

# Real-Time Focus Range Sensor

Shree K. Nayar, *Member, IEEE Computer Society*,  
Masahiro Watanabe, *Member, IEEE Computer Society*,  
and Minori Noguchi, *Member, IEEE Computer Society*

**Abstract**—Structures of dynamic scenes can only be recovered using a real-time range sensor. Depth from defocus offers an effective solution to fast and dense range estimation. However, accurate depth estimation requires theoretical and practical solutions to a variety of problems including recovery of textureless surfaces, precise blur estimation, and magnification variations caused by defocusing. Both textured and textureless surfaces are recovered using an illumination pattern that is projected via the same optical path used to acquire images. The illumination pattern is optimized to maximize accuracy and spatial resolution in computed depth. The relative blurring in two images is computed using a narrow-band linear operator that is designed by considering all the optical, sensing, and computational elements of the depth from defocus system. Defocus invariant magnification is achieved by the use of an additional aperture in the imaging optics. A prototype focus range sensor has been developed that has a workspace of 1 cubic foot and produces up to  $512 \times 480$  depth estimates at 30 Hz with an average RMS error of 0.2%. Several experimental results are included to demonstrate the performance of the sensor.

**Index Terms**—Depth from defocus, constant magnification defocusing, active illumination pattern, optical transfer function, image sensing, tuned focus operator, depth estimation, real-time range sensor.

## 1 INTRODUCTION

A pertinent problem in computational vision is the recovery of three-dimensional scene structure from two-dimensional images. Of all problems studied in vision, the above has by far attracted the most attention. This has resulted in a panoply of sensors and algorithms [8], [17] that can be broadly classified into two categories; passive and active. Passive techniques such as shape from shading and shape from texture attempt to extract structure from a single image. These algorithms are still under investigation and, given the assumptions they are forced to invoke, it is expected that they will prove complementary to other techniques but not serve as stand-alone strategies. Other passive methods such as stereo and structure from motion use multiple views to resolve shape ambiguities inherent in a single image. The primary bottleneck for these methods has proved to be correspondence and feature tracking. Recently, a large parallel architecture was developed to compute real-time depth maps using stereo [26].

The most popular range sensors in use today are based on time of flight or light striping [8]. In structured environments, where active radiation of a scene is feasible, light stripe methods offer a robust yet inexpensive solution to depth estimation. However, they have suffered from one inherent drawback, namely, speed. To achieve depth maps with sufficient spatial resolution, a large number (say,  $N$ ) of closely spaced stripes are used. If all stripes are projected

simultaneously it is impossible to associate a unique stripe with any given image point, a process that is necessary to compute depth by triangulation. The classical approach is to obtain  $N$  images, one for each stripe. If  $T_f$  is the time required to sense and digitize an image, the scanning of  $N$  stripes takes at least  $N \cdot T_f$ . Substantial improvements can be made by assigning gray codes to the stripes and scanning the entire collection of stripes in sets [23]. All the information needed is then acquired in  $\log_2(N)T_f$ , a significant improvement. An alternative approach uses color-coded stripe patterns [9]; this however is practical only in a gray-world that reflects all wavelengths of light. New hope for light stripe range finding has been instilled by advances in VLSI. Based on the notion of cell parallelism [25], a computational sensor has been developed where each sensor element records a stripe detection time-stamp as a single laser stripe sweeps the scene at high speed. Depth maps are produced in as little as one msec, though present day silicon packaging limits the total number of cells, and hence spatial depth resolution, to  $28 \times 32$  [24]. Future advances in VLSI are expected to yield high-resolution depth maps at unprecedented speeds.

In this paper, we present a range sensor based on focus analysis that produces a  $512 \times 480$  depth map at 30 Hz (video frame-rate). The sensor uses inexpensive off-the-shelf imaging and processing hardware and is shown to have an accuracy of approximately 0.2% [19]. Focus analysis has a major advantage over stereo and structure from motion; two or more images of a scene are taken under different optical settings but from the same viewpoint, as initially demonstrated in [10], [13], [12]. This circumvents the need for correspondence or feature tracking. The algorithm presented here uses only two scene images. These images correspond to different levels of focus and local frequency

- S.K. Nayar is with the Department of Computer Science, Columbia University, New York, USA. E-mail: nayar@cs.columbia.edu.
- M. Watanabe and M. Noguchi are with the Production Engineering Research Lab., Hitachi Ltd., Totsuka, Japan.

Manuscript received Jan. 13, 1995; revised Oct. 1, 1996. Recommended for acceptance by K. Boyer.

For information on obtaining reprints of this article, please send e-mail to: transpami@computer.org, and reference IEEECS Log Number P96098.

analysis implemented typically via linear operators yields depth estimates. However, differences between the two images tend to be very subtle and we believe that previous solutions to depth from defocus have met with limited success as they are based on rough approximations to the optical and sensing mechanisms involved in focus analysis. In contrast, our approach is based on a detailed physical modeling of all the optical, sensing, and computational elements at work; the optical transfer function, defocus, image sensing and sampling, and focus measure operators.

Depth from defocus shares one inherent weakness with stereo and motion; it requires that the scene has high frequency textures. A textureless surface appears the same focused or defocused and resulting images do not contain information necessary for depth computation. This has prompted us to develop a focus range sensor that uses active illumination. The key idea is to force a texture on the scene and then analyze the relative defocus of the texture in two images. Illumination projection has been suggested in the past [7], [11] for both depth from defocus and depth from pattern size distortion under perspective projection. However, these projected patterns were selected in an ad hoc fashion and do not guarantee high precision in computed depth. A critical problem therefore is determining an illumination pattern that would maximize the accuracy and robustness of depth from defocus. In this paper, a solution to this problem is arrived at through a detailed Fourier analysis of the entire depth from defocus system. First, theoretical models developed for each of the optical and computational elements of the system are expressed in spatial and Fourier domains. The derivation of the illumination pattern (or filter) is then posed as an optimization problem in Fourier domain. The optimal pattern is one that maximizes sensitivity of the focus measure to depth variations while minimizing the size of the focus operator to achieve high spatial resolution in computed depth.

A prototype real-time focus range sensor has been developed. It uses two CCD image detectors that view the scene through the same optical elements. The derived illumination pattern is fabricated using micro-lithography and incorporated into the sensor. The illumination pattern is projected onto the scene via the same optical path used to image the scene. This results in several advantages. It enables precise registration of the illumination pattern with the sampling grid of the image sensors. Light rays projected out through the imaging optics are subjected to similar geometric distortions as rays reflected back to the sensors. Therefore, despite ever-present lens distortions, the illumination pattern and the sensing grid of the detector are well registered. The coaxial illumination and imaging also results in a shadowless image; all surface regions that are visible to the sensor are also illuminated. Furthermore, since both images are acquired from the same viewing direction, the missing part or occlusion problem in stereo is avoided. Fig. 1 shows two brightness images and the computed depth map of a cup with milk flowing out of it. Structures of such dynamic scenes can only be recovered by a high-speed sensor. Several experiments have been conducted to evaluate the accuracy and real-time capability of the sensor.

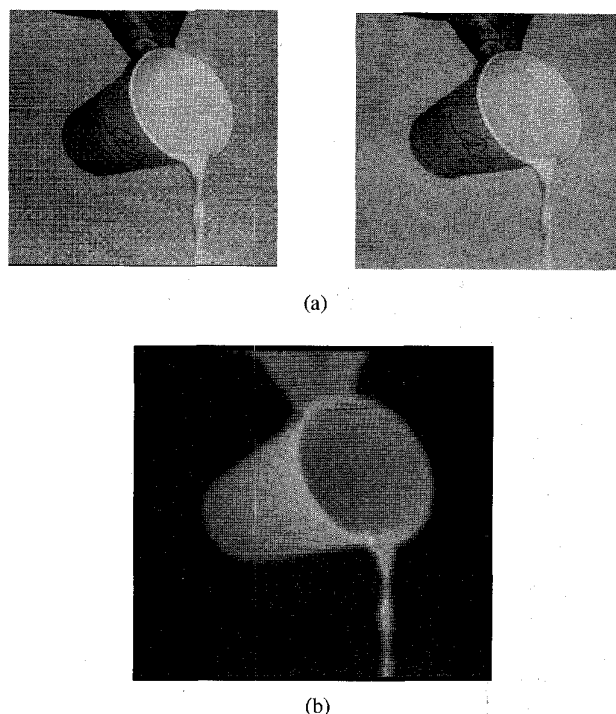


Fig. 1. a. Two images of a scene taken simultaneously using different focus settings; b. A depth map of the scene computed in 33 msec by the focus range sensor.

## 2 DEPTH FROM DEFOCUS

Fundamental to depth from defocus is the relationship between focused and defocused images [1]. Fig. 2 shows the basic image formation geometry. All light rays that are radiated by object point  $P$  and pass the aperture  $A$  are refracted by the lens to converge at point  $Q$  on the image plane. For a thin lens, the relationship between the object distance  $d$ , focal length of the lens  $f$ , and the image distance  $d_i$  is given by the Gaussian lens law:

$$\frac{1}{d} + \frac{1}{d_i} = \frac{1}{f}. \quad (1)$$

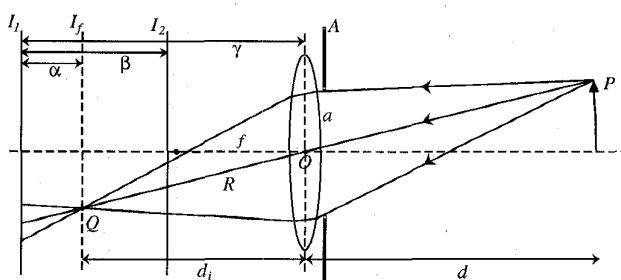


Fig. 2. Image formation and depth from defocus.

Each point on the object plane is projected onto a single point on the image plane, causing a clear or *focused* image  $I_f$  to be formed. If, however, the sensor plane does not coincide with the image plane and is displaced from it, the energy received from  $P$  by the lens is distributed over a patch

on the sensor plane. The result is a blurred image of  $P$ . It is clear that a single image does not include sufficient information for depth estimation as two scenes defocused to different degrees can produce identical images. A solution to depth is achieved by using two images  $I_1$  and  $I_2$  separated by a known physical distance  $\beta$  [10], [15]. The problem is reduced to analyzing the relative blurring of each scene point in the two images and computing the distance  $\alpha$  of its focused image. Then, using  $d_i = \gamma - \alpha$ , the lens law (1) yields depth  $d$  of the scene point. Simple as this procedure may appear, several technical problems emerge when implementing an algorithm of practical value.

- **Determining Relative Defocus:** In frequency domain, blurring can be viewed as low-pass filtering of scene texture. Relative blurring can thus in principle be estimated by frequency analysis. This problem is non-trivial since local scene texture includes frequencies with unknown magnitudes and phases. Since the effect of blurring is frequency dependent, it is not meaningful to investigate the net blurring of the entire collection of frequencies that constitute scene texture. This observation has forced investigators to use narrow-band filters that isolate more or less single frequencies and estimate their relative attenuation due to defocus in two or more images. Given that the dominant frequencies of the scene are unknown and possibly spatially varying, one is forced to use a large bank of tuned filters such as Gabor filters [14], [18] or hypergeometric filters [16]. Three problems surface with this approach.
  - 1) While it is rigorous, the necessity to use scores (at times, more than 100) filters makes it impractical for any real-time application without the use of expensive customized hardware.
  - 2) The filters are typically chosen by assuming the images to be continuous. Filter design for discrete images requires that the analysis be carried further to avoid undesirable artifacts in computed depth.
  - 3) Irrespective of the reliability of a filter in extracting focus measures, its output can be put to good use only if all optical and sensing elements of the depth from defocus system are accurately modeled. For instance, previous work has relied heavily on the Gaussian blur function, an approximation that may suffice for depth from focus<sup>1</sup> but limits the accuracy of depth from defocus.
- **Textureless Surfaces:** Depth from defocus shares a major weakness with stereo and structure from motion. If the imaged surface is textureless (a white sheet of paper, for instance) defocus and focus produce identical images and any number of filters would

1. All work in focus based depth computation can be broadly classified into depth from focus and depth from defocus. The former relies on a large number of images taken by varying  $\alpha$  in Fig. 2 in small increments (or through search) and uses a focus operator to detect the image of maximum focus for each scene point (see [16], [19], [27], [28], [29], [30], [31], [32]). In contrast, depth from defocus typically uses two images and estimates relative blurring to get depth (see [10], [12], [16], [13], [33], [34], [35], [36]).

prove ineffective in estimating relative blurring. A similar situation would arise in stereo or motion; correspondence and feature tracking would be ill-posed. Particularly in structured environments, this problem can be obviated by projecting an illumination pattern on the scene of interest, i.e., forcing scene texture [7], [10]. However, careful attention must be given to the pattern that is used, else the problem is at best reduced to applying depth from defocus to a scene with unknown texture. In our work we are interested in both textured and textureless scenes and hence adopt illumination projection. In contrast to previous work, however, we seek an optimal pattern that would ensure that all scene points have the same dominant texture, one that maximizes the spatial resolution and accuracy of computed depth. Derivation of the optimal projected pattern is posed as an optimization in Fourier domain.

- **Varying Magnification:** Lastly, the relation between magnification and focus is worth mentioning. In the imaging system shown in Fig. 2, the effective image location of point  $p$  moves along ray  $R$  as the sensor plane is displaced. This causes a shift in image coordinates of  $P$ . This variation in image magnification with defocus manifests as a mild correspondence-like problem in depth from defocus as the right set of points in images  $I_1$  and  $I_2$  are needed to estimate blurring. This problem has been underemphasized in previous work with the exception of [37] where a precise focus-magnification calibration of motorized zoom lenses is suggested and [28] where a registration-like correction in image domain is proposed. The calibration approach, while effective, is cumbersome and not viable for many off-the-shelf lenses. We use a simple but effective solution that is based on first principles of optics.

### 3 CONSTANT MAGNIFICATION DEFOCUS

We begin with the last of the problems raised in the above discussion; the variation of image magnification with defocus. We approach the problem from an optical perspective rather a computational one. Consider the image formation model shown in Fig. 3. The only modification made with respect to the model in Fig. 2 is the use of the external aperture  $A'$ . The aperture is placed at the *front-focal plane*, i.e., a focal length in front of the *principal point*  $O$  of the lens. This simple addition solves the prevalent problem of magnification variation with distance  $\alpha$  of the sensor plane from the lens. Simple geometrical analysis reveals that a ray of light  $R'$  from any scene point that passes through the center  $O'$  of aperture  $A'$  emerges parallel to the optical axis on the image side of the lens [2]. Furthermore, this parallel ray is the axis of a cone that includes all light rays radiated by the scene point, passed through by  $A'$ , and intercepted by the lens. As a result, despite blurring, the effective image coordinates of point  $P$  in both images  $I_1$  and  $I_2$  are the same, namely, the coordinate of its focused image  $Q$  on  $I_f$ .

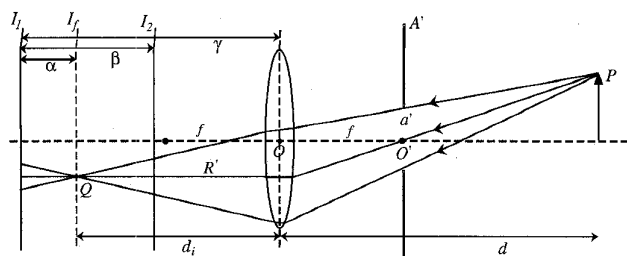


Fig. 3. Constant-magnification imaging system for depth from defocus is achieved by simply placing an aperture at the front-focal plane of the optics [2], [21].

This invariance of magnification to defocus holds true for any depth from defocus configuration (all values of  $\alpha$  and  $\beta$ ). It can also be shown that the constant-magnification property is unaffected by the aperture radius  $a'$  used. Furthermore, the lens law of (1) remains valid. The above optical configuration, called *telecentric*, is feasible not only in single lens systems but any compound lens system (see [2], [21]). Given an off-the-shelf lens, an aperture can be appended to the casing of the lens using the procedure described in [21]. While the nominal and effective *F-numbers* of the classical optics in Fig. 2 are  $f/a$  and  $d_i/a$ , respectively, they are both equal to  $f/a'$  in the telecentric case.

#### 4 MODELING

Effective solutions to both illumination projection and depth estimation require careful modeling and analysis of all physical phenomena involved in depth from defocus. The following five components are critical and will be modeled in this section:

- illumination pattern,
- optical transfer function,
- defocusing,
- image sensing, and
- focus operator.

These components, together, determine the relation between the depth of a scene point and its two focus measures. Since we have used the telecentric lens (Fig. 3) in our implementation, its parameters are used in developing each model. However, all of the following expressions can be made valid for a conventional lens system (Fig. 2) by simply replacing the factor  $\frac{f}{a'}$  by  $\frac{d_i}{a}$ .

##### 4.1 Illumination Pattern

Before the parameters of the illumination pattern can be determined, an illumination model must be defined. Such a model must be flexible in that it must subsume a large enough variety of possible illumination patterns. In defining the model, it is meaningful to take the characteristics of the other components into consideration. As we will describe shortly, the image sensor used has rectangular pixels arranged on a rectangular spatial grid. With this in mind, we define the following illumination model. The basic building block of the model is a rectangular illuminated patch, or cell, with uniform intensity:

$$i_c(x, y) = i_c(x, y; b_x, b_y) = {}^2\Pi\left(\frac{1}{b_x}x, \frac{1}{b_y}y\right) \quad (2)$$

where,  ${}^2\Pi()$  is the two-dimensional *Rectangular* function [3]. The *unknown* parameters of this illumination cell are  $b_x$  and  $b_y$ , the length and width of the cell.

This cell is assumed to be repeated on a two-dimensional grid to obtain a periodic pattern. This periodicity is essential since our goal is to achieve spatial invariance in depth accuracy, i.e., all image regions, irrespective of their distance from each other, must possess the same textural characteristics. The periodic grid is defined as:

$$i_g(x, y) = i_g(x, y; t_x, t_y) = {}^2\Pi\left(\frac{1}{2}\left(\frac{1}{t_x}x + \frac{1}{t_y}y\right), \frac{1}{2}\left(\frac{1}{t_x}x - \frac{1}{t_y}y\right)\right) \quad (3)$$

where,  ${}^2\Pi()$  is the two-dimensional *Shah* function [3], and  $2t_x$  and  $2t_y$  determine the periods of the grid in the  $x$  and  $y$  directions. Note that this grid is not rectangular but has vertical and horizontal symmetry on the  $x$ - $y$  plane. The final illumination pattern  $i(x, y)$  is obtained by convolving the cell  $i_c(x, y)$  with the grid  $i_g(x, y)$ :

$$i(x, y) = i(x, y; b_x, b_y, t_x, t_y) = i_c(x, y) * i_g(x, y) \quad (4)$$

The exact pattern is therefore determined by four parameters, namely,  $b_x, b_y, t_x$  and  $t_y$ . The above illumination grid is not as restrictive as it may appear upon initial inspection. For instance, the parameters  $b_x, b_y, 2t_x$  and  $2t_y$  can each be stretched to obtain repeated illumination and non-illumination stripes in the horizontal and vertical directions, respectively. Alternatively, they can also be adjusted to obtain a checkerboard illumination pattern with large or small illuminated patches. The exact values for  $b_x, b_y, t_x$  and  $t_y$  will be evaluated by the optimization procedure described later. In practice, the illumination pattern determined by the optimization is used to fabricate a filter with the same pattern.

The optimization procedure requires the analysis of each component of the system in spatial domain as well as frequency domain ( $u, v$ ). The Fourier transforms of the illumination cell, grid, and pattern are denoted as  $I_c(u, v), I_g(u, v)$ , and  $I(u, v)$ , respectively, and found to be:

$$I_c(u, v) = I_c(u, v; b_x, b_y) = b_x \frac{\sin(\pi b_x u)}{\pi b_x u} \cdot b_y \frac{\sin(\pi b_y v)}{\pi b_y v} \quad (5)$$

$$I_g(u, v) = I_g(u, v; t_x, t_y) = {}^2\Pi\left((t_x u + t_y v), (t_x u - t_y v)\right) \quad (6)$$

$$I(u, v) = I(u, v; b_x, b_y, t_x, t_y) = I_c(u, v) \cdot I_g(u, v) \quad (7)$$

##### 4.2 Optical Transfer Function

Adjacent points on the viewed surface reflect light waves that interfere with each other to produce diffraction effects. The angle of diffraction increases with the spatial frequency of surface texture. Since the lens aperture of the imaging system (Fig. 3) is of finite radius  $a'$ , it does not capture the higher order diffractions radiated by the surface (see [1] for details). This effect places a limit on the optical resolution of the imaging system characterized by the optical transfer function (OTF):

$$O(u, v; a', f) = \begin{cases} \left(\frac{a'}{f}\right)^2 (\gamma - \sin \gamma), & \sqrt{u^2 + v^2} \leq \frac{2a'}{\lambda f} \\ 0, & \sqrt{u^2 + v^2} > \frac{2a'}{\lambda f} \end{cases} \quad (8)$$

$$\text{where } \gamma = 2 \cos^{-1} \left( \frac{\lambda f \sqrt{u^2 + v^2}}{2a'} \right)$$

where  $(u, v)$  is the spatial frequency of the two-dimensional surface texture as seen from the image side of the lens,  $f$  is the focal length of the lens, and  $\lambda$  is the wavelength of incident light. It is clear from the above expression that only spatial frequencies below the limit  $\frac{2a'}{\lambda f}$  will be imaged by the optical system (Fig. 4). This in turn places restrictions on the frequency of the illumination pattern. Further, the above frequency limit can be used to "cut off" any desired number of higher harmonics produced by the illumination pattern.

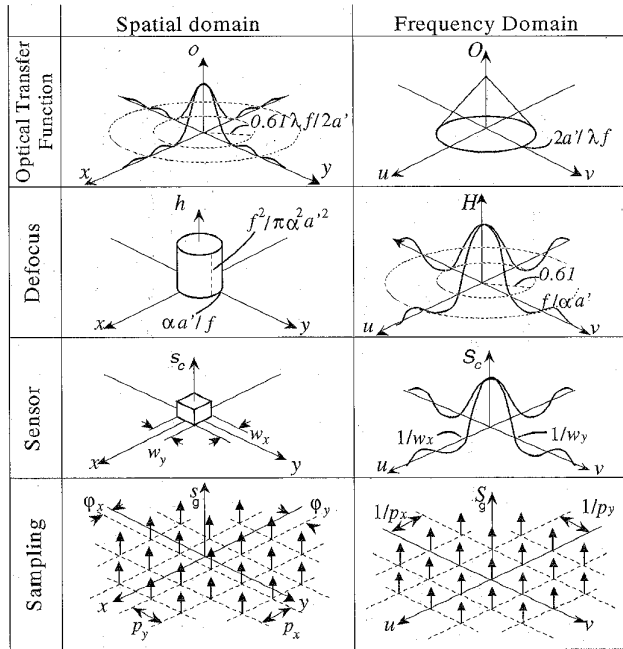


Fig. 4. Spatial and frequency models for the optical and sensing elements of depth from defocus.

### 4.3 Defocusing

The defocus function is described in detail in previous work (see [1], [4], for example). As in Fig. 3, let  $\alpha$  be the distance between the focused image of a surface point and its defocused image formed on the sensor plane. The light energy radiated by the surface point and collected by the imaging optics is uniformly distributed over a circular patch on the sensor plane. This patch, also called the *pillbox*, is the defocus function (Fig. 4):

$$h(x, y) = h(x, y; \alpha, a', f) = \frac{f^2}{2\pi a'^2 \alpha^2} \Pi \left( \frac{d}{2a\alpha} \sqrt{x^2 + y^2} \right) \quad (9)$$

where, once again,  $a'$  is the radius of the telecentric lens

aperture. In Fourier domain, the above defocus function is given by:

$$\begin{aligned} H(u, v) &= H(u, v; \alpha, a', f) \\ &= \frac{f}{2\pi a' \alpha \sqrt{u^2 + v^2}} J_1 \left( \frac{2\pi a' \alpha}{f} \sqrt{u^2 + v^2} \right) \end{aligned} \quad (10)$$

where  $J_1$  is the first-order Bessel function [1]. As is evident from the above expression, defocus serves as a low-pass filter. The bandwidth of the filter increases as  $\alpha$  decreases, i.e., as the sensor plane gets closer to the plane of focus. In the extreme case of  $\alpha = 0$ ,  $H(u, v)$  passes all frequencies without attenuation producing a perfectly focused image. Note that in a defocused image, all frequencies are attenuated at the same time. In the case of passive depth from focus or defocus, this poses a serious problem; different frequencies in an unknown scene are bound to have different (and unknown) magnitudes and phases. This again indicates that it would be desirable to have an illumination pattern that has a single dominant frequency, enabling robust estimation of defocus and hence depth.

### 4.4 Image Sensing

We assume the image sensor to be a typical CCD TV camera. Such a sensor can be modeled as a rectangular array of rectangular sensing elements (pixels). The quantum efficiency [5] of each pixel is assumed to be uniform over the area of the pixel. Let  $m(x, y)$  be the continuous image formed on the sensor plane. The finite pixel area has the effect of averaging the continuous image  $m(x, y)$ . In spatial domain, the averaging function is the rectangular cell:

$$s_c(x, y) = s_c(x, y; w_x, w_y) = \Pi \left( \frac{1}{w_x} x, \frac{1}{w_y} y \right) \quad (11)$$

where  $w_x$  and  $w_y$  are the length and width of the pixel, respectively. The discrete image is obtained by sampling the convolution of  $m(x, y)$  with  $s_c(x, y)$ . This sampling function is a rectangular grid:

$$\begin{aligned} s_g(x, y) &= s_g(x, y; p_x, p_y, \varphi_x, \varphi_y) \\ &= \frac{1}{p_x p_y} \Pi \left( \frac{1}{p_x} (x - \varphi_x), \frac{1}{p_y} (y - \varphi_y) \right) \end{aligned} \quad (12)$$

where  $p_x$  and  $p_y$  are spacings between discrete samples in the two spatial dimensions, and  $(\varphi_x, \varphi_y)$  is phase shift of the grid. The final discrete image is therefore:

$$m_d(x, y) = (s_c(x, y) * m(x, y)) \cdot s_g(x, y) \quad (13)$$

The parameters  $w_x, w_y, p_x,$  and  $p_y$  are all determined by the particular image sensor used. These parameters are therefore known and their values are substituted after the optimization is done. On the other hand, the phase shift  $(\varphi_x, \varphi_y)$  of the sampling function is with respect to the illumination pattern and will also be viewed as illumination parameters during optimization.

In Fourier domain, the above averaging and sampling functions are:

$$\begin{aligned}
 S_c(u, v) &= S_c(u, v; w_x, w_y) \\
 &= w_x \frac{\sin(\pi w_x u)}{\pi w_x u} \cdot w_y \frac{\sin(\pi w_y v)}{\pi w_y v} \quad (14)
 \end{aligned}$$

$$\begin{aligned}
 S_g(u, v) &= S_g(u, v; p_x, p_y, \varphi_x, \varphi_y) \\
 &= \text{III}(p_x u, p_y v) \cdot e^{-i2\pi(\varphi_x u + \varphi_y v)} \quad (15)
 \end{aligned}$$

The final discrete image is:

$$M_d(u, v) = (S_c(u, v) \cdot M(u, v)) * S_g(u, v) \quad (16)$$

#### 4.5 Focus Operator

Since defocusing has the effect of suppressing high-frequency components in the focused image, it is desirable that the focus operator respond to high frequencies in the image. For the purpose of illumination optimization, we use the Laplacian. However, the derived pattern will remain optimal for a large class of symmetric focus operators. In spatial domain, the discrete Laplacian is:

$$\begin{aligned}
 l(x, y) &= l(x, y; q_x, q_y) \\
 &= 4\delta x \cdot \delta y - [\delta x \cdot \delta y - q_y) + \delta x) \cdot \delta y + q_y) \\
 &\quad + \delta x - q_x) \cdot \delta y) + \delta x + q_x) \cdot \delta y)] \quad (17)
 \end{aligned}$$

Here,  $q_x$  and  $q_y$  are the spacings between neighboring elements of the discrete Laplacian kernel. In the optimization, these spacings will be related to the illumination parameters. The Fourier transform of the discrete Laplacian is:

$$\begin{aligned}
 L(u, v) &= L(u, v; q_x, q_y) \\
 &= 2(1 - \cos(2\pi q_x u)) * \delta u + 2(1 - \cos(2\pi q_y v)) * \delta v \quad (18) \\
 &= 4 - 2\cos(2\pi q_x u) - 2\cos(2\pi q_y v)
 \end{aligned}$$

The required discrete nature of the focus operator comes with a price. It tends to broaden the bandwidth of the operator. Once the pattern has been determined, the above filter will be tuned to maximize sensitivity to the fundamental illumination frequency while minimizing the effects of spurious frequencies caused either by the scene's inherent texture or image noise.

#### 4.6 Focus Measure

The focus measure is simply the output of the focus operator. It is related to defocus  $\alpha$  (and hence depth  $d$ ) via all of the components modeled above. Note that the illumination pattern ( $i_c * i_g$ ) is projected through optics that is similar to that used for image formation. Consequently, the pattern is also subjected to the limits imposed by the optical transfer function  $o$  and the defocus function  $h$ . Therefore, the texture projected on the scene is:

$$i(x, y; b_x, b_y, t_x, t_y) * o(x, y; a', f) * h'(x, y; \alpha', a', f) \quad (19)$$

where  $\alpha'$  represents defocus of the illumination itself that depends on the depth of the illuminated point. However, the illumination pattern once incident on a surface patch plays the role of surface texture and hence defocus  $\alpha'$  of illumination does not have any significant effect on depth estimation. The projected texture is reflected by the scene

and projected by the optics back onto the image plane to produce the discrete image:

$$\begin{aligned}
 &\{i(x, y; b_x, b_y, t_x, t_y) * o(x, y; a', f)^{*2} \\
 &\quad * h'(x, y; \alpha', a', f) * h(x, y; \alpha, a', f) \\
 &\quad * s_c(x, y; w_x, w_y)\} \cdot s_g(x, y; p_x, p_y, \varphi_x, \varphi_y) \quad (20)
 \end{aligned}$$

where  $o^{*2} = o * o$ . Strictly speaking, the above expression is valid only when all points in the image have the same depth (defocus). In general, depth varies over the scene and the resulting system is space-variant and thus cannot be expressed as a sequence of convolutions. However, for the purpose of discussion, we assume that depth is locally constant and hence the above expression is a valid local approximation for each point in the image.

The final focus measure function  $g(x, y)$  is the result of applying the discrete Laplacian to the above discrete image:

$$\begin{aligned}
 g(x, y) &= \{i(x, y; b_x, b_y, t_x, t_y) * o(x, y; a', f)^{*2} \\
 &\quad * h'(x, y; \alpha', a', f) * h'(x, y; \alpha, a', f) \\
 &\quad * s_c(x, y; w_x, w_y)\} s_g(x, y; p_x, p_y, \varphi_x, \varphi_y) \\
 &\quad * l(x, y; q_x, q_y) = \{(i * o^{*2} * h' * h * s_c) * l\} * s_g \quad (21)
 \end{aligned}$$

Since the distance between adjacent weights of the Laplacian kernel must be integer multiples of the period of the image sampling function  $s_g$ , the above expression can be rearranged as:

$$\begin{aligned}
 g(x, y) &= (i * o^{*2} * h' * h * s_c * l) \cdot s_g \\
 &= g_0 \cdot s_g \quad (22)
 \end{aligned}$$

where  $g_0 = i * o^{*2} * h' * h * s_c * l$ . The same can be expressed in Fourier domain as:

$$G(u, v) = (I \cdot O^2 \cdot H' \cdot H \cdot S_c \cdot L) * S_g = G_0 * S_g \quad (23)$$

The above expression gives us the final output of the focus operator for any value of the defocus parameter  $\alpha$ . It will be used in the following sections to determine the optimal illumination pattern and to estimate depth.

### 5 ILLUMINATION OPTIMIZATION

In our implementation, the illumination pattern is projected on the scene using a high power light source and a telecentric lens identical to the one used to image the scene. This allows us to assume that the projected illumination is the primary cause for surface texture and is stronger than the natural texture of the surface. Consequently, our results are applicable not only to textureless surfaces but also textured ones. The illumination optimization problem is formulated as follows: Establish closed-form relationships between the illumination parameters ( $b_x, b_y, t_x, t_y$ ), sensor parameters ( $w_x, w_y, p_x, p_y, \varphi_x, \varphi_y$ ), and discrete Laplacian parameters ( $q_x, q_y$ ) so as to maximize the sensitivity, robustness, and spatial resolution of the focus measure  $g(x, y)$ . High sensitivity implies that a small variation in the degree of focus results in a large variation in  $g(x, y)$ . This would ensure high depth estimation accuracy in the presence of image noise, i.e., high signal-to-noise ratio. By robustness we mean that all pixels with the same degree of defocus produce the same focus measure independent of their location on the image plane. This ensures that depth estimation accuracy is invariant to location on the image plane. Lastly, high spatial resolution

is achieved by minimizing the size of the focus operator. This ensures that rapid depth variations (surface discontinuities) can be detected with high accuracy.

In order to minimize smoothing effects and maximize spatial resolution of computed depth, the support (or span) of the discrete Laplacian must be as small as possible. This in turn requires the frequency of the illumination pattern be as high as possible. However, the optical transfer function described in Section 4.2 imposes limits on the highest frequency that can be imaged by the optical system. This maximum allowable frequency is  $\frac{2a'}{\lambda f}$ , determined by the numerical aperture of the telecentric lens. With this in mind, let us examine the Fourier transform of the illumination pattern. Since the pattern is periodic, its Fourier transform must be discrete. It may have a zero-frequency component, but this can be safely ignored since the Laplacian operator, being a sum of second-order derivatives, will eventually remove any zero-frequency component in the final image. Our objective then is to maximize the fundamental spatial frequency ( $1/t_x, 1/t_y$ ) of the illumination pattern. In order to maximize this frequency while maintaining high detectability, we must have

$$\sqrt{(1/t_x)^2 + (1/t_y)^2}$$

close to the optical limit  $\frac{2a'}{\lambda f}$ . This in turn pushes all higher frequencies in the illumination pattern outside the optical limit. What we are left with is a surface texture whose image has only the quadrupole fundamental frequencies ( $\pm 1/t_x, \pm 1/t_y$ ). As a result, these are the only frequencies we need consider in our analysis of the focus measure function  $G(u, v)$ .

Before we consider the final measure  $G(u, v)$ , we examine  $G_0(u, v)$  the focus measure prior to image sampling. For the reasons given above, the two-dimensional  $G_0(u, v)$  is reduced to four discrete spikes at  $(1/t_x, 1/t_y)$ ,  $(1/t_x, -1/t_y)$ ,  $(-1/t_x, 1/t_y)$ , and  $(-1/t_x, -1/t_y)$ . Since all components ( $I, O, H, S_c$ , and  $L$ ) of  $G_0$  are reflection symmetric about  $u = 0$  and  $v = 0$ , we have:

$$G_0\left(\frac{1}{t_x}, \frac{1}{t_y}\right) = G_0\left(\frac{1}{t_x}, -\frac{1}{t_y}\right) = G_0\left(-\frac{1}{t_x}, \frac{1}{t_y}\right) = G_0\left(-\frac{1}{t_x}, -\frac{1}{t_y}\right) \quad (24)$$

where:

$$\begin{aligned} G_0\left(\frac{1}{t_x}, \frac{1}{t_y}\right) &= I\left(\frac{1}{t_x}, \frac{1}{t_y}; b_x, b_y, t_x, t_y\right) \cdot O^2\left(\frac{1}{t_x}, \frac{1}{t_y}; a', f\right) \\ &\cdot H'\left(\frac{1}{t_x}, \frac{1}{t_y}; \alpha', a', f\right) \cdot H\left(\frac{1}{t_x}, \frac{1}{t_y}; \alpha, a', f\right) \\ &\cdot S_c\left(\frac{1}{t_x}, \frac{1}{t_y}; w_x, w_y\right) \cdot L\left(\frac{1}{t_x}, \frac{1}{t_y}; q_x, q_y\right). \quad (25) \end{aligned}$$

Therefore, in frequency domain the focus measure function prior to image sampling reduces to:

$$\begin{aligned} G_0(u, v) &= G_0\left(\frac{1}{t_x}, \frac{1}{t_y}\right) \\ &\cdot \left\{ \delta\left(u - \frac{1}{t_x}, v - \frac{1}{t_y}\right) + \delta\left(u + \frac{1}{t_x}, v - \frac{1}{t_y}\right) \right. \\ &\left. + \delta\left(u - \frac{1}{t_x}, v + \frac{1}{t_y}\right) + \delta\left(u + \frac{1}{t_x}, v + \frac{1}{t_y}\right) \right\} \quad (26) \end{aligned}$$

The function  $g_0(x, y)$  in image domain, is simply the inverse Fourier transform of  $G_0(u, v)$ :

$$g_0(x, y) = G_0\left(\frac{1}{t_x}, \frac{1}{t_y}\right) \cdot \left\{ 4 \cos 2\pi \frac{1}{t_x} x \cdot \cos 2\pi \frac{1}{t_y} y \right\} \quad (27)$$

Note that  $g_0(x, y)$  is the product of cosine functions weighted by the coefficient  $G_0(1/t_x, 1/t_y)$ . The defocus function  $h$  has the effect of reducing the coefficient  $G_0(1/t_x, 1/t_y)$  in the focus measure  $g_0(x, y)$ . Clearly, the sensitivity of the focus measure to depth (or defocus) is optimized by maximizing the coefficient  $G_0(1/t_x, 1/t_y)$  with respect to the unknown parameters of the system. This optimization procedure can be summarized as:

$$\frac{\partial}{\partial t_x} G_0\left(\frac{1}{t_x}, \frac{1}{t_y}\right) = 0, \quad \frac{\partial}{\partial t_y} G_0\left(\frac{1}{t_x}, \frac{1}{t_y}\right) = 0, \quad (28)$$

$$\frac{\partial}{\partial b_x} G_0\left(\frac{1}{t_x}, \frac{1}{t_y}\right) = 0, \quad \frac{\partial}{\partial b_y} G_0\left(\frac{1}{t_x}, \frac{1}{t_y}\right) = 0, \quad (29)$$

$$\frac{\partial}{\partial q_x} G_0\left(\frac{1}{t_x}, \frac{1}{t_y}\right) = 0, \quad \frac{\partial}{\partial q_y} G_0\left(\frac{1}{t_x}, \frac{1}{t_y}\right) = 0. \quad (30)$$

Since  $t_x$  and  $t_y$  show up in all the components in (25), the first two partial derivatives (28) are difficult to evaluate. Fortunately, the derivatives in (29) and (30) are sufficient to obtain relations between the system parameters. For details of the optimization procedure, we refer the reader to [20]. Maximum sensitivity and spatial resolution of the focus measure  $g(x, y)$  are achieved for the following illumination parameter values:

$$b_x = \frac{1}{2} t_x, \quad b_y = \frac{1}{2} t_y \quad (31)$$

$$q_x = \frac{1}{2} t_x, \quad q_y = \frac{1}{2} t_y \quad (32)$$

Next, we examine the spatial robustness of  $g(x, y)$ . Imagine the imaged surface to be planar and parallel to the image sensor. Then, we would like the image sampling to produce the same absolute value of  $g(x, y)$  at all discrete sampling points on the image. This entails relating the illumination and sensing parameters so as to facilitate careful sampling of the product of cosine functions in (27). Note that the final focus measure is:

$$\begin{aligned} g(x, y) &= g_0 \cdot s_g \\ &= G_0\left(\frac{1}{t_x}, \frac{1}{t_y}\right) \cdot \left\{ 4 \cos 2\pi \frac{1}{t_x} x \cdot \cos 2\pi \frac{1}{t_y} y \right\} \\ &\cdot {}^2\text{III}\left(\frac{1}{p_x}(x - \phi_x), \frac{1}{p_y}(y - \phi_y)\right) \quad (33) \end{aligned}$$

All samples of  $g(x, y)$  have the same *absolute* value when the two cosines in the above expression are sampled at their peak values. Such a sampling is possible when:

$$p_x = \frac{1}{2}t_x, \quad p_y = \frac{1}{2}t_y \tag{34}$$

and:

$$\varphi_x = 0, \quad \varphi_y = 0 \tag{35}$$

Alternatively, the cosines can be sampled with a period of  $\pi/2$  and phase shift of  $\pi/4$ . This yields the second solution:

$$p_x = \frac{1}{4}t_x, \quad p_y = \frac{1}{4}t_y, \tag{36}$$

$$\varphi_x = \pm \frac{1}{8}t_x, \quad \varphi_y = \pm \frac{1}{8}t_y. \tag{37}$$

The above equations give two solutions, both are checkerboard illumination patterns but differ in their fundamental frequencies, size of the illumination cell, and the phase shift with respect to the image sensor. Equations (31), (32), (34), and (35) yield the filter pattern shown in Fig. 5a. In this case the filter and detector are registered with zero phase shift, and the illumination cell has the same size and shape as the sensor elements (pixels). The second solution, shown in Fig. 5b, is obtained using the sampling solutions (36) and (37), yielding a filter pattern with illumination cell two times the size of the sensor element and phase shift of half the sensor element size. Exactly how such patterns can be projected and perfectly registered with the image detector will be described in the experimental section.

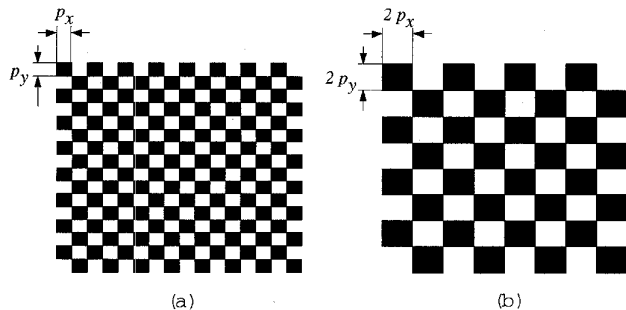


Fig. 5. Optimal illumination filter patterns: a.  $t_x = 2p_x, t_y = 2p_y, \varphi_x = 0, \varphi_y = 0$ ; and b.  $t_x = 4p_x, t_y = 4p_y, \varphi_x = 1/8 t_x, \varphi_y = 1/8 t_y$ . Hence,  $(t_x, t_y)$  is the illumination period,  $(p_x, p_y)$  is the pixel size, and  $(\varphi_x, \varphi_y)$  is the illumination phase shift with respect to the image sensing grid.

### 6 TUNED FOCUS OPERATOR

For the purpose of illumination optimization, we used the Laplacian operator. The resulting illumination pattern has only a single dominant absolute frequency,  $(1/t_x, 1/t_y)$ . Given this, we are in a position to further refine our focus operator so as to minimize the effects of all other frequencies caused either by the physical texture of the scene or image noise. To this end, let us consider the properties of the  $3 \times 3$  discrete Laplacian (see Figs. 6a and 6b). We see that though the Laplacian does have peaks exactly at  $(1/t_x, 1/t_y), (1/t_x, -1/t_y), (-1/t_x, 1/t_y),$  and  $(-1/t_x, -1/t_y)$ , it has a fairly broad bandwidth allowing other spurious frequencies to contribute to the focus measure  $G$  in (23). Here, we seek a

narrow-band operator with sharp peaks at the above four coordinates in frequency space.

Given that the operator must eventually be discrete and of finite support, there is a limit to the extent to which it can be tuned. To constrain the problem, we impose the following conditions:

- 1) To maximize spatial resolution in computed depth we force the operator kernel to be  $3 \times 3$ .
- 2) Since the fundamental frequency of the illumination pattern has a symmetric quadrupole arrangement, the focus operator must be rotationally symmetric. These two conditions force the operator to have the structure shown in Fig. 6c.
- 3) The operator must not respond to any DC component in image brightness.

This last condition is satisfied if the sum of all elements (see Fig. 6c) of the operator equals zero:

$$a + 4b + 4c = 0 \tag{38}$$

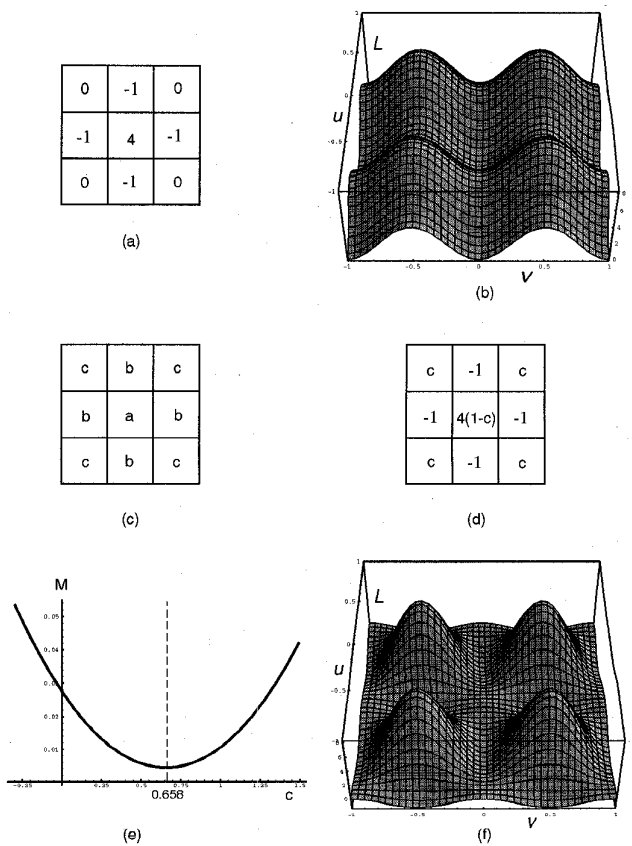


Fig. 6. a. The  $3 \times 3$  Laplacian and its b. Fourier transform; c. The kernel structure for a  $3 \times 3$  operator that is symmetric; d. The kernel of a  $3 \times 3$  operator that is insensitive to the zero frequency component (see text); e. The second moment  $M$  of each of the four operator peaks is minimized when  $c = 0.658$ ; f. Response of the tuned focus operator ( $c = 0.658$ ) has sharper peaks than the Laplacian.

It is also imperative that the response  $L(u, v)$  of the operator to the fundamental frequency not be zero:



$$L\left(\frac{1}{t_x}, \frac{1}{t_y}\right) = a + 2b \left( \cos 2\pi q_x \frac{1}{t_x} + \cos 2\pi q_y \frac{1}{t_y} \right) + 4c \cos 2\pi q_x \frac{1}{t_x} + \cos 2\pi q_y \frac{1}{t_y} \neq 0 \quad (39)$$

Given (32), the above reduces to:

$$a - 4b + 4c \neq 0 \quad (40)$$

Expressions (38) and (40) imply that  $b \neq 0$ . Without loss of generality, we set  $b = -1$ . Hence, (38) gives  $a = 4(1 - c)$ . Therefore, the tuned operator is determined by a single unknown parameter  $c$ , as shown in Fig. 6d. The problem then is to find  $c$  such that the operator's Fourier transform has a sharp peak at  $(1/t_x, 1/t_y)$ . A rough measure of sharpness is given by the second-order moment of the power  $\|L(u, v)\|^2$  with respect to  $(1/t_x, 1/t_y)$ :

$$M = \frac{1}{\left\|L\left(\frac{1}{t_x}, \frac{1}{t_y}\right)\right\|^2} \int_{u=0}^{\frac{2}{t_x}} \int_{v=0}^{\frac{2}{t_y}} \left[ \left(u - \frac{1}{t_x}\right)^2 + \left(v - \frac{1}{t_y}\right)^2 \right] \left\|L\left(u - \frac{1}{t_x}, v - \frac{1}{t_y}\right)\right\|^2 dv du = \frac{t_x^2 + t_y^2}{768\pi^2 t_x^3 t_y^3} (20\pi^2 c^2 + 6c^2 + 48c - 32\pi^2 c + 20\pi^2 - 93) \quad (41)$$

The above measure is minimized when  $\frac{\partial M}{\partial c} = 0$ , i.e., when  $c = 0.658$  as shown in Fig. 6e. The resulting tuned focus operator has the response shown in Fig. 6f, it has substantially sharper peaks than the discrete Laplacian. Given that the operator is  $3 \times 3$  and discrete, the sharpness of the peaks is limited. The above derivation brings to light the fundamental difference between designing tuned operators in continuous and discrete domains. In general, an operator that is deemed optimal in continuous domain is most likely suboptimal for discrete images. A quantitative comparison between the performance of the tuned and Laplacian operators will be presented in the section on experiments. A further refinement of the above tuned operator to make it insensitive to phase shifts in illumination is presented in [20].

## 7 DEPTH FROM TWO IMAGES

Depth estimation uses two images of the scene  $I_1(x, y)$  and  $I_2(x, y)$  that correspond to different effective focal lengths as shown in Fig. 3. Depth of each scene point is determined by estimating the displacement  $\alpha$  of the focused plane  $I_f$  for the scene point. The tuned focus operator is applied to both images to get focus measure images  $g_1(x, y)$  and  $g_2(x, y)$ . From (33) we see that:

$$\frac{g_1(x, y)}{g_2(x, y)} = \frac{G_0\left(\frac{1}{t_x}, \frac{1}{t_y}; \alpha\right)}{G_0\left(\frac{1}{t_x}, \frac{1}{t_y}; \alpha - \beta\right)} \quad (42)$$

From (23) we see that the only factor in  $G_0$  affected by parameter  $\alpha$  is defocus function  $H$ . Therefore:

$$\frac{g_1(x, y)}{g_2(x, y)} = \frac{H\left(\frac{1}{t_x}, \frac{1}{t_y}; \alpha\right)}{H\left(\frac{1}{t_x}, \frac{1}{t_y}; \alpha - \beta\right)} \quad (43)$$

Note that the above measure is not bounded. This poses a problem from a computational viewpoint which is easily remedied by using the following normalization:

$$q(x, y) = \frac{g_1(x, y) - g_2(x, y)}{g_1(x, y) + g_2(x, y)} = \frac{H\left(\frac{1}{t_x}, \frac{1}{t_y}; \alpha\right) - H\left(\frac{1}{t_x}, \frac{1}{t_y}; \alpha - \beta\right)}{H\left(\frac{1}{t_x}, \frac{1}{t_y}; \alpha\right) + H\left(\frac{1}{t_x}, \frac{1}{t_y}; \alpha - \beta\right)} \quad (44)$$

As shown in Fig. 7,  $q$  is a monotonic function of  $\alpha$  such that  $-p \leq q \leq p$ ,  $p \leq 1$ . In practice, the above relation can be pre-computed and stored as a look-up table that maps  $q$  computed at each image point to a unique  $\alpha$ . Since  $\alpha$  represents the position of the focused image, the lens law (1) yields the depth  $d$  of the corresponding scene point. Note that the tuned focus operator designed in the previous section is a linear filter, making it feasible to compute depth maps of scenes in real-time using off-the-shelf image processing hardware.

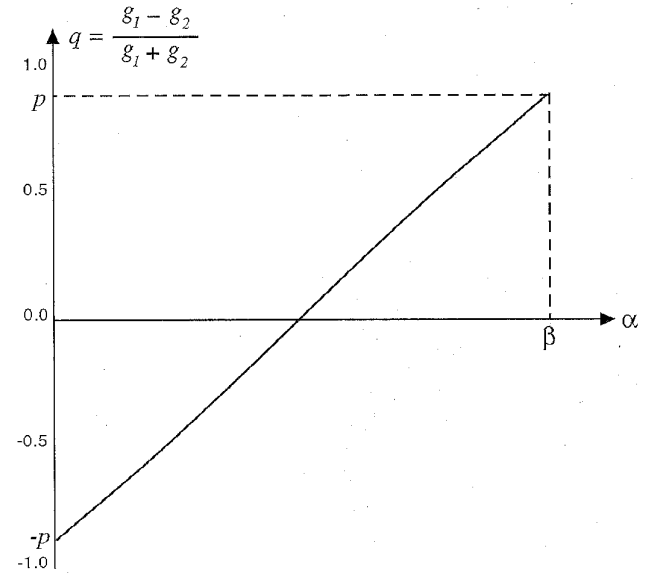


Fig. 7. Relation between focus measures  $g_1$  and  $g_2$  and the defocus parameter  $\alpha$ .

## 8 REAL TIME RANGE SENSOR

We have implemented the real-time focus range sensor shown in Fig. 8. The scene is imaged using a standard 12.5 mm Fujinon lens with an additional aperture added to convert it to telecentric (see [21] for details). Light rays passing through the lens are split in two directions using a beam-splitting prism. This produces two images that are simultaneously detected using two Sony XC-77RR eight bit CCD cameras. The positions of the two cameras are precisely fixed such that one obtains a near-focus image while the other a far-focus image. In this setup a physical displacement of 0.25 mm between the effective focal lengths of the two CCD cameras translates to a working range on the

object side of approximately 30 cm. This detectable range of the sensor can be varied either by changing the sensor displacement or the focal length of the imaging optics.

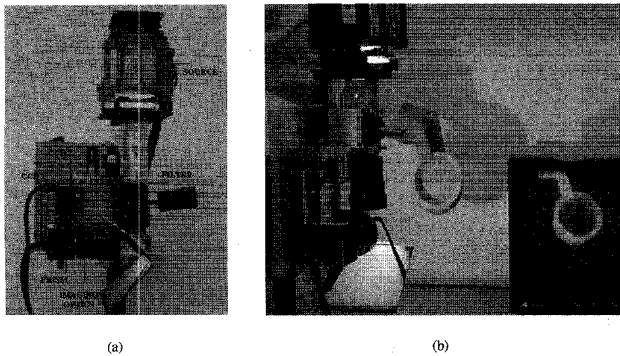


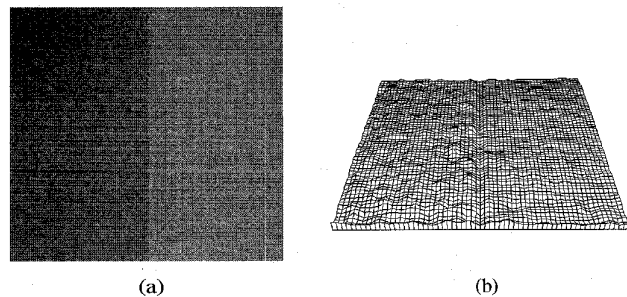
Fig. 8. a. The real-time focus range sensor and its key components. b. The sensor can produce depth maps up to  $512 \times 480$  in resolution at 30 Hz.

The illumination pattern shown in Fig. 5b was etched on a glass plate using microlithography, a process widely used in VLSI. The filter was then placed in the path of a 300 W Xenon arc lamp. The illumination pattern generated is projected using a telecentric lens identical to the one used for image capture. A half-mirror is used to ensure that the illumination pattern projects onto the scene via the same optical path used to acquire images. As a result, the pattern is almost perfectly registered with respect to the pixels of the two CCD cameras. Furthermore, a modification to the tuned focus operator of Section 6 is presented in [20] that makes it insensitive to slight misregistrations between the illumination pattern and image pixels. The above arrangement ensures that every scene point that is visible to the sensor is also illuminated by it, avoiding shadows and thus undetectable regions. If objects in the scene have a strong specular reflection component, cross-polarized filters can be attached to the illumination and imaging lens to filter out specularities and produce images that mainly include the diffuse reflection component.

Images from the two CCD cameras are digitized and processed using MV200 Datacube image processing hardware. The present configuration includes the equivalent of two eight bit digitizers, two A/D convertors, and one 12-bit convolver. This hardware enables simultaneous digitization of the two images, convolution of both images with the tuned focus operator, and the computation of a  $256 \times 240$  depth map, all within a single framerate of 33 msec with a lag of 33 msec. A look-up table is used to map each pair of focus measures ( $g_1$  and  $g_2$ ) to a unique depth estimate  $d$  (see [20] for implementation details). Alternatively, a  $512 \times 480$  depth map can be computed at the same rate if the two images are taken in succession. Simultaneous image acquisition is clearly advantageous since it makes the sensor less sensitive to variations in both illumination and scene structure between frames. With minor additions to the present processing hardware, it is easy to obtain  $512 \times 480$  depth maps at 30 Hz using simultaneous image grabbing. Depth maps produced by the sensor are visualized as wireframes at framerate using a DEC Alpha workstation.

## 9 EXPERIMENTS

Numerous experiments have been conducted to test the performance of the sensor. Here, we briefly summarize these results. Fig. 9a shows the near focused image of a planar surface, half of the surface is textureless while the other half has strong random texture. A computed depth map of the surface is shown in Fig. 9b. As expected the textureless area is estimated almost free of errors while the textured area has small errors due to texture frequencies that lie close to the illumination frequency. It may be noted that the texture used in this experiment includes a wide spectrum of frequencies. Most scenes have weaker textures and can be expected to produce even more accurate results. Several depth maps of the plane in Fig. 9a were computed by varying its position in the 30 cm workspace of the sensor and the average accuracy and repeatability of the sensor were estimated for both simultaneous and successive image grabbing configurations (see Table 9c). These results clearly demonstrate the superior performance of the sensor over previous implementations of depth from defocus. This improvement results from several factors including accurate modeling of sensor optics, the use of an optimized illumination pattern, and careful implementation of the sensor.



	Simultaneous Image Grab	Successive Image Grab
Depth Accuracy (rms)	0.24 %	0.34 %
Repeatability (rms)	0.23 %	0.29 %
Spatial Resolution	$256 \times 240$	$512 \times 480$
Speed	30 Hz	30 Hz
Delay	33 msec	33 msec

(c)

Fig. 9. a. Near focused image of a planar surface that includes highly textured and textureless areas; b. Depth of the surface computed using the focus range sensor; c. Performance characteristics of the sensor.

In Fig. 10, a simple example is shown to demonstrate that the tuned focus operator of Section 6 does outperform the Laplacian operator. Fig. 10a shows a planar surface with patches of different textures. Fig. 10b shows its depth map computed by the sensor when the focus operator used is a  $3 \times 3$  Laplacian. The broad frequency bands of the Laplacian allow some of the textures of the planar surface to induce large errors in computed depth. In contrast, the narrowband  $3 \times 3$  tuned operator produces a significantly improved depth map, as shown in Fig. 10c.

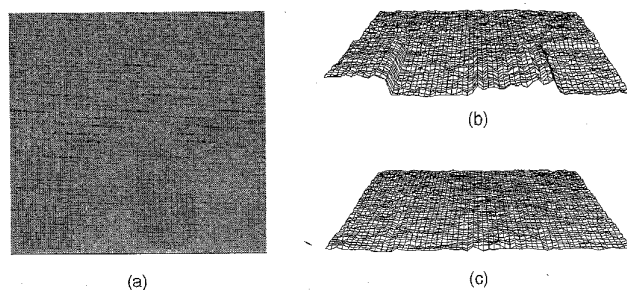


Fig. 10. a. Planar surface with textured patches; b. Depth map computed using a  $3 \times 3$  Laplacian focus operator; c. Depth map computed using the tuned focus operator. As expected, the broad frequency bands of the Laplacian make it more susceptible to depth errors induced by surface texture.

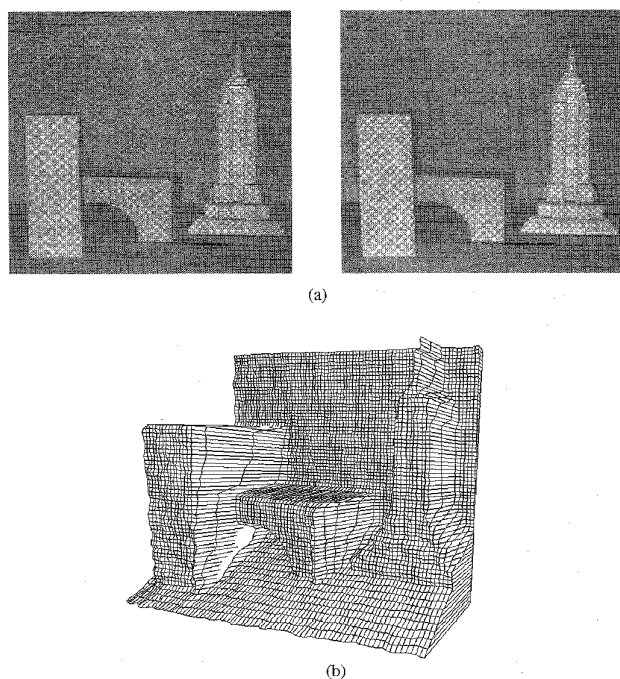


Fig. 11. a. Near and far focused images of a set of polyhedral objects; b. Computed depth map.

Fig. 11 shows a scene with polyhedral objects. The computed depth map in Fig. 11b is fairly accurate despite the complex textural properties of the objects. The only filtering that is applied to the depth map is a  $5 \times 5$  smoothing function to reduce high frequency noise in computed depth that results from the low signal-to-noise ratio of the CCD cameras and spurious frequencies caused by surface texture. All surface discontinuities and orientation discontinuities are well preserved. The recovered shapes are precise enough for a variety of visual tasks including recognition and inspection. In the case of dynamic scenes, structure can be estimated only by using a real-time sensor. Fig. 12 shows an object's depth map computed as it rotates on a motorized turntable. Such depth map sequences are useful in automatic CAD model generation from sample objects. Furthermore, real-time depth computation clearly enhances the capability of any vision system as it enables recovery of a deforming shape, precise tracking of moving objects, and robust navigation in dynamic scenes.

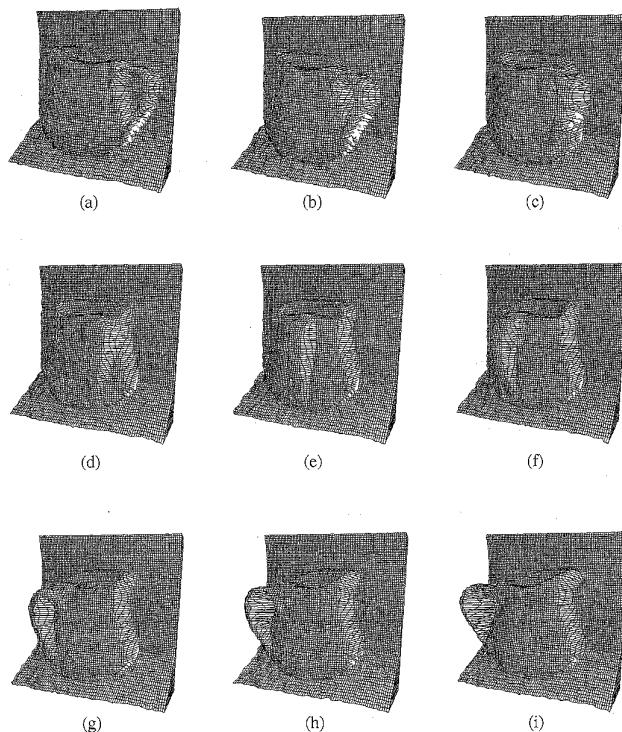


Fig. 12. Depth maps generated by the sensor at 30 Hz while an object rotates on a motorized turntable.

## 10 SUMMARY

We have reported theoretical results on a variety of issues related to depth estimation by focus analysis. Accurate modeling of optics and sensing were shown to be essential to precise depth estimation. Both textured and textureless surfaces are recovered by using an optimized illumination pattern that is registered with the image sensor. A telecentric optical configuration was used to achieve constant-magnification defocusing. All of these results were used to implement a real-time focus range sensor that produces high resolution depth maps at frame rate. This sensor is unique in its ability to produce fast, dense, and precise depth information at a very low cost. With time we expect the sensor to find applications ranging from visual recognition and robot control to automatic CAD model generation for vision and graphics. The obvious extension to this work is the development of a passive focus range finder for outdoor scenes. Such a sensor cannot afford the luxury of projected illumination. It must remain efficient while relying on the natural textures of scenes for depth estimation. Some progress in this regard has already been made [22].

## ACKNOWLEDGMENTS

This research was conducted at the Center for Research on Intelligent Systems, Department of Computer Science, Columbia University, New York, NY 10027, USA.

## REFERENCES

- [1] M. Born and E. Wolf, *Principles of Optics*. London: Pergamon, 1965.
- [2] R. Kingslake, *Optical System Design*. Academic Press, 1983.
- [3] R.N. Bracewell, *The Fourier Transform and Its Applications*. McGraw Hill, 1965.
- [4] B.K.P. Horn, *Robot Vision*. MIT Press, 1986.
- [5] B.K.P. Horn, *Focusing*, AI Lab, Memo 160. MIT, Cambridge, Mass, 1968.
- [6] T. Kanade, "Development of a Video-Rate Stereo Machine," *Proc. ARPA Image Understanding Workshop*, Nov. 1994, pp. 549-557.
- [7] B. Girod and S. Scherrock, "Depth from Focus of Structured Light," *Proc. SPIE: Optics, Illum., and Image Sng for Mach. Vis. IV*, vol. 1194, Nov. 1989, Philadelphia, Penn.
- [8] R.A. Jarvis, "A Perspective on Range Finding Techniques for Computer Vision," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 5, no. 2, pp. 122-139, March 1983.
- [9] K. L. Boyer and A. C. Kak, "Color-Encoded Structured Light for Rapid Active Sensing," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 9, no. 1, pp. 14-28, Jan. 1987.
- [10] A. Pentland, "A New Sense for Depth of Field," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 9, no. 4, pp. 523-531, July 1987.
- [11] A. Pentland, S. Scherrock, T. Darrell, and B. Girod, "Simple Range Cameras Based on Focal Error," *J. Optical Society of America*, vol. 11, no. 11, pp. 2,925-2,935, Nov. 1994.
- [12] A. Pentland, T. Darrell, M. Turk, and W. Huang, "A Simple, Real-Time Range Camera," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pp. 256-261, June 1989, San Diego, Calif.
- [13] M. Subbarao, "Parallel Depth Recovery by Changing Camera Parameters," *Proc. Int'l Conf. Computer Vision*, pp. 149-155, Dec. 1988.
- [14] M. Gokstorp, "Computing Depth from Out-of-Focus Blur Using a Local Frequency Representation," *Proc. Int'l Conf. Pattern Recognition*, Oct. 1994.
- [15] M. Subbarao and G. Surya, "Application of Spatial-Domian Convolution/Deconvolution Transform for Determining Distance from Image Defocus," *Proc. OE/BOSTON '92 SPIE Conf.*, vol. 1,822, Boston, Mass., Nov. 1992.
- [16] Y. Xiong and S.A. Shafer, *Moment and Hypergeometric Filters for High Precision Computation of Focus, Stereo and Optical Flow*. The Robotics Institute: Carnegie Mellon University, no. CMU-RI-TR-94-28, Pittsburgh, Penn., Sept. 1994.
- [17] P.J. Besl, *Range Imaging Sensors*. General Motors Research Laboratories no. GMR-6090, March 1988.
- [18] Y. Xiong and S.A. Shafer, "Variable Window Gabor Filters and Their Use in Focus and Correspondence," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pp. 668-671, June 1994.
- [19] S. K. Nayar, M. Watanabe, and M. Noguchi, "Real-Time Focus Range Sensor," *Proc. Int'l Conf. Computer Vision*, pp. 995-1,001, June 1995.
- [20] M. Watanabe, S.K. Nayar, and M. Noguchi, "Real-Time Implementation of Depth from Defocus," *Proc. SPIE Conf.*, Philadelphia, Penn., Oct. 1995.
- [21] M. Watanabe and S. K. Nayar, "Telecentric Optics for Computational Vision," *Proc. European Conf. Computer Vision*, Cambridge, U.K., Apr. 1996.
- [22] M. Watanabe and S.K. Nayar, "Minimal Operator Set for Passive Depth from Defocus," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, San Francisco, Calif., June 1996.
- [23] S. Inokuchi, K. Sato, and F. Matsuda, "Range Imaging System for 3D Object Recognition," *Proc. Seventh Int'l Conf. Pattern Recognition*, pp. 806-808, July 1984.
- [24] A. Gruss, S. Tada, and T. Kanade, "A VLSI Smart Sensor for Fast Range Imaging," *Proc. ARPA Image Understanding Workshop*, pp. 977-986, Apr. 1993, Washington, D.C.
- [25] T. Kanade, A. Gruss, and L.R. Carley, "A Very Fast VLSI Range-finder," *Proc. Int'l Conf. Robotics and Automation*, pp. 1,322-1,329, Apr. 1991, Sacramento, Calif.
- [26] T. Kanade, "Development of a Video-Rate Stereo Machine," *Proc. ARPA Image Understanding Workshop*, pp. 549-557, Nov. 1994.
- [27] E. Krotkov, "Focusing," *Int'l J. Computer Vision*, vol. 1, pp. 223-237, 1987.
- [28] T. Darrell and K. Wohn, "Pyramid Based Depth from Focus," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pp. 504-509, June 1988.
- [29] S.K. Nayar and Y. Nakagawa, "Shape from Focus," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 16, no. 8, pp. 824-831, Aug. 1994.
- [30] H.N. Nair and C.V. Stewart, "Robust Focus Ranging," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pp. 309-314, June 1988.
- [31] A. Krishnan and N. Ahuja, "Range Estimation From Focus Using a Non-Frontal Imaging Camera," *Proc. AAAI Conf.*, pp. 830-835, July 1993.
- [32] N. Asada, H. Fujiwara, and T. Matsuyama, "Edge and Depth from Focus," *Proc. Asian Conf. Computer Vision* pp. 83-86, Osaka, Japan, Nov. 1993.
- [33] P. Grossman, "Depth from Focus," *Pattern Recognition*, vol. 9, no. 1, pp. 63-69, 1987.
- [34] V.M. Bove, Jr., "Discrete Fourier Transform Based Depth-From-Focus," *Proc. OSA Topical Meeting Image Understanding and Machine Vision*, 1989.
- [35] J. Ens and P. Lawrence, "A Matrix Based Method for Determining Depth from Focus," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pp. 600-609, June 1991.
- [36] M. Gokstorp, "Computing Depth from Out-of-Focus Blur Using a Local Frequency Representation," *Proc. Intl. Conf. Pattern Recognition*, Oct. 1994.
- [37] R.G. Willson and S.A. Shafer, "Modeling and Calibration of Automated Zoom Lenses," PhD dissertation, The Robotics Institute, Carnegie Mellon University," CMU-RI-TR-94-03, Jan. 1994.



**Shree K. Nayar** is an associate professor at the Department of Computer Science, Columbia University. He received his PhD degree in electrical and computer engineering from the Robotics Institute at Carnegie Mellon University in 1990. His primary research interests are in computational vision and robotics, with emphasis on physical models for early visual processing, sensors and algorithms for shape recovery, pattern reaming and recognition, vision based manipulation and tracking, and the use of machine

vision for computer graphics and virtual reality.

Dr. Nayar has authored and coauthored papers that have received the David Marr Prize at the 1995 International Conference on Computer Vision (ICCV'90) held in Boston, Siemens Outstanding Paper Award at the 1994 IEEE Computer Vision and Pattern Recognition Conference (CVPR'94) held in Seattle, 1994 Annual Pattern Recognition Award from the Pattern Recognition Society, Best Industry Related Paper Award at the 1994 International Conference on Pattern Recognition (ICPR'94) held in Jerusalem, and the David Marr Prize at the 1990 International Conference on Computer Vision (ICCV'90) held in Osaka. He holds several U.S. and international patents for inventions related to computer vision and robotics.



**Masahiro Watanabe** received BS and MS degrees in precision engineering from the University of Tokyo in 1986 and 1988, respectively. He received the Annual Award of Japan Society of Precision Engineering for his paper, Precise Positioner Utilizing Rapid Deformations of Piezoelectric Elements, in March 1989. He joined the Production Engineering Research Lab., Hitachi Ltd., in 1988, where he conducted research on computer integrated manufacturing until 1991. Since then, he has been working on the development of optical measurement systems for LSI fabrication. He was a visiting scientist at Columbia University in 1994 and 1995 for the joint research program on range sensing between Hitachi Ltd. and Columbia University. He is a member of SPIE and JSPE.



**Minori Noguchi** received the BS degree in precision engineering from the University of Tokyo in 1982, after which he joined the Production Engineering Research Laboratory, Hitachi Ltd. He has conducted research on a wide range of topics related to metrology and inspection, including wafer cooling for LSI fabrication, trench depth measurement for SI wafers, particle inspection on reticula, particle inspection on mirror wafers, high resolution imaging with conjugated spatial filters, and particle inspection on

patterned wafers using frequency variable spatial filters. He was a visiting scientist in 1993 at Columbia University, where he developed optimal filters for shape from focus in collaboration with Prof. Shree K. Nayar. For this work, he received the Best Industry Related Paper Award at the 1994 International Conference on Pattern Recognition. He is a member of the IEEE, SPIE, and IAPR.