

From Shading to Local Shape

Ying Xiong, Ayan Chakrabarti, Ronen Basri, Steven J. Gortler, David W. Jacobs, and Todd Zickler

Abstract—We develop a framework for extracting a concise representation of the shape information available from diffuse shading in a small image patch. This produces a mid-level scene descriptor, comprised of local shape distributions that are inferred separately at every image patch across multiple scales. The framework is based on a quadratic representation of local shape that, in the absence of noise, has guarantees on recovering accurate local shape and lighting. And when noise is present, the inferred local shape distributions provide useful shape information without over-committing to any particular image explanation. These local shape distributions naturally encode the fact that some smooth diffuse regions are more informative than others, and they enable efficient and robust reconstruction of object-scale shape. Experimental results show that this approach to surface reconstruction compares well against the state-of-art on both synthetic images and captured photographs.

Index Terms—Shape from shading, local shape descriptors, statistical models, 3D reconstruction

1 INTRODUCTION

RECOVERING shape from diffuse shading is point-wise ambiguous because each surface normal can lie anywhere on a cone of directions. Surface normals are uniquely determined only where they align with the light direction which, at best, occurs at only a handful of singular points. A common strategy for reducing the ambiguity is to pursue global reconstructions of large, pre-segmented regions, with the hope that many point-wise ambiguities will collaboratively resolve, or that shape information will successfully propagate from identifiable singular points and occluding contours.

Global strategies are difficult to apply in natural scenes because diffuse shading is typically intermixed with other phenomena such as texture, gloss, shadows, translucency, and mesostructure. Occluding contours and singular points are hard to detect in these scenes; and shading-based shape propagation breaks down unless occlusions, gloss, texture, etc. are somehow analyzed and removed by additional visual reasoning. Moreover, most global strategies do not provide spatial uncertainty information to accompany their output reconstructions, and this limits their use in providing feedback to improve top-down scene analysis, or in co-computing with other necessary bottom-up processes that perform complimentary analysis of other phenomena.

As illustrated in Fig. 1, this paper develops a framework for leveraging diffuse shading more broadly and robustly by developing a richer description of what it says *locally* about shape. We show that point-wise ambiguity can be

systematically reduced by jointly analyzing intensities in small image patches, and that some of these patches are inherently more informative than others. Accordingly, we develop an algorithm that produces for any image patch a concise distribution of surface patches that are likely to have created it. We propose these dense, local shape distributions as a new mid-level scene representation that provides useful local shape information without over-committing to any particular image explanation. Finally, we show how these local shape distributions can be combined to recover global object-scale shape.

Our framework is developed in three parts:

- 1) *Local uniqueness.* We provide uniqueness results for jointly recovering shape and lighting from a small image patch. By considering a world in which the shape of each small surface patch is exactly the graph of a quadratic function, we prove two generic facts: i) when the light direction is known, quadratic shape is uniquely determined; and ii) when the light is unknown, it is determined up to a four-way choice. We also catalog the degenerate cases, which correspond to special shapes, or conspiracies between the light and shape. These results are of direct interest to those studying the mathematics of shape from shading.
- 2) *Local shape distributions.* We introduce a computational process that takes an image patch at any scale and produces a compact distribution of quadratic shapes that are likely to have produced it. At the core of this process is our observation that all likely shapes corresponding to a (noisy) image patch lie close to a one-dimensional manifold embedded in the five-dimensional space of quadratic shapes. This part of the paper is of broad interest because these local, multi-scale shape distributions may be useful as intermediate scene descriptors for various visual tasks.
- 3) *Reconstruction.* We present a simple and effective bottom-up reconstruction system for inferring object-scale shape from a single image of a

• Y. Xiong, A. Chakrabarti, S.J. Gortler, and T. Zickler are with the Harvard School of Engineering and Applied Sciences, Cambridge, MA 02138. E-mail: {yxiong, zickler}@seas.harvard.edu, ayanc@eecs.harvard.edu, sjg@cs.harvard.edu.

• R. Basri is with the Weizmann Institute of Science, Rehovot 76100, Israel. E-mail: ronien.basri@weizmann.ac.il.

• D.W. Jacobs is with the Department of Computer Science, University of Maryland, College Park, MD 20742. E-mail: djacobs@cs.umd.edu.

Manuscript received 12 Oct. 2013; revised 27 Mar. 2014; accepted 19 May 2014. Date of publication 25 July 2014; date of current version 5 Dec. 2014.

Recommended for acceptance by C.-K. Tang.

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.2014.2343211

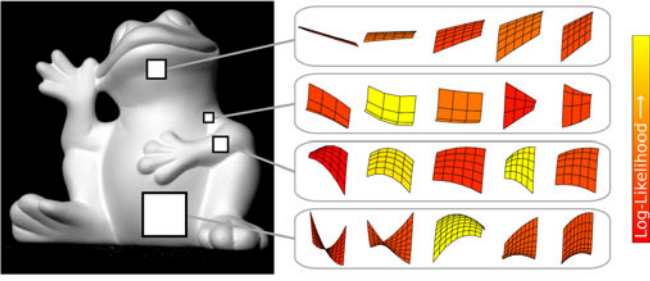


Fig. 1. We infer from a Lambertian image patch a concise representation for the distribution of quadratic surfaces that are likely to have produced it. These distributions naturally encode different amounts of shape information based on what is locally available in the patch, and can be unimodal (row 2 & 4), multi-modal (row 3), or near-uniform (row 1). This inference is done across multiple scales.

predominantly textureless and diffusely-shaded surface. This reconstruction system uses as input our local shape distributions inferred from dense, overlapping patches at multiple scales. It is conceptually simple, computationally parallelizable, and robust to non-idealities like shadows, texture, highlights, etc., with reconstruction accuracy that compares well to the state of the art. This system is of direct interest to those studying algorithms for traditional shape from shading, and it is structured in a modular way that provides a step toward co-computation with other reconstructive processes that also analyze other phenomena.

These three parts are tightly bound together. The uniqueness results in Section 3 show that the quadratic model is a particularly convenient representation for small surface patches. In the absence of noise, both shape and lighting are locally revealed, and local shape is generally unique when lighting is known. Building on this, Section 4 examines how uniqueness breaks down in the presence of noise. While very different quadratic shapes can produce equally-likely local intensity patterns, we find that all highly-likely shapes lie close to a one-dimensional sub-manifold. Then, Section 5 shows how to infer a dense set of sample shapes along this sub-manifold, thereby taking an image patch and producing a one-dimensional shape distribution. Finally, Section 6 shows how these multi-scale local distributions can enable robust global reconstruction of shape, by naturally encoding the fact that some smooth diffuse regions are more informative than others.

The project page associated with this paper [1] provides separate implementations of our algorithms for inferring local distributions (Section 5) and global shape (Section 6). These are highly parallelized and can be executed on a single machine, a local cluster of machines, or a cluster from a standard utility computing service.

2 RELATED WORK

Background on shape inference from diffuse shading can be found in several reviews and surveys [2], [3], [4]. An important question is whether shape is uniquely determined by a noiseless image, which has been addressed by a variety of PDE-based formulations. For example, Oliensis considered

C^2 surfaces and showed that shape can be uniquely determined for the entire image by singular points and occluding boundaries together [5], and in many parts of the image by singular points alone [6]. For the more general class of C^1 surfaces, Prados and Faugeras [7] employed a smoothness constraint to prove uniqueness properties in a more general perspective setup [8], [9] given appropriate boundary conditions. In this paper, we use a more restrictive local surface model but prove local uniqueness without any boundary conditions or knowledge of singular points. This generalizes previous studies of local uniqueness, which have considered locally-spherical [10] and fronto-parallel [11] surfaces.

Global uniqueness analyses have inspired global propagation and energy-based methods for global shape inference (e.g. [2], [12], [13]), some of which rely on identifying occluding boundaries and/or singular points. While most methods do not typically provide any measurement of uncertainty in their output, progress toward representing shape ambiguity was made by Ecker and Jepson [14], who use a polynomial formulation of global shape from shading to numerically generate distinct global surfaces that are equally close to an input image. In this paper, we study uniqueness and uncertainty at the local level, and infer distributions over candidate local shapes.

Our work is related to patch-based approaches that use synthetically-generated reference databases. The idea there is to reconstruct depth (or other scene properties [15]) by synthesizing a database of aligned image and depth-map pairs, and then finding and stitching together depth patches from this database to match the input image and be spatially consistent. Hassner and Basri [16] obtain plausible results this way when the input image and the database are of similar object categories, and Huang et al. [17] pursue a similar goal for textureless objects using a database of rendered Lambertian spheres. Cole et al. [18] focus on patches located at detected key-points near an object's occlusion boundaries, combining shading and contour cues. We also describe global shape as a mosaic of per-patch depth primitives, but instead of relying on primitives from a pre-chosen set of 3D models, we consider a continuous five-parameter family of depth primitives corresponding to graphs of quadratic functions at multiple scales.

One of our main motivations is the long-term goal of enabling better co-computation with other bottom-up and top-down visual processes, and by providing useful local shape information without choosing any single image interpretation, our distributions are consistent with Marr's principle of least commitment [19]. We focus on diffuse shading on textureless surfaces, leaving for future work the task of merging with bottom-up processes for other cues like occluding contours (e.g., [18], [20]), texture, gloss, and so on. Our belief that this will be useful is bolstered by promising results achieved by recent global approaches to such combined reasoning [21].

In independent work, Kunsberg and Zucker [22] have recently derived local uniqueness results that are related to, and consistent with, our results in Section 3. Their elegant analysis, which uses differential geometry and applies to continuous images, is complimentary to the discrete and

algebraic approach employed in this paper. Kunsberg and Zucker also observe that the analysis of shading in patches instead of at isolated points is consistent with early processing in the visual cortex, and they discuss the possibility of local shading distributions being computed there. Indeed, the notion of such distributions is compatible with evidence that humans perceive shape in some diffuse regions more accurately than others [11].

3 QUADRATIC-PATCH SHAPE FROM SHADING

We begin by analyzing the ability to uniquely determine the shape and lighting of a local patch from a Lambertian shading image in the absence of noise. The key assumption in our analysis is that depth of the patch can be *exactly* expressed as the graph of a quadratic function. While subsequent sections consider deviations from this idealized setting, the following analysis characterizes the inherent ambiguity under a local quadratic patch model.

We model the depth $z(x, y)$ of a local surface patch as a quadratic function defined by coefficient vector $a \in \mathbb{R}^5$ up to a constant offset:¹

$$z(x, y; a) = a_1 x^2 + a_2 y^2 + a_3 xy + a_4 x + a_5 y. \quad (1)$$

In matrix form, this is $z = [x, y]H[x, y]^T + J[x, y]^T$ with

$$H = \begin{bmatrix} a_1 & a_3/2 \\ a_3/2 & a_2 \end{bmatrix} \quad (2)$$

the Hessian matrix and $J = [a_4, a_5]$ the Jacobian of the depth function. The un-normalized surface normal to this patch at each location (x, y) is then given by

$$n(x, y; a) = [n_x(x, y; a), n_y(x, y; a), 1]^T, \quad (3)$$

where

$$n_x(x, y; a) \triangleq -\frac{\partial z}{\partial x} = -2a_1 x - a_3 y - a_4, \quad (4)$$

$$n_y(x, y; a) \triangleq -\frac{\partial z}{\partial y} = -2a_2 y - a_3 x - a_5. \quad (5)$$

In matrix form, this is $n(x, y; a) = A[x, y, 1]^T$ with

$$A \triangleq \begin{bmatrix} -2a_1 & -a_3 & -a_4 \\ -a_3 & -2a_2 & -a_5 \\ 0 & 0 & 1 \end{bmatrix} \quad (6)$$

the shape matrix corresponding to quadratic shape a .

The intensity $I(x, y; a)$ of this patch, observed from viewing direction $v = [0, 0, 1]^T$ under a directional light source $l = [l_x, l_y, l_z]^T$, is

$$I(x, y; a) = \frac{l^T n(x, y; a)}{\|n(x, y; a)\|}, \quad (7)$$

assuming spatially-uniform Lambertian reflectance and that no part of the patch is in shadow, i.e., $l^T n(x, y) > 0, \forall (x, y)$.

1. Local shading for the special case $a_4 = a_5 = 0$ is described in [11], and a more restrictive, locally-spherical model $z(x, y) = \sqrt{r^2 - x^2 - y^2}$ is analyzed in [10].

Here, the magnitude $\|l\|$ of the light vector represents the product of the surface albedo and the light strength, and it is not assumed to be equal to one. Re-arranging, the intensity I at each point (x, y) induces a quadratic constraint on its surface normal [14]:

$$I^2 n^T n = n^T l l^T n \Rightarrow n^T (l l^T - I^2 \mathbb{I}_{3 \times 3}) n = 0, \quad (8)$$

where $\mathbb{I}_{3 \times 3}$ is the identity matrix. This further induces a related constraint on shape parameters a :

$$[a^T \ 1](D^T(l l^T - I^2 \mathbb{I}_{3 \times 3})D) \begin{bmatrix} a \\ 1 \end{bmatrix} = 0, \quad (9)$$

where we use the matrix $D \in \mathbb{R}^{3 \times 6}$ to re-write the relationship between n and a in (3)-(5) as $n = D[a^T \ 1]^T$.

Every pixel (x, y) in an image patch gives one such constraint on shape parameters a , and shape from shading for quadratic patches rests on solving this system of polynomial equations. Our immediate goal is to determine whether the shape a and lighting l can be uniquely determined from these local constraints.

3.1 Uniqueness of Simultaneous Shape and Light

We assume that the local patch is sufficiently large to contain a minimum number of *non-degenerate* pixel locations, where the condition for non-degeneracy is defined as follows:

Definition 1. For a patch $\Omega = \{(x_i, y_i)\}_{i=1}^N$, we define the matrix $V_\Omega \in \mathbb{R}^{N \times 15}$ such that each row v_i of V_Ω consists of all fourth-order and lower terms of x_i and y_i :

$$v_i = \begin{bmatrix} x_i^4, x_i^3 y_i, \dots, x_i^p y_i^q, \dots, x_i, y_i, 1 \end{bmatrix}_{p,q \geq 0, p+q \leq 4}. \quad (10)$$

A patch Ω is considered non-degenerate if the matrix V_Ω has rank 15.

Note that rectangular grids of pixels that are 5×5 or larger will be non-degenerate under the above definition.

Theorem 1. Given intensities $I(x, y)$ in an image patch Ω collected at a set of non-degenerate locations not in shadow, if any quadratic-patch/lighting pair (a, l) that satisfies the set of polynomial equations (9) has a surface Hessian with eigenvalues that are not equal in magnitude, then there are no more than four distinct surfaces that can create the same image. Each of these surfaces is associated with a unique lighting when the Hessian of any solution is non-singular, and a one-dimensional family of lighting vectors otherwise.

This theorem states that given measurements of intensity from a quadratic surface patch, there generically exists four physical explanations, each comprised of a shape a , a light direction $l/\|l\|$, and a scalar $\|l\|$ encoding the product of albedo and light strength.

Before proceeding to the proof, we introduce a lemma that relates to equations with ratios of quadratic terms. We define $\bar{x} \triangleq [x \ y \ 1]^T$, so that the normals are given by $n(x, y; a) = A\bar{x}$, and the intensity constraint (9) becomes

$$I_{\bar{x}}^2 = \left(\frac{l^T \bar{n}}{\|\bar{n}\|} \right)^2 = \frac{\bar{x}^T A^T l l^T A \bar{x}}{\bar{x}^T A^T A \bar{x}}. \quad (11)$$

Using this notation, we can state the following lemma, which is proven in the supplementary material:

Lemma 1. Let A and \tilde{A} correspond to two matrices of the form in (6), and l and \tilde{l} to two lighting vectors. If

$$\frac{\bar{x}^T A^T l^T A \bar{x}}{\bar{x}^T A^T A \bar{x}} = \frac{\bar{x}^T \tilde{A}^T \tilde{l}^T \tilde{A} \bar{x}}{\bar{x}^T \tilde{A}^T \tilde{A} \bar{x}}, \quad \forall \bar{x} \in \Omega, \quad (12)$$

and if $\text{Rank}(V_\Omega) = 15$, $\text{Rank}(A) \geq 2$, and $l^T A \bar{x} > 0, \forall \bar{x} \in \Omega$ (i.e., no point is in shadow), then

$$A^T l^T A = \tilde{A}^T \tilde{l}^T \tilde{A}, \quad A^T A = \tilde{A}^T \tilde{A}. \quad (13)$$

Moreover, if $\text{Rank}(A) = 2$, then $\text{Rank}(\tilde{A}) = 2$ and both A and \tilde{A} share a common null space.

Proof of Theorem 1. Suppose there exists a solution (a, l) that produces the observed set of intensities in the patch Ω , and the Hessian matrix of surface a has eigenvalues of un-equal magnitude. We will prove that if there exists another solution (\tilde{a}, \tilde{l}) , such that

$$\frac{\bar{x}^T A^T l^T A \bar{x}}{\bar{x}^T A^T A \bar{x}} = l_x^2 = \frac{\bar{x}^T \tilde{A}^T \tilde{l}^T \tilde{A} \bar{x}}{\bar{x}^T \tilde{A}^T \tilde{A} \bar{x}}, \quad \forall \bar{x} \in \Omega_i, \quad (14)$$

then \tilde{a} must be related to a in one of four specific ways.

Since a is not planar (otherwise the Hessian would have both eigenvalues equal to zero), the corresponding matrix A is at least rank 2, and we can apply Lemma 1:

$$\tilde{A}^T \tilde{l}^T \tilde{A} = A^T l^T A, \quad \tilde{A}^T \tilde{A} = A^T A. \quad (15)$$

We define a new matrix B satisfying $\tilde{A} = BA$. Specifically, when A is full rank we set $B = \tilde{A}A^{-1}$; and when $\text{Rank}(A) = 2$, we set $B = (\tilde{A} + vv^T)(A + vv^T)^{-1}$ with v a vector in the common null-space of A and \tilde{A} , i.e., $Av = \tilde{A}v = 0$. We will show that there are only four possibilities for the matrix B .

Note that A and \tilde{A} are affine matrices (last rows are both $[0, 0, 1]$). Moreover, in the rank 2 case, the last entry of v will be 0 and $A + vv^T$ will also be an affine matrix. Therefore, A^{-1} (if A is full rank) and $(A + vv^T)^{-1}$ (if A is rank 2) are affine. Hence, B is also an affine matrix:

$$B = \begin{bmatrix} b_{11} & b_{12} & b_{13} \\ b_{21} & b_{22} & b_{23} \\ 0 & 0 & 1 \end{bmatrix}. \quad (16)$$

From (15), we have $B^T B = \mathbb{I}_{3 \times 3}$, i.e.,

$$b_{13}^2 + b_{23}^2 + 1 = 1 \implies b_{13} = b_{23} = 0. \quad (17)$$

The orthogonality of B further restricts its top-left block to be either a 2D rotation matrix

$$B = \begin{bmatrix} \cos \varphi & -\sin \varphi & 0 \\ \sin \varphi & \cos \varphi & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad (18)$$

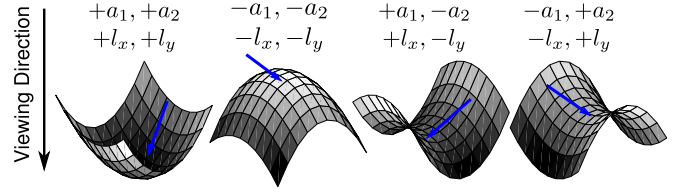


Fig. 2. Four quadratic-patch/lighting configurations that produce the same image (left is $a = [1, 1/2, 0, 0, 0]$, $l = [2/3, 1/3, 2/3]$). The lighting is shown as blue arrows. The left pair and right pair are each convex-concave.

or an “anti-rotation” matrix

$$B = \begin{bmatrix} \cos \varphi & \sin \varphi & 0 \\ \sin \varphi & -\cos \varphi & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad (19)$$

for $\varphi \in [-\pi, \pi]$.

From $\tilde{A} = BA$ and the fact that the (1,2)-entry and (2,1)-entry of \tilde{A} matrix should be the same (since $a_{12} = a_{21} = -a_3$, $\tilde{a}_{12} = \tilde{a}_{21} = -\tilde{a}_3$), we have

$$2a_1b_{21} + a_3b_{22} = a_3b_{11} + 2a_2b_{12}. \quad (20)$$

This implies that when B is of the form in (18)

$$(a_1 + a_2)\sin \varphi = 0, \quad (21)$$

and when B is of the form in (19)

$$(a_1 - a_2)\sin \varphi = a_3 \cos \varphi. \quad (22)$$

Since the Hessian of a defined in (2) has eigenvalues of un-equal magnitude, $a_1 + a_2 \neq 0$, and either $a_1 \neq a_2$, or $a_3 \neq 0$. This leaves only four possible solutions for B :

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \begin{bmatrix} \cos \varphi_0 & \sin \varphi_0 & 0 \\ \sin \varphi_0 & -\cos \varphi_0 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \begin{bmatrix} -1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \begin{bmatrix} -\cos \varphi_0 & -\sin \varphi_0 & 0 \\ -\sin \varphi_0 & \cos \varphi_0 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad (23)$$

where $\varphi_0 = \arctan \frac{a_3}{a_1 - a_2}$. Thus $\tilde{A} = BA$ can relate to A in only four possible ways.

Next, we consider the lighting \tilde{l} associated with each shape \tilde{A} . Equation (15) implies $\tilde{A}^T \tilde{l} = A^T l$ or $\tilde{A}^T \tilde{l} = -A^T l$ but the latter has shadows, so

$$A^T l = \tilde{A}^T \tilde{l} = A^T B^T \tilde{l}. \quad (24)$$

When A is full rank, (24) implies a unique \tilde{l} given by

$$\tilde{l} = (B^T)^{-1} l = B l. \quad (25)$$

If $\text{Rank}(A) = 2$, we define l_\perp as the component of l in the null space of A^T . Then, from (24), we have

$$B^T \tilde{l} = l + c l_\perp \implies \tilde{l} = B(l + c l_\perp), \quad (26)$$

where c is a scalar. In this case there is a 1D family of \tilde{l} for each of the four shapes \tilde{A} . \square

Fig. 2 provides an example of the four choices of shape/light pairs in the generic, non-cylindrical case when both

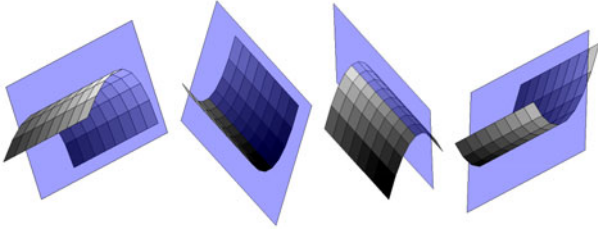


Fig. 3. Lighting solutions in the cylinder case, when one of the eigenvalues of the surface Hessian is zero. There is a 1D family of lighting (any lighting direction in the blue plane with appropriate strength) for each of the four shapes that can produce the same image.

eigenvalues of the surface Hessian are non-zero. Without loss of generality, we consider a rotated co-ordinate system where $a_3 = 0$, i.e., the x - and y -axes are aligned with the eigenvectors of the surface Hessian. Then, the four solutions from (23) are:

$$([a_1, a_2, 0, a_4, a_5], [l_x, l_y, l_z]), \quad (27)$$

$$([-a_1, -a_2, 0, -a_4, -a_5], [-l_x, -l_y, l_z]), \quad (28)$$

$$([a_1, -a_2, 0, a_4, -a_5], [l_x, -l_y, l_z]), \quad (29)$$

$$([-a_1, a_2, 0, -a_4, a_5], [-l_x, l_y, l_z]). \quad (30)$$

The first choice is the surface/lighting pair (a, l) that actually induced the image. The second corresponds to the well-known convex-concave ambiguity [10], and is obtained by reflecting both the light and the normals across the view direction. The last two choices (29)-(30) correspond to performing the reflection separately along each of the eigenvector directions of the Hessian matrix. These form a second concave-convex pair.

When one of the Hessian eigenvalues is zero (say $a_2 = 0$ in our rotated co-ordinate system), the patch surface is a cylinder and it is possible to construct a 1D family of lights for each of the four surfaces:

$$\tilde{l} = \text{diag}\{\text{sign}(\tilde{a}_1 a_1), \text{sign}(\tilde{a}_5 a_5), 1\} (l + c \cdot [0, 1, a_5]^T) \quad (31)$$

for any $c \in \mathbb{R}$ such that no pixel is in shadow. Fig. 3 shows an example of four cylindrical surfaces and associated families of lights that can produce the same image.

Theorem 1 applies when the Hessian eigenvalues of any solution shape are not equal in magnitude. What happens when shape solutions have Hessian eigenvalues that are of equal magnitude? There are two distinct cases. The first is when the Hessian is zero and the true surface is planar. In

this case every surface normal in the patch is identical, and the well-known point-wise cone ambiguity applies to the patch as a whole: The observed image can be explained by a one-parameter family of planar surfaces for *every* light l .

In the second case, the true surface is not planar but the magnitudes of the two eigenvalues of the Hessian matrix are equal. Unlike the planar ambiguity, there is not an infinite number of surfaces that can combine with every lighting. But as depicted in Fig. 4, there is still an infinite number of allowable patch/lighting pairs. We note that all quadratic surfaces in this category can be expressed as either one of two following forms:

$$a = [r \cos \theta, -r \cos \theta, 2r \sin \theta, p \cos \theta - q \sin \theta, p \sin \theta + q \cos \theta], \quad (32)$$

$$a = [\lambda r, \lambda r, 0, \lambda p, -\lambda q], \quad (33)$$

where $\theta \in (-\pi, \pi]$, $\lambda \in \{-1, +1\}$, $r \in \mathbb{R}^+$, and $p, q \in \mathbb{R}$. Given fixed values of r, p and q , these surfaces generate identical images when paired with lighting

$$l = [l_x \cos \theta - l_y \sin \theta, l_x \sin \theta + l_y \cos \theta, l_z], \quad (34)$$

for surfaces (32), or with

$$l = [\lambda l_x, -\lambda l_y, l_z], \quad (35)$$

for surfaces (33), with fixed values of l_x, l_y, l_z .

3.2 Unique Shape when Light Is Known

Theorem 2. Given intensities $I(x, y)$ at a non-degenerate set of locations Ω , a known light l , and a quadratic patch a that satisfies the set of equations in (9), if the planar component $[l_x, l_y]$ of the light is non-zero (i.e., l is not equal to the viewing direction) and not an eigenvector of the Hessian of a , then the solution a is unique.

Proof of Theorem 2. Without loss of generality, we choose a co-ordinate system where $a_3 = 0$. Note that for any such choice l_x and l_y will both be non-zero, unless $[l_x, l_y]$ is zero or an eigenvector of the surface Hessian, which is ruled out by the statement of the theorem.

If the Hessian of a has eigenvalues with unequal magnitudes, then it is easy to see that each of the four possible solutions from Theorem 1 has distinct light from (25) and (26), and therefore for a fixed light, the shape is unique. A Hessian with *equal* eigenvalues is ruled out since then every light-direction would be an eigenvector. When the eigenvalues have equal magnitudes but opposite signs, a must be of the form in (32) with $\theta = 0$ or π

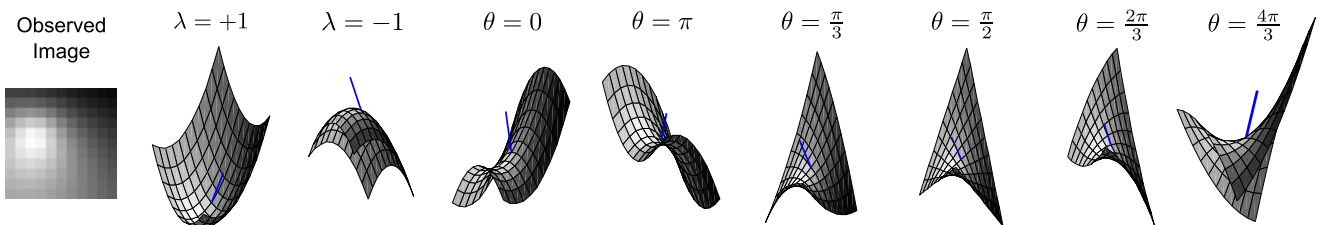


Fig. 4. When Hessian eigenvalues are equal in magnitude, there is a continuous family of patch/lighting pairs (given by (32) and (33)) that produce the same image. Note that the first four pairs above are analogous to Fig. 2.



Fig. 5. Left: Two quadratic surfaces that produce the same image when the light is aligned with one of their common Hessian eigenvectors. For other view and light configurations (e.g., right) their images are distinct.

(since $a_3 = 0$) and $r = |a_1| = |a_2|$. In this case too, we see that each member of the continuous family of solutions—with $\theta \in (-\pi, \pi]$ for surface (32) and light (34), or $\lambda \in \{-1, +1\}$ for surface (33) and light (35)—has a distinct light-direction. \square

When the conditions in Theorem 2 are not satisfied, there are shape ambiguities as follows. First, planar patches have Hessians with zero eigenvalues so that every l is an eigenvector; this leads to an infinite set of planar shape explanations for any given light. Second, when the light and view directions are the same, there are generically four shape solutions analogous to Fig. 2 or, in the case of equal eigenvalue magnitudes, a continuous family of solutions analogous to Fig. 4. Finally, when the true surface is not planar but the azimuthal component of the light $[l_x, l_y]$ happens to be aligned with one of the Hessian eigenvectors, it is possible to construct a second solution by performing a reflection of the normals across that eigenvector direction. Fig. 5 demonstrates this with photographs of two 3D-printed surfaces that are distinct but related by a horizontal reflection of their normals.

4 AMBIGUITY IN THE PRESENCE OF NOISE

The uniqueness results from the previous section suggest that among the many possible models one could use for local shapes—such as splines, linear subspaces, exemplar dictionaries [17], or continuous functions with smoothness constraints as in [21]—the quadratic function model may be particularly useful. However, before we can use this model for inference, we must understand the effects of deviations, such as intensity noise and higher-order (non-quadratic) components of local shape. To this end, we provide some intuition about the types of quadratic shapes that *almost* satisfy the polynomial system (9) and thus become likely explanations in the presence of noise. These intuitions motivate a statistical inference technique that will be introduced in Section 5.

In the rest of this paper, we assume that the light direction $l/\|l\|$ and the albedo/light-strength product $\|l\|$ are known. Then, the polynomial system (9) relating the quadratic parameters a to the observed intensities I can be understood as combining two types of constraints on the patch normals $n = [n_x, n_y, 1]$. First, each pixel's normal is constrained by its intensity to a light-centered circle of directions as per (7). This is shown in the left of Fig. 6, where the circle of directions is parameterized by “azimuthal” angle

$$\theta = \arctan\left(\frac{n_x l_y - n_y l_x}{l_x^2 + l_y^2 - l_z(n_x l_x + n_y l_y)}\right). \quad (36)$$

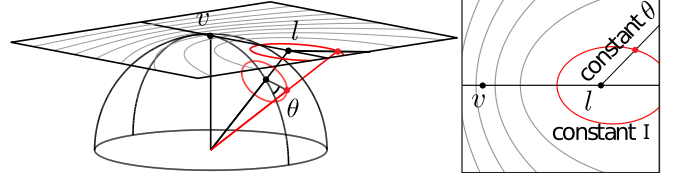


Fig. 6. The light-centered cone of possible surface normals at any image point projects radially to a conic on the projective plane. We parameterize these conics by the radial projection of spherical angle θ .

The second type of constraint comes from the quadratic shape model, which induces a joint geometric constraint on the set of surface normals that belong to the patch. This joint constraint has an intuitive interpretation when we represent the normals, light, and view as points on the plane defined by $n_z = 1$ (the so-called projective plane [23]). This representation is constructed by radially-projecting the hemisphere of directions onto the plane as shown in Fig. 6. The view is the origin of the plane, the light projects to another planar point $(l_x, l_y)/l_z$, and each pixel's θ -parameterized circle of normal azimuthal directions projects to a conic section, still parameterized by θ . The set of normals that lie on different conics but have the same azimuthal angle θ form a ray (right of Fig. 6), and an inversion in the sign of θ corresponds to a reflection of the surface normal across light point.

Using this representation, Fig. 7 visualizes the two types of constraints (under a light with $l_y = 0$) for 25 normals at a 5×5 grid of (x, y) pixel locations. In addition to each pixel's normal being constrained to its conic, the set of normals is collectively constrained, via (6), to be a symmetric affine grid. Therefore, solving the polynomial system for quadratic coefficients a amounts to finding a symmetric affine grid that aligns properly with the per-pixel conics. Theorem 2 tells us there is only one grid that aligns perfectly, but as shown in the figure, there will be other grids that come close. When there is noise, the shapes corresponding to all of these grids become likely explanations, even though they are physically quite different from one another. To avoid over-committing, local inference systems must output distributions of shapes that encode this fact.

Then, a natural question is: do we need to search the entire five-dimensional space of quadratic parameters a to

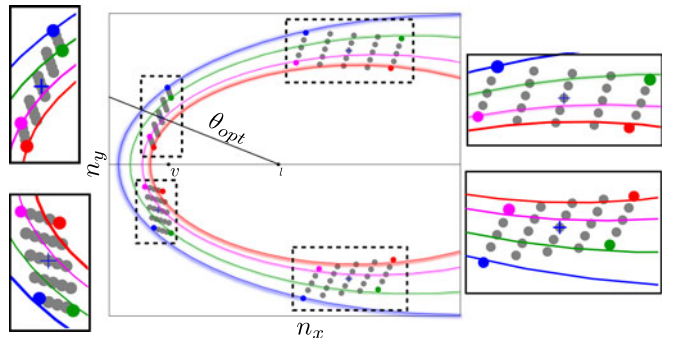


Fig. 7. Exact and approximate solutions for quadratic shape. Each color corresponds to a pixel in the patch (four are shown in the plot), whose intensity defines a conic curve that the normal vector should lie on. The normal vectors for a quadratic patch should form an affine grid on the projective plane, and good-fit shapes have grids that are well-aligned with the corresponding conics. The top left grid corresponds to an exact fit.

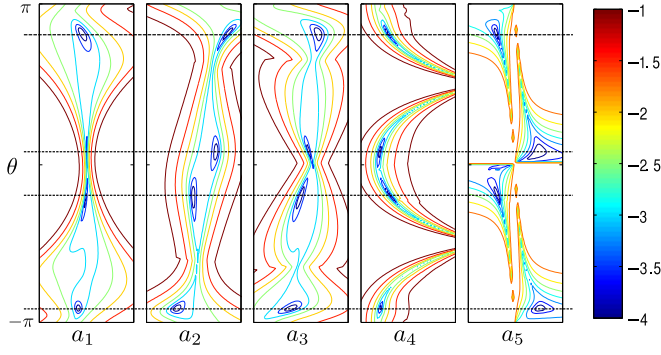


Fig. 8. Iso-contours of RMS intensity error for renderings of best-fit shape parameters $(a_1, a_2, a_3, a_4, a_5)$ when θ is fixed. Close fits occur at very different orientations (four modes here), but for any fixed orientation θ the remaining shape parameters are very constrained.

find all the likely approximate solutions? To answer this question, we note that these approximate solutions are intuitively expected to arise from the degenerate cases detailed in Theorem 2. For example, we find that these solutions often occur in pairs corresponding to reflections across the light direction (i.e., across the x -axis in Fig. 7), which would correspond to a second exact solution if the light were a eigenvector of the surface Hessian. Remember that the most ambiguous degeneracy is the one induced by the true surface being planar, when all the conics overlap and there is a continuous set of solutions whose normals can be parametrized by a single angle θ as per (36). Based on this intuition, we define $\theta(a)$ as the first-order orientation of the shape a to be the angle of the center normal, and find empirically that it is sufficient to search along only a one-dimensional manifold parametrized by this angle.

In Fig. 7, this search can be understood as fixing the value of $\theta(a)$, and warping an affine grid by optimizing the parameters a_1, a_2, a_3, a_4, a_5 to fit the conic intensity constraints. We see that this leaves very little play in the parameters, so the shapes a of possible solutions are highly constrained once $\theta(a)$ is fixed. This effect is further visualized in Fig. 8, which shows contours of constant RMS intensity difference—equally spaced in value on a logarithmic scale—between the observed intensities and the Lambertian renderings of best-fit shapes obtained by fixing $\theta(a)$ and one coefficient (say, a_1) and then optimally fitting the others (say, a_2, a_3, a_4, a_5). The four “close fits” appear as the four modes in these plots, where the value of $\theta(a)$ strongly constrains each coefficient of low-error shapes a .

5 LOCAL SHAPE DISTRIBUTIONS

Armed with intuition about the characteristics of approximate solutions for the quadratic-patch model, we now develop a method for inferring shape distributions at any local image patch of any size. The output for each image patch is a set of quadratic shapes of the same size that correspond to a discrete sampling along a θ -parametrized one-dimensional manifold, as well as a probability distribution over this set of quadratic shapes. The previous sections have demonstrated that shading in some image patches is inherently more informative than others. Our goal is to create a compact description of this ambiguity in each local region at multiple scales, thereby providing a

useful mid-level representation of “intrinsic” scene information for vision.

5.1 Computing Quadratic Shape Proposals

Given the intensities $I_o(x, y)$ at a patch $(x, y) \in \Omega$, we first generate a set of quadratic proposals for the shape of that patch, and based on the intuition from the previous section, we index these proposals angularly in reference to the light l . Consider a discrete set of uniformly-spaced values θ^j , $j \in \{1, \dots, J\}$ over $(-\pi, \pi]^2$, and for each angle θ^j we find the corresponding quadratic shape a^j that best explains the observed intensities $I_o(x, y)$ in terms of minimum sum of squared errors:

$$a^j = \arg \min_{a: \theta(a) = \theta^j} \sum_{(x, y) \in \Omega} \|I_o(x, y) - I(x, y; a)\|^2, \quad (37)$$

where $I(x, y; a)$ is defined as per (7).

Let $(0, 0)$ be the center of the patch. Then since $\theta(a_i)$ is fixed, the quadratic coefficients a_4 and a_5 of a^j only have one degree of freedom, and can be re-parametrized in terms of a single variable $r \in \mathbb{R}^+$ that indexes points along the constant θ ray on the projective plane:

$$a_4 = -\frac{l_x}{l_z} - r \left(-\frac{l_x}{l_z} \cos \theta^j + l_y \sin \theta^j \right), \quad (38)$$

$$a_5 = -\frac{l_y}{l_z} - r \left(-\frac{l_y}{l_z} \cos \theta^j - l_x \sin \theta^j \right). \quad (39)$$

Therefore, the non-linear least-squares minimization in (37) is over the four variables $a_{1:3}, r$, and can be efficiently carried out with Levenberg-Marquardt [24]. We found empirically that it is insensitive to initialization, and use $[0, 0, 0, r_0]$ in our experiments, where r_0 is chosen such that the center pixel lies on the corresponding conic.

This minimization occurs independently and in parallel for every patch in an image, and it can therefore be parallelized over an arbitrary number of CPU cores, on a single machine or a cluster of machines, as required for increasing image sizes. Our reference implementation considers $J = 21$ quantized angles for each patch, and takes one minute on an eight-core machine for inference on all overlapping 5×5 patches in a 128×128 image.

5.2 Computing Shape Likelihoods

Next, we define a probability distribution over these shape proposals by computing the likelihoods for the observed intensities being generated by each proposed shape a^j . We introduce a model for the deviation between observed intensity $I_o(x, y)$ and expected intensity $I(x, y; a)$ from a proposal a at each location (x, y) :

$$I_o(x, y) | a \sim \mathcal{N}(I(x, y; a); \sigma_i^2 + \sigma_z^2(x, y; a)). \quad (40)$$

This is a Gaussian distribution conditioned on a , where the variance at each pixel (x, y) is a sum of two terms. The first is

2. For some patches, we consider closer-spaced samples over a shorter interval when values close to $\pm\pi$ do not correspond to physically feasible estimates for shape.

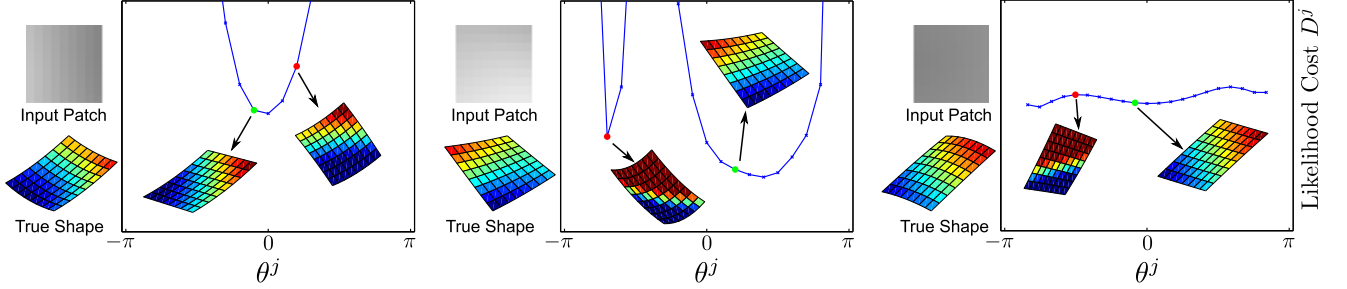


Fig. 9. Shape likelihood distributions inferred from image patches. Graphs show likelihood cost D^j over first-order orientation θ^j , each of which is associated with a shape proposal a^j .

additive *i.i.d.* intensity noise σ_i^2 induced, for example, by sensor noise. The second is a function of a and varies spatially across the patch, capturing the fact that the veridical shape may exhibit higher-order (non-quadratic) variations at this patch's scale. It is the expected variance in intensity $\sigma_z^2(x, y; a)$ induced by higher-order components of shape that exist on top of the shape predicted by a at the current scale.

To compute $\sigma_z^2(x, y; a)$, we model the deviations of the true normals $(\tilde{n}_x, \tilde{n}_y)$ from those predicted by a as *i.i.d.* Gaussian random variables:

$$\tilde{n}_x(x, y) \sim \mathcal{N}(n_x(x, y; a); \sigma_{n0}^2), \quad (41)$$

$$\tilde{n}_y(x, y) \sim \mathcal{N}(n_y(x, y; a); \sigma_{n0}^2), \quad (42)$$

where σ_{n0}^2 is the expected normal variance of these deviations, which is set to 10^{-6} in our experiment. Then, we compute $\sigma_z^2(x, y; a)$ as the expected variance in intensity over the distribution of \tilde{n}_x, \tilde{n}_y :

$$\begin{aligned} \sigma_z^2(x, y; a) &= \mathbb{E}_{\tilde{n}_x, \tilde{n}_y} \left\| I(x, y; a) - \frac{l_x \tilde{n}_x + l_y \tilde{n}_y + l_z}{\sqrt{\tilde{n}_x^2 + \tilde{n}_y^2 + 1}} \right\|^2 \\ &= \mathbb{E}_{\tilde{n}_x, \tilde{n}_y} \left\| \frac{l_x n_x(a) + l_y n_y(a) + l_z}{\sqrt{n_x^2(a) + n_y^2(a) + 1}} - \frac{l_x \tilde{n}_x + l_y \tilde{n}_y + l_z}{\sqrt{\tilde{n}_x^2 + \tilde{n}_y^2 + 1}} \right\|^2. \end{aligned} \quad (43)$$

We find that for lights not aligned with the view, i.e., $|l_z| < 1$, this expression can be reliably approximated as:

$$\sigma_z^2(x, y; a) \approx \frac{(l_x^2 + l_y^2) \sigma_{n0}^2}{n_x^2(x, y; a) + n_y^2(x, y; a) + 1}. \quad (44)$$

Intuitively, this says that the observed intensity is less sensitive to perturbations in the normal $[n_x, n_y]$ when the surface is tilted away from the viewing direction.

Putting everything together, we compute a cost D^j for every proposal a^j , defined as the negative log-likelihood of all observed intensities under the above model:

$$\begin{aligned} D^j = -\log p(I_o | a^j) &= \sum_{(x,y) \in \Omega} \frac{1}{2} \left[\log(\sigma_i^2 + \sigma_z^2(x, y; a^j)) \right. \\ &\quad \left. + \frac{(I_o(x, y) - I(x, y; a^j))^2}{\sigma_i^2 + \sigma_z^2(x, y; a^j)} \right]. \end{aligned} \quad (45)$$

5.3 Evaluation

We evaluate the accuracy of the proposed local shape distributions using images of a set of six random surfaces synthetically rendered (with the light at an elevation of 60 degree), where each surface is created by generating a 5×5 grid of random depth values, and then smoothly interpolating these to form a 128×128 surface (see [1], and Fig. 12 for an example).

Fig. 9 shows likelihood distributions and proposed shapes for representative image patches from this synthetic dataset. Empirically, we find that the distributions D^j can have a single peak (left), be multi-modal (center), or nearly uniform (right); re-enforcing our intuition that shading in some image patches is more informative than others. Note that given the correct value of θ^j (green dot in the figure), the corresponding estimated shape proposal a^j yields an accurate reconstruction of the true shape in all three cases shown.

We perform a quantitative evaluation of accuracy using all overlapping patches from all six random surfaces and for three different patch-sizes (roughly 80 k -90 k patches in total for each size). We are interested in knowing: (i) how often the veridical shape is among the set of shape proposals a^j at a patch; and (ii) whether the cost D^j is an informative metric for determining which proposals a^j are most accurate. To this end, for each image patch in the evaluation set, we sort the proposed quadratic shapes according to their likelihood costs D^j , compute for each proposal the mean angular difference between its surface normals the veridical ones, and record for increasing values $N = \{1, \dots, J = 21\}$ the lowest mean angular error among the N most-likely shape proposals. Fig. 10 shows the statistics of these errors across all test patches for increasing values of N . Although the most-likely shape proposal (i.e. $N = 1$) is often reasonably close to true shape, the error quantiles decrease significantly more as we consider larger sets of likely proposals. This emphasizes the value of maintaining full distributions of local shape as a mid-level scene representation, as opposed to “over-committing” to only one (often sub-optimal) shape proposal for each patch through a process of hard local decision-making.

Fig. 10 also provides some insight about the effect of patch-size, and it shows that patches at multiple scales tend to be complimentary. Smaller patches are more likely to have lower errors when considering the full set of proposals ($N = J = 21$), since the veridical shape is much more likely to be exactly quadratic at smaller scales. But, as evidenced

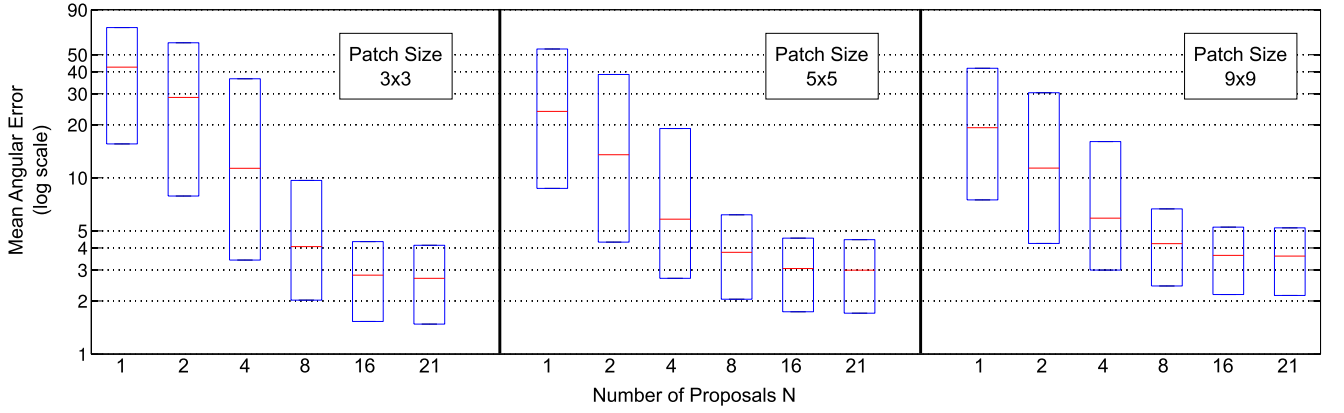


Fig. 10. Local shape accuracy. We show quantiles (25 percent, median, 75 percent) of each patch’s mean normal angular error, for the best estimate amongst the N most-likely shape proposals for each patch, for different values of N and for different patch sizes. The quantiles for $N = 1$ correspond to making a hard decision at each patch, and errors for $N = 21$ correspond to the best estimate amongst the full set of proposals.

by the relatively smaller error quantiles for lower values of N , larger patches tend to be more informative, with their likelihoods D^j being better predictors of *which* of the proposals a^j is the true one.

6 SURFACE RECONSTRUCTION

To demonstrate the utility of our theory and local distributions for higher-level scene analysis, we consider the application of reconstructing object-scale surface shape when the light l is known. The local representations provide concise summaries of the shape information available in each image patch, and they do this without “over-committing” to any one local explanation. This allows us to achieve reliable performance with very a simple algorithm for global reasoning that infers object-scale shape through simple iterations between: 1) choosing one likely shape proposal for each local patch; and 2) fitting a global smooth surface to the set of chosen per-patch proposals.

Formally, our goal is to estimate the depth map $Z(x, y)$ from an observed intensity image $I(x, y)$, with known lighting l and under the assumption that the surface is predominantly texture-less and diffuse, i.e., the shading equation (7) holds at most pixels. We first compute local distributions as described

in the previous section, by dividing the surface into a mosaic of overlapping patches of different sizes. We let $p \in \{1, \dots, P\}$ index all patches (across different patch-sizes), with Ω_p corresponding to the pixel locations, and $\{a_p^j, D_p^j\}_j$ denoting the local shape proposals and distribution for each patch p .

In addition to the J proposals at each patch, we use an approach similar to [18] and include a dummy proposal $\{a_p^{J+1} = \phi, D_p^{J+1} = D_\phi\}$ in the distribution for every patch. This serves to make the surface reconstruction robust to outliers, such as when the local patch deviates significantly from a quadratic approximation (e.g. sharp edges or depth discontinuities), or when the observed intensities vary from the diffuse model in (7), e.g. specularities, shadows, or albedo variations due to texture.

We formulate the reconstruction problem as simultaneously finding a depth estimate Z and a labeling $L_p \in \{1, \dots, J+1\}, \forall p$ that minimize the cost function:

$$C(Z, \{L_p\}, \lambda) = \sum_{p=1}^P \left(\lambda D_p^{L_p} + \sum_{(x,y) \in \Omega_p} \delta(Z, a_p^{L_p}, x, y) \right), \quad (46)$$

where λ is a scalar weight, and δ measures the agreement between the local shape proposal $a_p^{L_p}$ and Z at (x, y) :

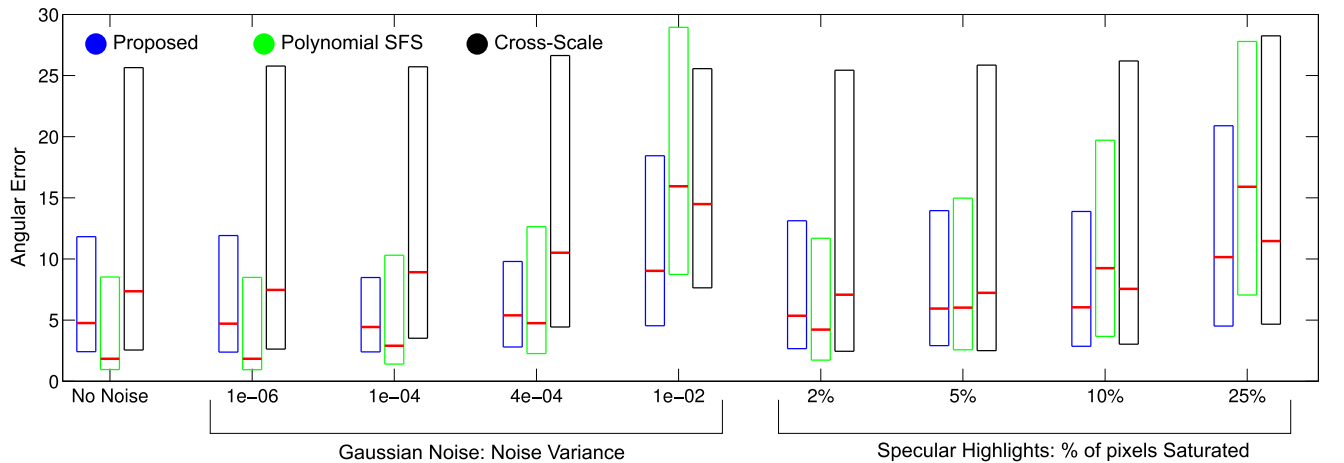


Fig. 11. Surface reconstruction accuracy for different methods on synthetic images of random surfaces. Shown here are quantiles (25 percent, median, 75 percent) of angular errors of individual normals across all surfaces for images rendered with different degrees of additive Gaussian noise, and specular highlights.

$$\delta(Z, a, x, y) = \left\| \begin{matrix} \nabla_x Z(x, y) - n_x(x, y; a) \\ \nabla_y Z(x, y) - n_y(x, y; a) \end{matrix} \right\|^2, \quad (47)$$

if $a \neq \phi$, and 0 otherwise.

We use an iterative algorithm to minimize the cost C , alternating updates to Z and $\{L_p\}$ based on the current estimate of the other. Given the current estimate Z^* of the depth map at each iteration, we update the label L_p of every patch independently (and in parallel):

$$L_p \leftarrow \arg \min_{L \in \{1, \dots, J+1\}} \lambda D_p^L + \sum_{(x,y) \in \Omega_p} \delta(Z^*, a_p^L, x, y). \quad (48)$$

Similarly, given a labeling $\{L_p^*\}$ (with corresponding shape proposals $\{a_p^*\}$), we compute the depth map Z that minimizes the cost in (46) as:

$$\begin{aligned} Z &\leftarrow \arg \min_Z \sum_p \sum_{(x,y) \in \Omega_p} \delta(Z, a_p^*, x, y) \\ &= \arg \min_Z \sum_{x,y} w^*(x, y) \left\| \begin{matrix} \nabla_x Z(x, y) - n_x^*(x, y) \\ \nabla_y Z(x, y) - n_y^*(x, y) \end{matrix} \right\|^2, \end{aligned} \quad (49)$$

where $w^*(x, y)$ is the number patches that include (x, y) and have not been labeled as outliers, and $n_x^*(x, y), n_y^*(x, y)$ their mean normal estimates:

$$\begin{aligned} \Omega^{-1}(x, y) &= \{p : (x, y) \in \Omega_p, a_p^* \neq \phi\}, \\ w^*(x, y) &= |\Omega^{-1}(x, y)|, \\ n_x^*(x, y) &= \frac{1}{w^*(x, y)} \sum_{p \in \Omega^{-1}(x, y)} n_x(x, y; a_p^*), \\ n_y^*(x, y) &= \frac{1}{w^*(x, y)} \sum_{p \in \Omega^{-1}(x, y)} n_y(x, y; a_p^*). \end{aligned} \quad (50)$$

The computation in (49) could be carried out exactly and efficiently using the Frankot-Chellappa algorithm [25] if all $w^*(x, y)$ were equal. But this is not the case since $w^*(x, y)$ will be lower near the boundary and in regions where some patches have been detected as outliers. Nevertheless, we find that [25] provides an acceptable approximation in these cases. We use [25] throughout the alternating iterations until Z and $\{L_p\}$ converge,³ and then we run a limited number of additional iterations using conjugate-gradient to compute step (49) exactly.

To speed up convergence, we smooth the estimate of Z in the first few iterations with a Gaussian filter with variance σ and set $\lambda' = \sigma^2 \lambda$, starting from an initial value σ_0 that is decreased by a constant factor σ_f till it reaches 1 (at which point we stop smoothing). We also initially run the algorithm over only the valid proposals at each patch till convergence, and then introduce the dummy proposal ϕ . We set the parameters λ and D_ϕ automatically based on the input distributions— λ is set to 1/4th the reciprocal of the mean of the differences between the minimum and median likelihood costs across all patches at the smallest scale, and D_ϕ is set to $10\lambda^{-1}$. The reconstruction from local patch proposals

3. We simply set $n_x^*(x, y) = n_y^*(x, y) = 0$ when $w^*(x, y) = 0$.

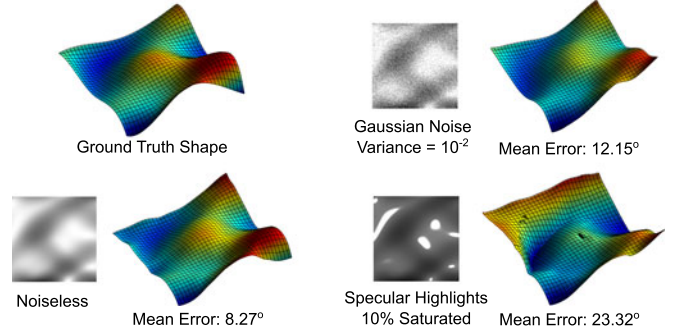


Fig. 12. Reconstructions by the proposed method on different rendered images of a synthetic surface.

takes 40 seconds on average on an eight-core machine for 128×128 images with local distributions at four scales (not including the computation time for estimating local proposals, which is reported in Section 5.1).

6.1 Evaluation

We first quantitatively evaluate the proposed reconstruction algorithm, under known lighting, with the random surfaces described in Section 5.3. We render images with different amounts of additive white Gaussian noise, as well as with specular highlights. For the latter, we use the Beckmann and Spizzichino [26] model and consider different values of “surface smoothness” to get images with increasing numbers of saturated pixels. We mask out pixels that are saturated during estimation, but note that many nearby unclipped pixels will also include a non-zero specular component that violates the diffuse shading model.

Fig. 11 summarizes the performance—using local distributions of 3×3 , 5×5 , 9×9 , and 17×17 overlapping patches—and compares it to two state-of-the-art methods. The first is the iterative algorithm proposed by Ecker and Jepson [14] (labeled “Polynomial SFS”). The second (labeled “Cross-scale”) is the shape from shading component of the SIRFS method [21], i.e., where we treat the light and shading-image as given, and do not use contour information (so as to evaluate the shading cue in isolation). The cross-scale method uses an over-complete, multi-scale representation of the global depth map and minimizes the rendering error along with the likelihood of the recovered shape under a prior model. For both methods we use implementations provided by their authors, and for the cross-scale method, we use the author-provided prior parameters that were trained on the MIT intrinsic image database [27].⁴

We see from Fig. 11 that while the polynomial SFS method performs the best in the noiseless case, the proposed algorithm is more robust to both Gaussian noise and the structured artifacts from specular highlights. The cross-scale method is also reasonably robust to these effects due to its use of a shape prior, but in general has higher errors. Fig. 12 provides example reconstructions for the proposed method for one surface—the full set of reconstructions are available at [1].

Next, we evaluate all algorithms using photographs of seven relatively-diffuse objects, captured with a Canon EOS

4. We also evaluated the cross-scale method with a prior trained on the random surfaces, but this did not improve performance.


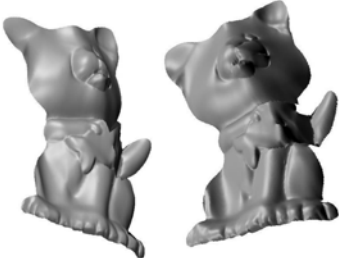
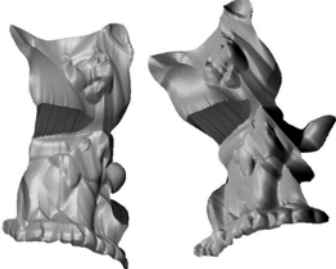
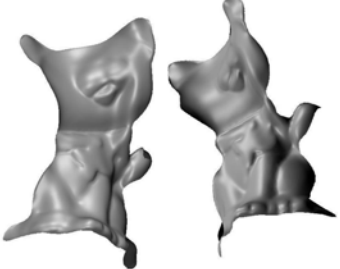

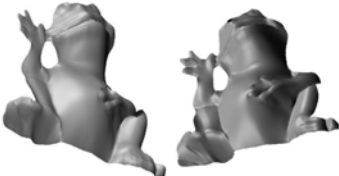
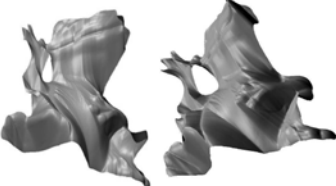
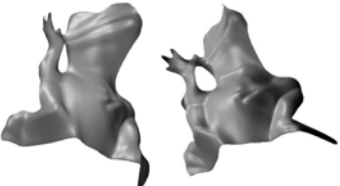

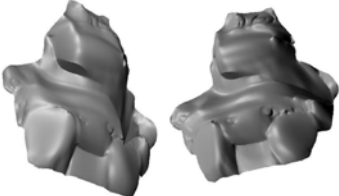
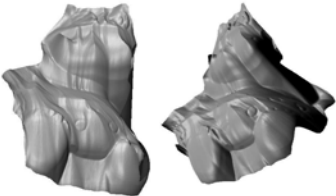

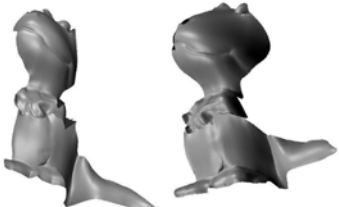
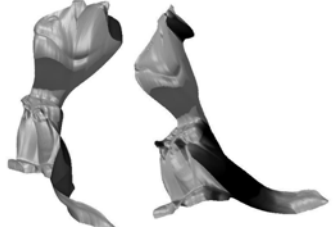






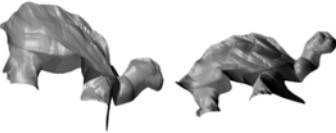
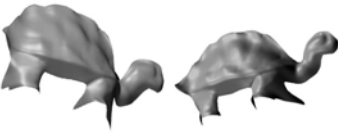


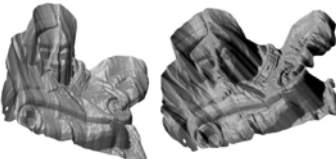
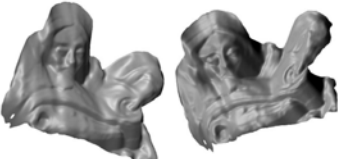
Input Image	Proposed	Polynomial SFS	Cross-Scale
			
Resolution 640×500	Median Angular Error 14.83°	Median Angular Error 24.81°	Median Angular Error 20.02°
			
Resolution 590×690	Median Angular Error 11.80°	Median Angular Error 20.77°	Median Angular Error 19.86°
			
Resolution 580×580	Median Angular Error 20.25°	Median Angular Error 17.50°	Median Angular Error 21.00°
			
Resolution 720×660	Median Angular Error 12.70°	Median Angular Error 22.33°	Median Angular Error 23.26°
			
Resolution 550×760	Median Angular Error 15.29°	Median Angular Error 15.58°	Median Angular Error 13.17°
			
Resolution 450×850	Median Angular Error 17.90°	Median Angular Error 14.50°	Median Angular Error 11.96°
			
Resolution 790×1070	Median Angular Error 28.13°	Median Angular Error 29.21°	Median Angular Error 25.80°

Fig. 13. Surface reconstruction on real captured data. We show two novel view points for each reconstruction, and the median angular error between estimated surface normal vectors and ground truth surface normal vectors. A more interactive visualization is available at [1].

40D camera under directional lighting, with two chrome spheres in the scene to measure light direction. These photographs contain non-idealities such as mutual illumination, self-shadowing, and slight texture. For each object under a

fixed viewpoint, we took twenty images with varying light directions, with which we can recover the normal vectors as well as depth map by photometric stereo to a high accuracy. We use this recovered shape as ground truth for our

evaluation. All the captured images, calibration information, and recovered normal and depth maps are available on our project page [1].

For each object, we choose a single image as input to evaluate the performance of different SFS frameworks. The 99th percentile intensity value of the image is assumed to correspond to the albedo times light intensity and used for image normalization; and since these images are larger, we use local distributions at two additional patch-scales: 33×33 and 65×65 . The surfaces reconstructed using the different methods and measured light direction are shown in Fig. 13 along with median angular error values. We find that in most cases, the proposed algorithm provides a better reconstructions of object-scale shape than the baselines.

7 DISCUSSION

Our theoretical analysis shows that in an idealized quadratic world, local shape can be recovered uniquely in almost every local image patch, without the use of singular points, occluding contours, or any other external shape information. Beyond this idealized world, our evaluations on synthetic and captured images suggest that one can infer, efficiently and in parallel, concise multi-scale local shape distributions that are accurate and useful for global reasoning.

There are many viable directions for interesting future work. Foremost among these is the joint estimation of shape, lighting, and albedo. The reconstruction algorithms proposed in this paper are limited to the case when lighting is known, but the uniqueness results in Section 3.1 suggest that simultaneous reconstruction of shape and lighting may also be possible. Theorem 1 tells us that, in an idealized quadratic world, there are generically four lights l that can explain each local patch, and that these quadruples of possible lights will vary from patch to patch according to the directions of each patch's Hessian eigenvectors. Intuitively, one might infer the true light (along with its reflection across the view, which is always equally-likely) as the one that is common to all or most of the per-patch quadruples.⁵ Practically speaking, it is likely that for a reconstruction algorithm to handle unknown lighting, it will need to jointly reason about shape, lighting, and varying albedo, in the same spirit as Barron and Malik [21]; and that such reasoning will benefit from an analysis of the joint ambiguities that are induced by noise and non-quadratic shape, similar to what was done for shape alone in Sections 4 and 5.

Also, while we provide a means to extract a single estimate of the global surface from local shape distributions, one could also imagine using reasoning about consistency and outliers to allow the full distributions of neighboring patches to collaboratively refine themselves. This could be useful, for instance, when the object boundaries in a scene are not known a-priori. These refined local distributions may then be able to identify depth discontinuities in the scene, and help segment out individual objects for shape recovery.

5. We have experimented with a direct implementation of this intuition that does a brute-force search only on lighting direction, assuming a known constant light-strength and albedo, and with pooling local estimates without considering consistency or noise. This method worked reasonably well in many cases, but was computationally expensive and not entirely robust.

Finally, it will be interesting to pursue combining our shading-based local distributions with complementary reasoning about contours, shading keypoints [28], texture, gloss, shadows, and so on—treating these as additional cues for shape, as well as to better identify outliers to our smooth diffuse shading model. We also believe it is worth integrating these local shape distributions into processes for higher-level vision tasks such as pose estimation, object recognition, and multi-view reconstruction, where one can imagine additionally using top-down processing to aid local inference, for example by exploiting priors on local quadratic shapes that are based on object identity or scene category.

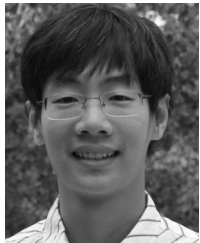
ACKNOWLEDGMENTS

The authors would like to thank the associate editor and reviewers for their valuable comments. YX, AC, DWJ, and TZ acknowledge support from the National Science Foundation under Grants no. 1212928, 0926148, and 0915977. RB and DWJ were supported in part by the US-Israel Binational Science Foundation under Grant no. 2010331. RB also acknowledges support from the Citigroup Foundation. Some of this work was performed while TZ was at the Weizmann Institute of Science as a Feinberg Foundation Visiting Faculty Program Fellow.

REFERENCES

- [1] *Project Page*. (2014). [Online]. Available: <http://vision.seas.harvard.edu/qsf/>
- [2] J. D. Durou, M. Falcone, and M. Sagona, "Numerical methods for shape-from-shading: A new survey with benchmarks," *Comput. Vis. Image Understanding*, vol. 109, pp. 22–43, 2008.
- [3] B. K. P. Horn and M. J. Brooks, *Shape from Shading*. Cambridge, MA, USA: MIT Press, 1989.
- [4] R. Zhang, P. Tsai, J. E. Cryer, and M. Shah, "Shape from shading: A survey," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 21, no. 8, pp. 690–706, Aug. 1999.
- [5] J. Oliensis, "Uniqueness in shape from shading," *Int. J. Comput. Vis.*, vol. 6, pp. 75–104, 1991.
- [6] J. Oliensis, "Shape from shading as a partially well-constrained problem," *CVGIP: Image Understanding*, vol. 54, pp. 163–183, 1991.
- [7] E. Prados and O. Faugeras, "A generic and provably convergent shape-from-shading method for orthographic and pinhole cameras," *Int. J. Comput. Vis.*, vol. 65, pp. 97–125, 2005.
- [8] E. Prados and O. Faugeras, "Unifying approaches and removing unrealistic assumptions in shape from shading: Mathematics can help," in *Proc. 8th Eur. Conf. Comput. Vis.*, 2004, pp. 141–154.
- [9] E. Prados and O. Faugeras, "Shape from shading: A well-posed problem?" in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recog.*, 2005, pp. 870–877.
- [10] A. P. Pentland, "Local shading analysis," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. PAMI-6, no. 2, pp. 170–187, Mar. 1984.
- [11] J. Wagemans, A. J. Van Doorn, and J. J. Koenderink, "The shading cue in context," *i-Perception*, vol. 1, pp. 159–177, 2010.
- [12] M. K. Johnson and E. H. Adelson, "Shape estimation in natural illumination," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2011, pp. 2553–2560.
- [13] Q. Zhu and J. Shi, "Shape from shading: Recognizing the mountains through a global view," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recog.*, 2006, pp. 1839–1846.
- [14] A. Ecker and A. D. Jepson, "Polynomial shape from shading," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2010, pp. 145–152.
- [15] W. T. Freeman, E. C. Pasztor, and O. T. Carmichael, "Learning low-level vision," *Int. J. Comput. Vis.*, vol. 40, pp. 25–47, 2000.
- [16] T. Hassner and R. Basri, "Example based 3D reconstruction from single 2D images," in *Proc. CVPR Workshop "Beyond Patches"*, 2006, p. 15.

- [17] X. Huang, J. Gao, L. Wang, and R. Yang, "Exemplar-based shape from shading," in *Proc. 6th Int. Conf. 3-D Digital Imaging Model.*, 2007 pp. 349–356.
- [18] F. Cole, P. Isola, W. T. Freeman, F. Durand, and E. H. Adelson, "Shapecollage: Occlusion-aware, example-based shape interpretation," in *Proc. 12th Eur. Conf. Comput. Vis.*, 2012, pp. 665–678.
- [19] D. Marr, "Early processing of visual information," *Philosophical Trans. Roy. Soc. London. B, Biological Sci.*, vol. 275, pp. 483–519, 1976.
- [20] P. S. Huggins, H. F. Chen, P. N. Belhumeur, and S. W. Zucker, "Finding folds: On the appearance and identification of occlusion," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recog.*, 2001, pp. II-718–II-725.
- [21] J. T. Barron and J. Malik, "Shape, illumination, and reflectance from shading," EECS Dept., Univ. California, Berkeley, CA, USA, Tech. Rep. UCB/EECS-2013-117, 2013.
- [22] B. Kunsberg and S. W. Zucker, "Characterizing ambiguity in light source invariant shape from shading," arXiv:1306.5480v1 [cs.CV], Jun. 2013.
- [23] P. Tan, L. Quan, and T. Zickler, "The geometry of reflectance symmetries," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 33, no. 12, pp. 2506–2520, Dec. 2011.
- [24] D. W. Marquardt, "An algorithm for least-squares estimation of nonlinear parameters," *J. Soc. Ind. Appl. Math.*, vol. 11, pp. 431–441, 1963.
- [25] R. T. Frankot and R. Chellappa, "A method for enforcing integrability in shape from shading algorithms," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 10, no. 4, pp. 439–451, Jul. 1988.
- [26] P. Beckmann and A. Spizzichino, *The Scattering of Electromagnetic Waves from Rough Surfaces*. New York, NY, USA: Pergamon, 1963.
- [27] R. Grosse, M. K. Johnson, E. H. Adelson, and W. T. Freeman, "Ground truth dataset and baseline evaluations for intrinsic image algorithms," in *Proc. Int. Conf. Comput. Vis.*, 2009, pp. 2335–2342.
- [28] J. Haddon and D. Forsyth, "Shape representations from shading primitives," in *Proc. 5th Eur. Conf. Comput. Vis.*, 1998, pp. 415–431.



Ying Xiong received the BEng degree from Tsinghua University in 2010, and the SM degree from Harvard University in 2012. He is currently working toward the PhD degree at Harvard University under the supervision of Professor Todd Zickler. He was also a research assistant at NEC-Labs America and a software engineering intern at Google Inc. His research interests are in developing physics-based algorithms and techniques for computer vision.



Ayan Chakrabarti received the BTech and MTech degrees in electrical engineering from the Indian Institute of Technology Madras, Chennai, India, in 2006, and the SM and PhD degrees in engineering sciences from Harvard University, Cambridge, MA, in 2008 and 2011, respectively. He is currently a postdoctoral fellow at Harvard University, where his research focuses on the development of statistical models and inference algorithms for applications in computer vision and computational photography.



Ronen Basri received the BSc degree in mathematics and computer science from Tel Aviv University in 1985 and the PhD degree from the Weizmann Institute of Science in 1991. From 1990 to 1992, he was a postdoctoral fellow at the Department of Brain and Cognitive Science and the Artificial Intelligence Laboratory, Massachusetts Institute of Technology. Since then, he has been affiliated with the Department of Computer Science and Applied Mathematics, Weizmann Institute of Science, where he currently holds the

position of a professor. His research has focused on the areas of object recognition, shape reconstruction, lighting analysis, and image segmentation. His work deals with the design of algorithms, analysis, and implications to human vision.



Steven J. Gortler received the BA degree in computer science and applied mathematics from Queens College/CUNY in 1988 and the PhD degree in computer science from Princeton University in 1995. For two years, he was a postdoctoral researcher at the University of Washington and Microsoft Research. Since 1996, he has been affiliated with Harvard University, where he currently holds the position of a professor. He received the National Science Foundation Career Awards and a Research Fellowship from the Alfred P. Sloan Foundation. In 2002, he was received the ACM SIGGRAPH Significant New Researcher Award. He is interested in computer graphics, computer vision, and a menagerie of geometric problems.



David W. Jacobs received the BA degree from Yale University in 1982. He is a professor at the Department of Computer Science, University of Maryland with a joint appointment in the University's Institute for Advanced Computer Studies. From 1985 to 1992, he attended MIT, where he received the MS and PhD degrees in computer science. From 1982 to 1985, he was at Control Data Corporation on the development of data base management systems, and attended graduate school in computer science at New York University. From 1992 to 2002, he was a research scientist and then a senior research scientist at the NEC Research Institute. In 1998, he spent a sabbatical at the Royal Institute of Technology (KTH) in Stockholm, and in 2008 spent a sabbatical at the Ecole normale supérieure de Cachan. In 2002, he joined the CS Department at the University of Maryland. His research has focused on human and computer vision, especially in the areas of object recognition and perceptual organization. He has also published articles in the areas of motion understanding, memory and learning, computer graphics, human computer interaction, and computational geometry. He has served as an associate editor of the *IEEE Transactions on Pattern Analysis and Machine Intelligence*, and has assisted in the organization of many workshops and conferences, including serving as Program co-chair for CVPR 2010. He and his co-authors received honorable mention for the Best Paper Award at CVPR 2000. He also co-authored a paper that received the best student paper award at UIST 2003. With researchers at Columbia University and the Smithsonian Institution, he created Leafsnap, an app that uses computer vision for plant species identification. Leafsnap has been downloaded over a million times, and has been used in biodiversity studies and in many classrooms. He and his collaborators have been received the 2011 Edward O. Wilson Biodiversity Technology Pioneer Award for the development of Leafsnap.



Todd Zickler received the BEng degree in honors electrical engineering from McGill University in 1996 and the PhD degree in electrical engineering from Yale University in 2004. He joined the School of Engineering and Applied Sciences, Harvard University, as an assistant professor in 2004 and was appointed a professor of electrical engineering and computer science in 2011. He spent a sabbatical year at the Weizmann Institute of Science in 2012. He is the director of the Harvard Computer Vision Laboratory and member of the Graphics, Vision and Interaction Group at Harvard. His research is focused on modeling the interaction between light and materials, and developing systems to extract scene information from visual data. His work is motivated by applications in face, object, and scene recognition; image-based rendering; image retrieval; image and video compression; robotics; and human-computer interfaces. He received the National Science Foundation Career Award and a Research Fellowship from the Alfred P. Sloan Foundation.

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