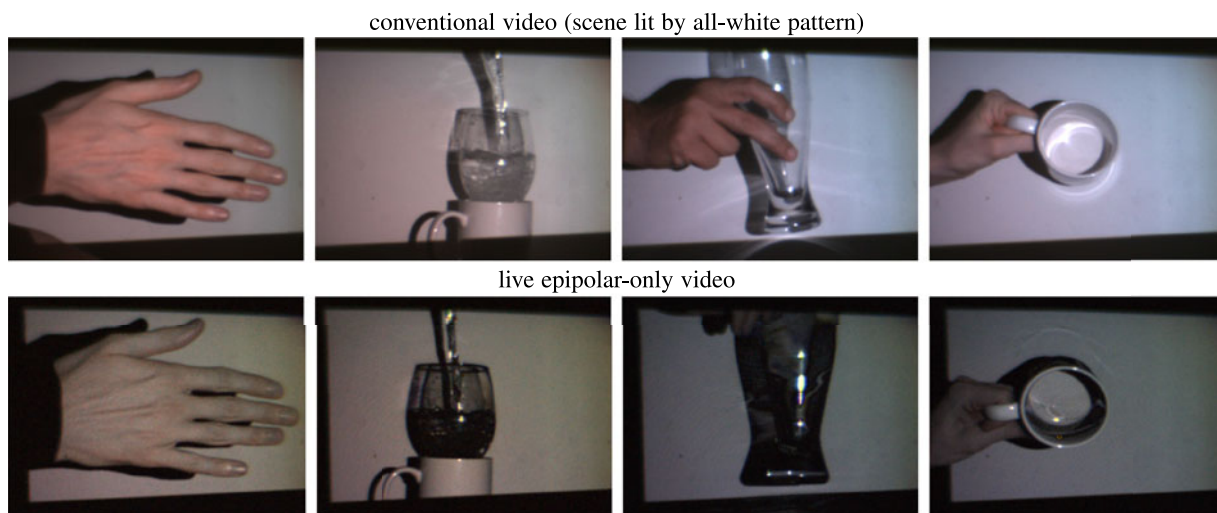


Fig. 13. Frames from conventional and epipolar-only video, corresponding to the live indirect video frames of Fig. 1. From left-to-right: (1) A hand; note the significant difference in apparent color of the hand in the indirect-only and epipolar-only components, due to sub-surface absorption and direct surface reflection, respectively. (2) Pouring water into a glass demonstrates our ability to successfully image highly-complex, time-varying phenomena. The water appears dark in the epipolar-only image because the light refracted by the water does not satisfy epipolar constraints. (3) This beer glass appears essentially opaque in the epipolar-only component because the light transmitted through the glass undergoes refraction, yielding non-linear paths that almost never lie on a single epipolar plane. (4) A mug. Artifacts appear on the white background, because the mug moved quickly during acquisition and the frame-differencing we do for epipolar-only imaging caused ghosting.

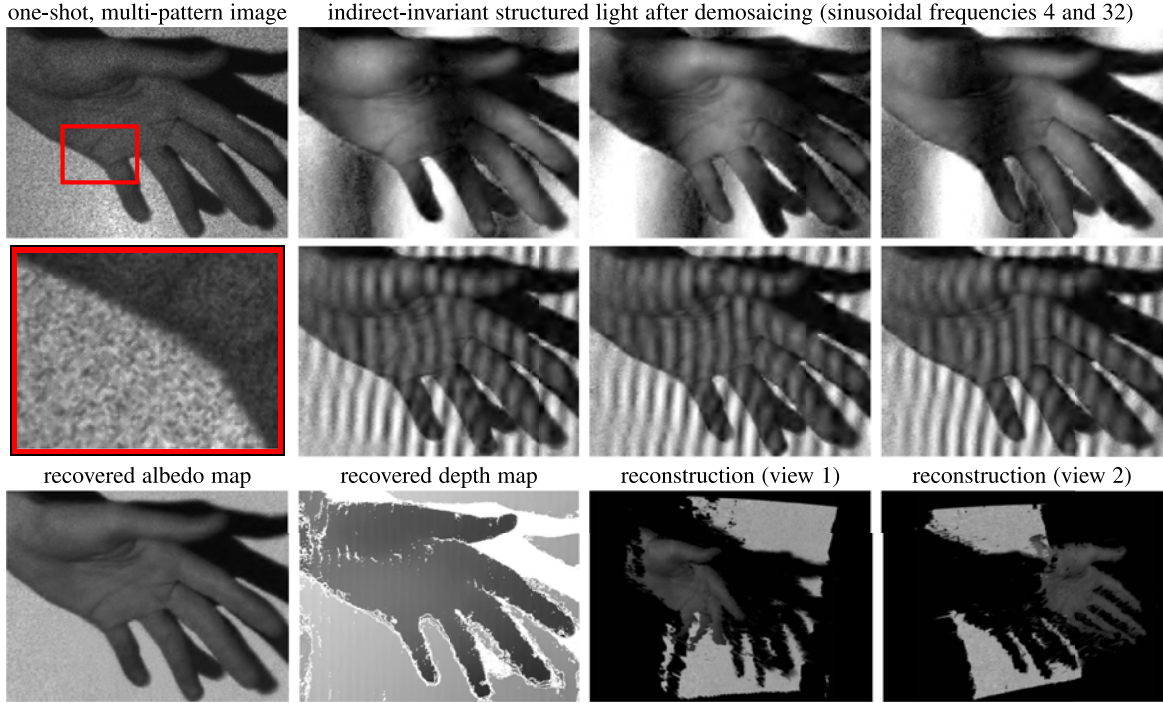


Fig. 14. Reconstructing dense depth and albedo from a video (frame 131 of 169) of a moving hand, recorded live using one-shot, indirect-invariant, multi-pattern imaging. Our demosaicing algorithm recovers 6 full-resolution indirect-invariant structured light images of the hand, for six sinusoidal patterns. These images yield albedo and depth maps on the bottom, and texture-mapped geometry shown from two viewpoints.

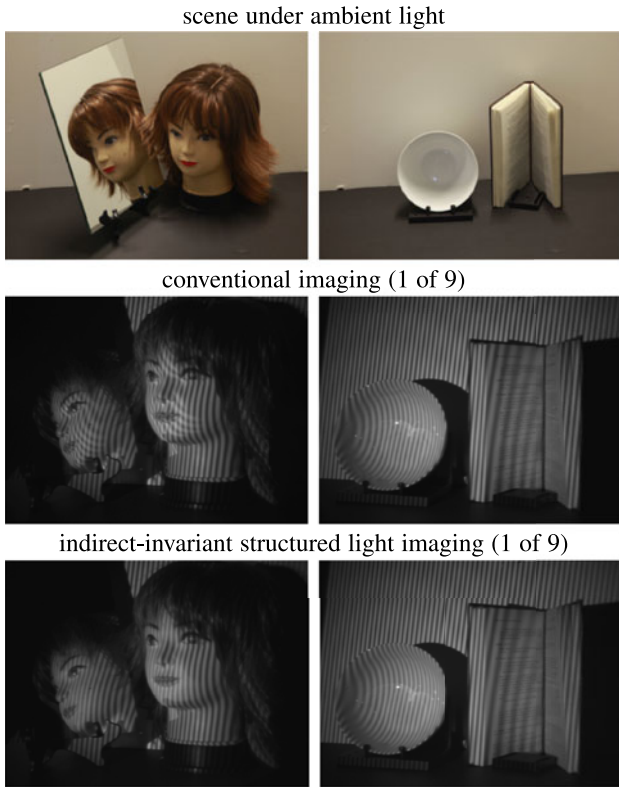


Fig. 15. We imaged the scene on the top row in two ways: (1) projecting 9 phase-shifted patterns directly onto it and (2) capturing indirect-invariant structured light images for the same patterns. Exposure time was held fixed, giving an SNR advantage to conventional projection which does not mask pixels. For the middle row, note that the conventional image contains “double fringes” from secondary reflections whereas the indirect-invariant one does not; this “double fringe” effect occurs because of the interference between the phase-shifted pattern transmitted through the direct channel and the same pattern specularly reflected by the mirror.

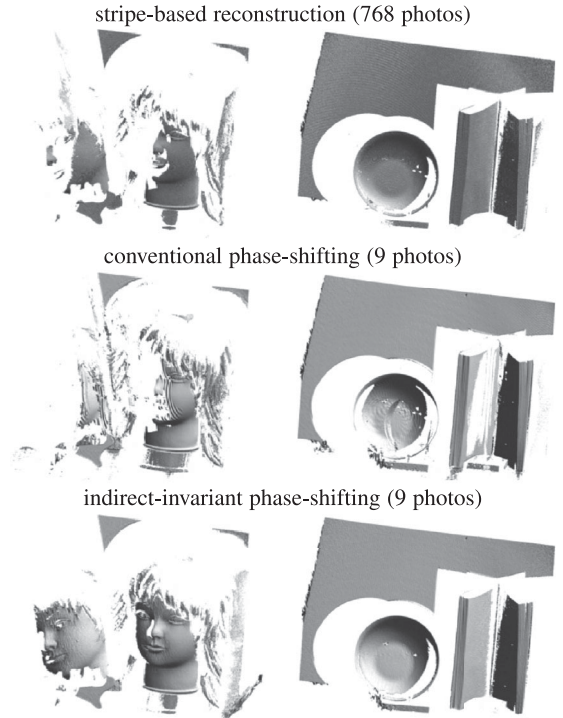


Fig. 16. 3D results for the scenes shown in Fig. 15. *Row 1:* Reconstruction results obtained by sweeping a vertical stripe across the scene, as done in conventional triangulation-based 3D laser scanning (768 images total). *Row 2:* The conventional phase-shift reconstruction procedure takes 9 input images acquired by conventional projection of phase-shifted patterns to compute raw 3D points. *Row 3:* Our indirect-invariant phase-shift results use the same 9 patterns and reconstruction algorithm as in conventional phase-shifting, combined with indirect-invariant structured light imaging. The phase-shift reconstruction algorithm fails catastrophically for the conventionally-acquired images, whereas with SLT imaging it is able to reconstruct even the hidden side of the face, from the mirror’s indirect view.

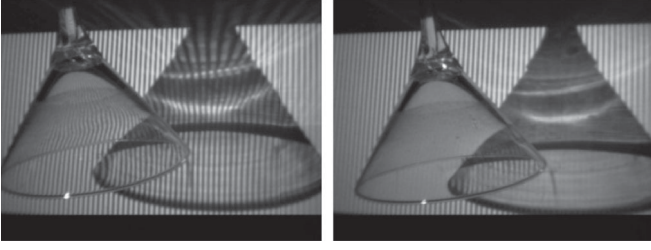


Fig. 17. Conventional (left) and indirect-invariant (right) structured light video frames of a martini glass lit by a high-frequency pattern. Note that the pattern appears to have disappeared from regions where light undergoes refraction, with the exception of degenerate regions where light paths are doubly refracted.



Fig. 18. Snapshots from a conventional video (left) and raw live indirect video (right) of a tight-fitting latex glove. Note the seemingly transparent nature of the glove in indirect-only mode. The hand is dark only in regions where the glove makes contact with the skin; otherwise, light scatters within cavities between the glove and skin.

apply the algorithm outlined in Section 5 independently to each frame of the raw live video stream.

7 EXPERIMENTAL RESULTS

Indirect-only and epipolar-only imaging. We used our low-speed, low-cost prototype with a total of $T = 96$ patterns/masks per frame for indirect-only and epipolar-only imaging. For calibration, we computed the epipolar geometry between the two DMDs by first relating them to the image plane. Overall resolution was equal to the resolution of our DMDs, *i.e.*, 608×684 . See Figs. 1, 12, 13, and 18 for examples of indirect- and epipolar-only images.

Indirect-invariant structured light. We used high-end DMDs and a monochrome camera to capture $S = 9$ indirect-invariant structured light frames (Fig. 15) with $T = 800$ patterns/masks per frame. We then supply these frames as input to a reconstruction algorithm to compute 3D shape (Fig. 16), and compare our approach to conventional approaches. The effective DMD resolution was approximately 484×364 . The scenes occupied a 40^3cm^3 volume about 70 cm away from the camera. To show the effectiveness of SLT imaging, we chose the most basic pattern and technique—phase-shifting with nine sinusoids total, at frequencies 1, 8 and 64. We also demonstrate an example of indirect-invariant structured light on a refractive object in Fig. 17.

Dense depth and albedo from one shot. We used $S = 6$ sinusoids at frequencies 4 and 32 for the experiment in Fig. 14, and a random assignment of pixels to sinusoids, rather than the regular assignment illustrated in Fig. 8. We recorded multi-pattern, indirect-invariant video at 28 fps and reconstructed each frame independently by (1) solving for the six demosaiced patterns using SPGL1 [28] for optimization and the JPEG2000 wavelet basis, and (2) using them to get per-pixel depth and albedo. Our reconstruction procedure took, as input, a raw frame cropped to a $1,024 \times 1,024$ region of interest (with an effective DMD resolution of 256×256), and recovered the depth and albedo in under 8 minutes per frame on an Apple iMac with a 2.8 GHz Intel Core i7 processor and 16 GB of memory.

8 CONCLUDING REMARKS

We believe that optical-domain processing—and SLT imaging in particular—offers a powerful new way to analyze the appearance of complex scenes, and to boost

the abilities of existing reconstruction algorithms. Although our focus was mainly on monochromatic light and conventional cameras, SLT imaging depends on neither; integrating this framework with other imaging dimensions (polarization, wavelength, time, etc.) is a promising direction. Last but not least, although our prototypes rely on DMD masks and several optical components, these would be rendered unnecessary if per-pixel processing was implemented directly on the sensor [29], [30]. We are looking forward to the wide availability of such technologies.

APPENDIX A

PROOFS OF PROPOSITIONS 1 AND 2

A.1 Proof of Proposition 1

Proposition 1. *If \mathbf{T}^{EI} and \mathbf{T}^{NE} are discretized forms of transport functions that are square-integrable and positive over the rectified projector and image planes, then*

$$\lim_{\epsilon \rightarrow 0} \frac{\mathbf{T}^{\text{EI}} \mathbf{p}}{\mathbf{T}^{\text{NE}} \mathbf{p}} = \mathbf{0} \quad (18)$$

where division is entrywise, $\mathbf{0}$ is a vector of zeros, and ϵ is the pixel size for discretization.

Proof sketch. We begin by identifying the rectified projector and image planes with the continuous domain $\mathcal{D} = [-1, 1] \times [-1, 1] \subset \mathbb{R}^2$. Let $p = (p_x, p_y)$ be a point on the projector plane and let $\mathbf{I}_\epsilon(p)$ be an indicator function over \mathcal{D} that specifies the spatial extent of the discrete epipolar line through the origin:

$$\mathbf{I}_\epsilon(p) = \begin{cases} 1 & \text{if } |p_x| \leq \frac{\epsilon}{2} \text{ and } |p_y| \leq 1 \\ 0 & \text{otherwise.} \end{cases} \quad (19)$$

In the continuous setting, light transport from the projector plane to the image plane is described by the *light transport equation* [31]. Given an image point $i \in \mathcal{D}$ on the epipolar line through the origin, this equation describes the total radiance transported to i from points on the projector plane:

$$\mathcal{I}(i) = \underbrace{\mathcal{T}(\hat{p}, i) \mathcal{P}(\hat{p})}_{\text{direct}} + \underbrace{\int_{\mathcal{D}-\{\hat{p}\}} \mathcal{T}(p, i) \mathcal{P}(p) dp}_{\text{indirect}}, \quad (20)$$

where \hat{p} is the projector point in stereo correspondence with image point i ; $\mathcal{P}(p)$ is the radiance along the ray through projector point p ; and $\mathcal{T}(p, i)$ is the *transport function* describing the proportion of radiance from p that gets transported to i .

We prove the continuous form of the ratio in Eq. (18) for an image point i ; this point is taken to be inside a discrete image pixel of dimension $\epsilon \times \epsilon$ on the epipolar line through the origin.

More specifically, we consider the epipolar indirect, total indirect, and non-epipolar indirect contributions at i :

$$\mathcal{I}^{\text{EI}}(i) = \int_{\mathcal{D}-\{\hat{p}\}} \mathbf{I}_\epsilon(p) \mathcal{T}(p, i) \mathcal{P}(p) dp \quad (21)$$

$$\mathcal{I}^{\text{I}}(i) = \int_{\mathcal{D}-\{\hat{p}\}} \mathcal{T}(p, i) \mathcal{P}(p) dp \quad (22)$$

$$\mathcal{I}^{\text{NE}}(i) = \mathcal{I}^{\text{I}}(i) - \mathcal{I}^{\text{EI}}(i). \quad (23)$$

We now show that for any $\delta > 0$, there is an $\epsilon > 0$ such that

$$\left| \frac{\mathcal{I}^{\text{EI}}(i)}{\mathcal{I}^{\text{NE}}(i)} \right| < \delta. \quad (24)$$

Since $\mathcal{T}()$ is square-integrable for domain $\mathcal{D} - \{\hat{p}\}$, we can apply the Cauchy-Schwarz inequality [32] to Eq. (21) to get an upper bound on the epipolar indirect contributions:

$$\begin{aligned} \mathcal{I}^{\text{EI}}(i) &\leq \left\{ \int_{\mathcal{D}-\{\hat{p}\}} \mathbf{I}_\epsilon(p) dp \right\}^{\frac{1}{2}} \left\{ \int_{\mathcal{D}-\{\hat{p}\}} [\mathcal{T}(p, i) \mathcal{P}(p)]^2 dp \right\}^{\frac{1}{2}} \\ &= (2\epsilon)^{\frac{1}{2}} \left\{ \int_{\mathcal{D}-\{\hat{p}\}} [\mathcal{T}(p, i) \mathcal{P}(p)]^2 dp \right\}^{\frac{1}{2}}. \end{aligned} \quad (25)$$

By combining Eqs. (22), (23) and (25) we also get a lower bound on the non-epipolar contributions:

$$\begin{aligned} \mathcal{I}^{\text{NE}}(i) &\geq \int_{\mathcal{D}-\{\hat{p}\}} \mathcal{T}(p, i) \mathcal{P}(p) dp \\ &- (2\epsilon)^{\frac{1}{2}} \left\{ \int_{\mathcal{D}-\{\hat{p}\}} [\mathcal{T}(p, i) \mathcal{P}(p)]^2 dp \right\}^{\frac{1}{2}}. \end{aligned} \quad (26)$$

Equation (24) now follows by choosing ϵ to be

$$\epsilon = \frac{1}{2} \left(\frac{\delta}{2 + \delta} \right)^2 \frac{\left\{ \int_{\mathcal{D}-\{\hat{p}\}} \mathcal{T}(p, i) \mathcal{P}(p) dp \right\}^2}{\int_{\mathcal{D}-\{\hat{p}\}} [\mathcal{T}(p, i) \mathcal{P}(p)]^2 dp}. \quad (27)$$

Substituting Eq. (27) into Eqs. (25) and (26) we get

$$\left| \frac{\mathcal{I}^{\text{EI}}(i)}{\mathcal{I}^{\text{NE}}(i)} \right| \leq \frac{\frac{\delta}{2+\delta}}{1 - \frac{\delta}{2+\delta}} = \frac{\delta}{2} < \delta. \quad (28)$$

□

A.2 Proof of Proposition 2

We prove Proposition 2 for scenes consisting of a finite collection of objects, each of which is an open set in \mathbb{R}^3 bounded by a smooth generic surface [33], [34].

Proposition 2. *Two generic n -bounce specular transport paths that originate from corresponding epipolar lines do not intersect for $n > 1$.*

Proof. For simplicity, we reverse the direction of light travel through image pixels, treating the camera as a second projector that also sends light onto the scene.

Let $\mathcal{L}, \mathcal{L}'$ be a pair of corresponding epipolar lines on the (continuous) projector and image planes, respectively, and let $p \in \mathcal{L}$ and $i \in \mathcal{L}'$ be points on them.

Suppose that the light originating at p and i undergoes $n \geq 1$ consecutive specular bounces upon entering the scene. Furthermore, suppose that the associated transport paths are generic, i.e., they remain stable under infinitesimal perturbations of the scene's surfaces and of the points p and i . To prove the proposition, we show that the following cannot hold simultaneously:

- (1) the transport paths through p and i intersect at their $(n+1)$ th bounce, i.e., their $(n+1)$ th bounce occurs at the same surface point in the scene; and
- (2) this intersection is generic, i.e., it occurs for all points p, i in an open interval $\mathcal{Q} \subset \mathcal{L}$ and $\mathcal{Q}' \subset \mathcal{L}'$, respectively.

In particular, let $l_n(p)$ be the 3D ray that light follows after n specular bounces from projector point p . Similarly, let $l'_n(i)$ be the corresponding 3D ray for image point i . Since the transport paths through p and i are generic, the mappings $p \mapsto l_n(p)$ and $i \mapsto l'_n(i)$ are smooth functions for some open neighborhood $\mathcal{Q} \subset \mathcal{L}$ and $\mathcal{Q}' \subset \mathcal{L}'$ of p and i , respectively. These mappings define a pair of ruled surfaces in \mathbb{R}^3 : intuitively, as point p ranges over \mathcal{Q} , the 3D ray $l_n(p)$ twists and translates in space, tracing a ruled surface.

For transport paths through p and i to have their $(n+1)$ th bounce in common, three surfaces must meet at a point: ruled surface $l_n(\mathcal{Q})$, ruled surface $l'_n(\mathcal{Q}')$, and a surface in the scene. This, however, is *not* a generic condition because surfaces $l_n(\mathcal{Q})$ and $l'_n(\mathcal{Q}')$ transversally intersect along a curve and this curve will transversally intersect the scene's surfaces at isolated points [34]. □

APPENDIX B

EXPANDED DERIVATIONS OF SELECTED EQUATIONS

B.1 Derivation of Eq. (11)

Combining Eqs. (6) and (9) we have

$$\mathbf{i}_e = \frac{1}{T} \sum_{t=1}^T \sum_{f=1}^E \overline{\mathbf{q}_e(t)} \circ [\mathbf{T}_{ef} \mathbf{q}_f(t)], \quad (29)$$

where the $1/T$ factor captures the fact that each term in the sum is allocated $1/T$ of the total exposure time. We now split the sum into its epipolar and non-epipolar terms

$$\begin{aligned} \mathbf{i}_e &= \frac{1}{T} \sum_{t=1}^T \left\{ \overline{\mathbf{q}_e(t)} \circ [\mathbf{T}_{ee} \mathbf{q}_e(t)] \right. \\ &\quad \left. + \sum_{\substack{f=1 \\ f \neq e}}^E \overline{\mathbf{q}_e(t)} \circ [\mathbf{T}_{ef} \mathbf{q}_f(t)] \right\}, \end{aligned} \quad (30)$$

and observe that the first term is always a vector of zeros. So,

$$\mathbf{i}_e = \frac{1}{T} \sum_{t=1}^T \sum_{\substack{f=1 \\ f \neq e}}^E \overline{\mathbf{q}_e(t)} \circ [\mathbf{T}_{ef} \mathbf{q}_f(t)]. \quad (31)$$

Letting $T \rightarrow \infty$ and applying the Central Limit Theorem to Eq. (31) we get the expected image $\mathcal{E}[\mathbf{i}_e]$ for epipolar line e :

$$\mathcal{E}[\mathbf{i}_e] = \mathcal{E} \left[\sum_{\substack{f=1 \\ f \neq e}}^E \overline{\mathbf{q}_e} \circ [\mathbf{T}_{ef} \mathbf{q}_f] \right] \quad (32)$$

$$= \mathcal{E}[\overline{\mathbf{q}_e}] \circ \sum_{\substack{f=1 \\ f \neq e}}^E \mathbf{T}_{ef} \mathcal{E}[\mathbf{q}_f] \quad (33)$$

$$= 0.25 \sum_{\substack{f=1 \\ f \neq e}}^E \mathbf{T}_{ef} \mathbf{1}, \quad (34)$$

where Eq. (34) follows from the fact that epipolar lines e and f are distinct and thus their corresponding random vectors \mathbf{q}_e and \mathbf{q}_f are independent.

B.2 Derivation of Eq. (15)

Combining Eqs. (6) and (9) for the indirect-invariant mask and pattern we have:

$$\mathbf{i}_e = \frac{1}{T} \sum_{t=1}^T \sum_{f=1}^E \mathbf{m}_e(t) \circ \left\{ \mathbf{T}_{ef} [\mathbf{m}_f(t) \circ \mathbf{r}_f(t) + \overline{\mathbf{m}_f(t)} \circ \overline{\mathbf{r}_f(t)}] \right\}. \quad (35)$$

We split the sum into its epipolar and non-epipolar terms,

$$\begin{aligned} \mathbf{i}_e = \frac{1}{T} \sum_{t=1}^T & \left\{ \mathbf{m}_e(t) \circ \mathbf{T}_{ee} [\mathbf{m}_e(t) \circ \mathbf{r}_e(t)] \right. \\ & + \mathbf{m}_e(t) \circ \mathbf{T}_{ee} [\overline{\mathbf{m}_e(t)} \circ \overline{\mathbf{r}_e(t)}] \\ & + \sum_{\substack{f=1 \\ f \neq e}}^E \mathbf{m}_e(t) \circ \mathbf{T}_{ef} [\mathbf{m}_f(t) \circ \mathbf{r}_f(t)] \\ & \left. + \sum_{\substack{f=1 \\ f \neq e}}^E \mathbf{m}_e(t) \circ \mathbf{T}_{ef} [\overline{\mathbf{m}_f(t)} \circ \overline{\mathbf{r}_f(t)}] \right\}, \end{aligned} \quad (36)$$

and note that the second term of Eq. (36) is always a vector of zeros. Letting $T \rightarrow \infty$ and applying the Central Limit Theorem to Eq. (36) we get the expected image for epipolar line e :

$$\begin{aligned} \mathcal{E}[\mathbf{i}_e] = \mathcal{E} & \left[\mathbf{q}_e \circ [\mathbf{T}_{ee} (\mathbf{q}_e \circ \mathbf{r}_e)] \right. \\ & \left. + \sum_{\substack{f=1 \\ f \neq e}}^E \mathbf{q}_e \circ [\mathbf{T}_{ef} (\mathbf{q}_f \circ \mathbf{r}_f + \overline{\mathbf{q}_f} \circ \overline{\mathbf{r}_f})] \right]. \end{aligned} \quad (37)$$

Now, \mathbf{q}_e is a random binary vector whose probability of being either $\mathbf{1}$ or $\mathbf{0}$ is 0.5. Using this fact as well as \mathbf{q}_e 's independence from all other random vectors, the expectation in Eq. (37) becomes

$$\mathcal{E}[\mathbf{i}_e] = 0.5 \mathbf{T}_{ee} \mathcal{E}[\mathbf{r}_e] + 0.5 \sum_{\substack{f=1 \\ f \neq e}}^E \mathbf{T}_{ef} \mathcal{E}[\mathbf{q}_f \circ \mathbf{r}_f + \overline{\mathbf{q}_f} \circ \overline{\mathbf{r}_f}]. \quad (38)$$

Finally, using the definition of binary random vector \mathbf{r}_f in Eq. (14) the expectation becomes

$$\begin{aligned} \mathcal{E}[\mathbf{i}_e] = 0.5 \mathbf{T}_{ee} \mathbf{p}_e + 0.5 \sum_{\substack{f=1 \\ f \neq e}}^E \mathbf{T}_{ef} \{ \text{Prob}[\mathbf{q}_f = \mathbf{1}] \mathbf{p}_e \\ + \text{Prob}[\mathbf{q}_f = \mathbf{0}] (\mathbf{1} - \mathbf{p}_e) \} \end{aligned} \quad (39)$$

$$= 0.5 \mathbf{T}_{ee} \mathbf{p}_e + 0.25 \sum_{\substack{f=1 \\ f \neq e}}^E \mathbf{T}_{ef} \mathbf{1}. \quad (40)$$

ACKNOWLEDGMENTS

We are grateful for the support of the Natural Sciences and Engineering Research Council of Canada under the PGS-D, RGPIN, RTI, SGP and GRAND NCE programs.

REFERENCES

- [1] P. Debevec, et al., "Acquiring the reflectance field of a human face," in *Proc. 27th Annu. Conf. Comput. Graph. Interactive Tech.*, 2000, pp. 145–156.
- [2] S. K. Nayar, et al., "Fast separation of direct and global components of a scene using high frequency illumination," in *Proc. SIGGRAPH*, 2006, pp. 935–944.
- [3] G. Godin, M. Rioux, J. Beraldin, and M. Levoy, "An assessment of laser range measurement on marble surfaces," presented at the 5th Conf. Optical 3D Measurement Techniques, Vienna, Austria, 2001.
- [4] B. Curless and M. Levoy, "Better optical triangulation through spacetime analysis," in *Proc. 5th Int. Conf. Comput. Vision*, 1995, pp. 987–994.
- [5] J. Salvi, S. Fernandez, T. Pribanic, and X. Llado, "A state of the art in structured light patterns for surface profilometry," *Pattern Recogn.*, vol. 43, no. 8, pp. 2666–2680, 2010.
- [6] M. Gupta and S. Nayar, "Micro phase shifting," in *Proc. IEEE Comput. Soc. Conf. Comput. Vision Pattern Recog.*, 2012, pp. 1–8.
- [7] S. K. Nayar, K. Ikeuchi, and T. Kanade, "Shape from inter-reflections," *Int. J. Comput. Vision*, vol. 6, no. 3, pp. 173–195, 1991.
- [8] R. Hartley and A. Zisserman, *Multiple View Geometry in Computer Vision*. Cambridge, U.K.: Cambridge Univ. Press, Dec. 2000.
- [9] M. O'Toole, J. Mather, and K. N. Kutulakos. (2014). Supplementary materials [Online]. Available: <http://www.dgp.toronto.edu/~motoole/slt>
- [10] M. O'Toole, R. Raskar, and K. N. Kutulakos, "Primal-dual coding to probe light transport," in *Proc. SIGGRAPH*, 2012, pp. 1–11.
- [11] A. Velten, et al., "Femto-photography: capturing and visualizing the propagation of light," in *Proc. SIGGRAPH*, 2013, pp. 1–8.
- [12] S. Achar, S. T. Nuske, and S. G. Narasimhan, "Compensating for motion during direct-global separation," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2013, pp. 1481–1488.
- [13] M. Gupta, et al., "Structured light 3D scanning in the presence of global illumination," in *Proc. IEEE Conf. Comput. Vision Pattern Recog.*, 2011, pp. 713–720.
- [14] V. Couture, N. Martin, and S. Roy, "Unstructured light scanning to overcome interreflections," in *Proc. IEEE Int. Conf. Comput. Vision*, 2011, pp. 1895–1902.
- [15] T. Chen, H.-P. Seidel, and H. P. A. Lensch, "Modulated phase-shifting for 3D scanning," in *Proc. IEEE Conf. Comput. Vision Pattern Recog.*, 2008, pp. 1–8.
- [16] L. Zhang, B. Curless, and S. M. Seitz, "Rapid shape acquisition using color structured light and multi-pass dynamic programming," in *Proc. 1st Int. Symp. 3D Data Process., Vis. Transmiss.*, 2002, pp. 24–36.
- [17] H. Kawasaki, et al., "Dynamic scene shape reconstruction using a single structured light pattern," in *Proc. IEEE Conf. Comput. Vision Pattern Recog.*, 2008, pp. 1–8.

- [18] C. Hernandez, et al., "Non-rigid photometric stereo with colored lights," in *Proc. IEEE 11th Int. Conf. Comput. Vision*, 2007, pp. 1–8.
- [19] G. Fyffe, X. Yu, and P. Debevec, "Single-shot photometric stereo by spectral multiplexing," in *Proc. IEEE Int. Conf. Computat. Photography*, 2011, pp. 1–6.
- [20] R. Heintzmann, et al., "A dual path programmable array microscope (PAM)," *J. Microscopy*, vol. 204, pp. 119–135, 2001.
- [21] J. Mertz, "Optical sectioning microscopy with planar or structured illumination," *Nature Methods*, vol. 8, no. 10, pp. 811–819, 2011.
- [22] T. Wilson, R. Juskaitis, and M. Neil, "Confocal microscopy by aperture correlation," *Opt. Lett.*, vol. 21, no. 3, p. 1879, 1996.
- [23] R. Ng, R. Ramamoorthi, and P. Hanrahan, "All-frequency shadows using non-linear wavelet lighting approximation," in *Proc. SIGGRAPH*, 2003, pp. 376–381.
- [24] H. Jensen, et al., "A practical model for subsurface light transport," in *Proc. 28th Annu. Conf. Comput. Graph. Interactive Tech.*, 2001, pp. 511–518.
- [25] J. Zhong, "Binary ranks and binary factorizations of nonnegative integer matrices," *Electron. J. Linear Algebra*, vol. 23, pp. 540–552, 2012.
- [26] D. Reddy, A. Veeraraghavan, and R. Chellappa, "P2C2: Programmable pixel compressive camera for high speed imaging," in *Proc. IEEE Conf. Comput. Vision Pattern Recog.*, 2011, pp. 329–336.
- [27] Y. Hitomi, J. Gu, M. Gupta, T. Mitsunaga, and S. K. Nayar, "Video from a single coded exposure photograph using a learned over-complete dictionary," in *Proc. Int. Conf. Comput. Vision*, 2011, pp. 287–294.
- [28] E. van den Berg and M. P. Friedlander. (2007). SPGL1: A solver for large-scale sparse reconstruction, [Online]. Available: <https://www.math.ucdavis.edu/~mpf/spgl1/>
- [29] M. W. Kelly and M. H. Blackwell, "Digital-pixel FPAs enhance infrared imaging capabilities," *Laser Focus World*, vol. 49, no. 1, pp. 90, 2013.
- [30] G. Wan, et al., "CMOS image sensors with multi-bucket pixels for computational photography," *IEEE J. Solid State Circuits*, vol. 47, no. 4, pp. 1031–1042, 2012.
- [31] J. T. Kajiya, "The rendering equation," in *Proc. 13th Annu. Conf. Computer Graph. Interactive Tech.*, 1986, pp. 143–150.
- [32] W. Rudin, *Principles of Mathematical Analysis*. New York, NY, USA: McGraw-Hill, 1976.
- [33] J. J. Koenderink, *Solid Shape*. Cambridge, U.K.: Cambridge Univ. Press, 1990.
- [34] V. Guillemin and A. Pollack, *Differential Topology*. Englewood Cliffs, NJ, USA: Prentice-Hall, 1974.



Matthew O'Toole received the BS degree in computer science and mathematics from the University of British Columbia, Vancouver, BC, Canada, in 2007, and the MS degree in computer science from the University of Toronto, Toronto, ON, Canada, in 2009. In 2011, he was a visiting student at the Massachusetts Institute of Technology. He is currently working toward the PhD degree in computer science at the University of Toronto, supported in part by the National Sciences and Engineering Research Council of Canada. His work received a David Marr Prize Honorable Mention award at ICCV 2007, and the Best Paper Honorable Mention award at CVPR 2014.



John Mather received the BS degree in computer science from the University of Toronto, Toronto, ON, Canada, in 2015. From 2012 to 2015, he was a research assistant at the Dynamic Graphics Project at the University of Toronto. He is a recipient of the Best Paper Honorable Mention award at CVPR 2014. He is currently a software developer for the animation and visual effects industry.



Kiriakos N. Kutulakos received the BS degree from the University of Crete, Greece, in 1988, and the PhD degree from the University of Wisconsin-Madison, Wisconsin, in 1994, both in Computer Science. He is a professor of computer science at the University of Toronto, Toronto, ON, Canada. In addition to the University of Toronto, he has held appointments at the University of Rochester (1995–2001) and Microsoft Research Asia (2004–05 and 2011–12). He is the recipient of an Alfred P. Sloan Fellowship, an Ontario Premier's Research Excellence Award, a David Marr Prize in 1999, a David Marr Prize Honorable Mention in 2005, and three other paper awards (CVPR 1994, ECCV 2006, CVPR 2014). He was an associate editor of the *IEEE Transactions on Pattern Analysis and Machine Intelligence* from 2005 to 2010 and served as Program Co-Chair of CVPR 2003, ICCP 2010, and ICCV 2013. He is a member of the IEEE.

► For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/publications/dlib.