# CODAC: A COMPRESSIVE DEPTH ACQUISITION CAMERA FRAMEWORK

*Ahmed Kirmani, Andrea Colaço, Franco N. C. Wong, and Vivek K Goyal*

Research Laboratory of Electronics
Massachusetts Institute of Technology

## ABSTRACT

Light detection and ranging (LIDAR) systems use time of flight (TOF) in combination with raster scanning of the scene to form depth maps, and TOF cameras instead make TOF measurements in parallel by using an array of sensors. Here we present a framework for depth map acquisition using neither raster scanning by the illumination source nor an array of sensors. Our architecture uses a spatial light modulator (SLM) to spatially pattern a temporally-modulated light source. Then, measurements from a single omnidirectional sensor provide adequate information for depth map estimation at a resolution equal that of the SLM. Proof-of-concept experiments have verified the validity of our modeling and algorithms.

***Index Terms***— compressed sensing, depth maps, LIDAR, ranging, time of flight

## 1. INTRODUCTION

Sensing 3D scene structure is an integral part of applications ranging from 3D microscopy [1] to geographical surveying [2], and it is now increasingly of interest for consumer applications. While 2D imaging is a mature technology, 3D acquisition techniques have room for significant improvements in spatial resolution, range accuracy, and cost effectiveness. In comparison to stereo disparity, depth-from-focus, depth-from-shape, and depth-from-motion, active range acquisition systems such as LIDAR systems [3] and TOF cameras [4] are more robust against noise [5], work in real-time at video frame rates, and acquire range information from a single viewpoint with little dependence on scene reflectance or texture. Both LIDAR and TOF cameras operate by measuring the time difference of arrival between a transmitted pulse and the scene reflection. A LIDAR system consists of a pulsed illumination source such as a laser, a mechanical 2D laser scanning unit, and a single time-resolved photodetector or avalanche photodiode [3]. The TOF camera illumination unit is composed of an array of omnidirectional, modulated, infrared light emitting diodes (LEDs) [4]. The reflected light from the scene—with time delay proportional to distance—is focused at a 2D array of TOF range sensing pixels. A major shortcoming of LIDAR systems and TOF cameras is low spatial resolution, or the inability to resolve sharp spatial features in the scene.

This paper presents a new compressive depth acquisition camera (CoDAC) architecture as an alternative to both LIDAR and TOF camera systems. Like a LIDAR system, CoDAC uses only a single omnidirectional sensor. However, rather than raster scanning the scene to obtain spatial resolution, CoDAC uses a spatial light modulator (SLM) to serially illuminate subsets of the scene (see Fig. 1(a)). Spatial resolution equal to that of the SLM is achieved despite using fewer SLM patterns than the number of pixels in the SLM.

CoDAC is reminiscent of compressed sensing (CS) systems for photographic (intensity or reflectance) image capture [6], and it certainly has similarities. However, conventional compressed sensing techniques [7] are not applicable here because the information of interest (depths of scene points) appears nonlinearly in the measurements. For example, see the expression for $r(t)$ in Fig. 1(a), where the depths $d_A$, $d_B$, and $d_C$ are nonlinearly combined in all discrete samples of $r(t)$. A main contribution is thus a reconstruction technique that works despite the nonlinear mixing of scene depths in the measured values. Our method exploits approximation of the scene geometry as piecewise planar. Furthermore, since a scene depth is generally more compressible than scene reflectance or texture, we expect a smaller number of measurements to suffice; this is indeed the case, as our number of measurements is 1 to 5% of the number of pixels as compared to 10 to 40% for reflectance imaging [6].

In a previous application of the CS framework to LIDAR systems [8], the authors use spatial patterning of measurements provided by a digital micromirror device. Incident reflected light is measured with a photon-counting detector and gated to collect photons arriving from *a priori* chosen range intervals. The use of impulsive illumination and range gating make this a conventional CS problem in that the quantities of interest (reflectances as a function of spatial position, *within a depth range*) are combined linearly in the measurements. This approach achieves 3D imaging with a single sensor, but it has two major disadvantages: acquiring a complete scene depth map requires a full range sweep; and there is no method to distinguish between objects at different depths within a chosen range interval.

CoDAC has two acquisition stages and a two-step reconstruction procedure. The paper is organized according to these steps:
• Step 1, discussed in Section 3, uses no spatial patterning, i.e., a fully-transparent SLM configuration. Under the assumption that the scene is approximately piecewise planar, the continuous-time light intensity signal at the single photodetector is approximated well in a certain parametric class. Estimation of the parameters of the signal implies recovery of the range of depth values present in the scene. The use of a parametric signal modeling and recovery framework [9] enables us to achieve high depth resolution relative to the speed of the time sampling at the photodetector. After discretizing the depths identified in this step, the remaining problem is to find correspondences between spatial locations and depths to form the depth map.
• Step 2, discussed in Section 4, uses many pseudorandom binary patterns on the SLM. The assumption that the scene is approximately piecewise planar translates to the Laplacian of the depth map being approximately sparse. We introduce a novel convex optimization problem that finds the depth map consistent with the measure-
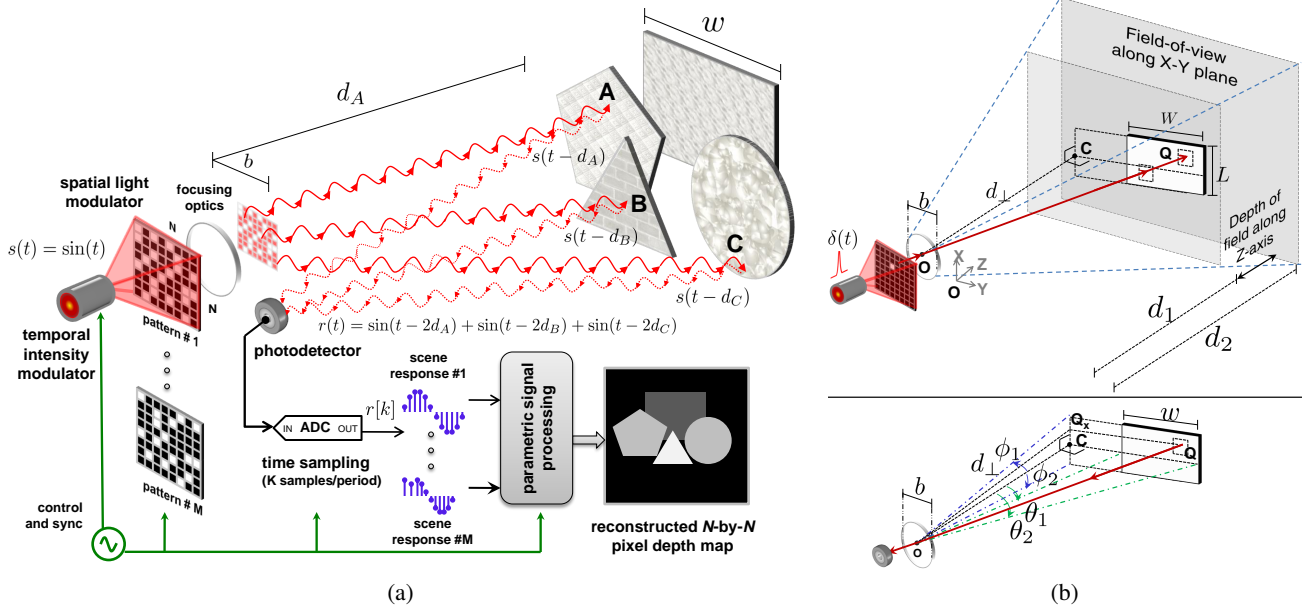
**Fig. 1**. (a) Proposed architecture for acquiring depth maps of scenes constituted of piecewise-planar facets. The scene is in far field, i.e., the baseline $b$ and the dimensions of each planar facet $w$ are much smaller than the distance between the imaging device and the scene. A light source with periodically-varying intensity $s(t)$ illuminates an $N \times N$-pixel SLM. The scene is serially illuminated with $M$ spatial patterns. For each patterned illumination the reflected light is focused at the photodetector and $K$ digital time samples are recorded. The $M \cdot K$ samples are computationally processed using parametric signal processing to reconstruct an $N \times N$-pixel depth map of the scene. (b) [TBD]

ments that approximately minimizes the number of nonzero entries in the Laplacian of the depth map. Solving this optimization problem yields the desired depth map.

To convey the main ideas despite space limitations, Sections 3 and 4 use restrictions to simple rectangular planar facets and omit many details. More general scenes are discussed in [10]. The full paper [10] also includes details on proof-on-concept experiments and extensions to scenes with multiple planar facets that are not necessarily rectangular, textured scenes, and non-impulsive illumination.

## 2. ANALYSIS FOR A SINGLE RECTANGULAR FACET

Consider the setup shown in Fig. 1(b). A chosen SLM pattern is focused on the scene using a focusing system. The center of the focusing system is denoted by $O$ and is also the origin for a 3D coordinate system $(X, Y, Z)$. All angles and distances are measured with respect to this global coordinate system. The focusing optics for the SLM illumination unit are chosen such that it has a depth-of-field (DOF) between distances $d_1$ and $d_2$ ($d_1 < d_2$) along the $Z$ dimension and a square field-of-view (FOV) along the $X$-$Y$ axes. Thus, the dimensions of a square SLM pixel projected onto the scene remains constant within the DOF and across the FOV. We denote the dimensions of an SLM pixel within the DOF by $\Delta \times \Delta$. An SLM with higher spatial resolution corresponds to a smaller value of $\Delta$. We also assume that the scene lies within the DOF so that all planar facets in the scene are illuminated by projection pixels of the same size. We consider only binary patterns, i.e., each SLM pixel is chosen to be either completely opaque or fully transparent.

The light reflected from the scene is focused at the photodetector. We assume that the baseline separation $b$ between the focusing optics of the detector and the SLM illumination optics is very small

compared to the distance between the imaging device and the scene. Thus, we may conveniently model $O$ as the effective optical center of the entire imaging setup (illumination and detector).

*Notation.* Let $OC$ be the line that lies in the $Y$-$Z$ plane and is also perpendicular to the rectangular facet. The following parameters completely specify the rectangular facet (see Fig. 1(b)):

- $d_\perp$ denotes the length of the line $OC$.
- $\phi_1$ and $\phi_2$ are angles between line $OC$ and the extreme rays connecting the vertical corners of the rectangular facet to $O$. Also let $|\phi_1 - \phi_2| = \delta\phi$. Clearly, $\delta\phi$ is related to $L$.
- $\theta_1$ and $\theta_2$ are angles between line $OC$ and the extreme rays connecting the horizontal corners of the rectangular facet to $O$. Let $|\theta_1 - \theta_2| = \delta\theta$. Clearly, $\delta\theta$ is related to $W$.
- $\alpha$ is the angle between $OC$ and the $Z$ axis in the $Y$-$Z$ plane.

*Assumptions.* We assume that the scene is in the far field, i.e., the dimensions of the rectangular facet are small compared to the distance between the scene and the imaging device, or $W \ll d_1$ and $L \ll d_1$. This implies that $\delta\phi$ and $\delta\theta$ are small angles and the radial fall-off attenuation of light arriving from different points on the rectangular facet is approximately the same for all the points. We also assume that the rectangular facet is devoid of texture and reflectance patterns. Finally, we normalize to unit speed of light.

## 3. RESPONSE TO FULLY-TRANSPARENT SLM PATTERN

Let $Q$ be a point on the rectangular planar facet at an angle of $\theta$ and $\phi$ with respect to the line $OC$ as shown in Fig. 1(b). A unit-intensity illumination pulse, $s(t) = \delta(t)$, that originates at the source at time $t = 0$ will be reflected from $Q$, attenuated due to scattering, and arrive back at the detector delayed in time by an amount proportional
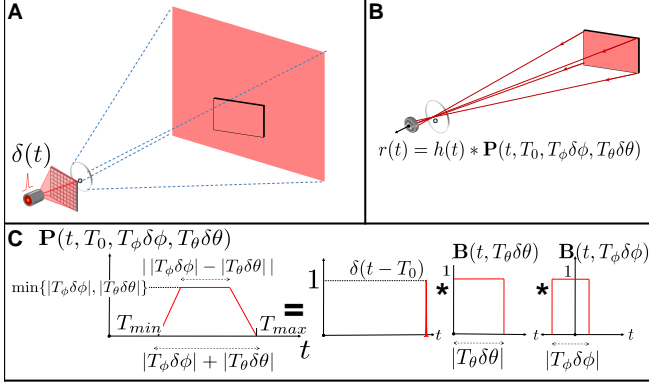
**Fig. 2**. (**A**) All-ones scene illumination. (**B**) Scene response to all-ones scene illumination. (**C**) Diagrammatic explanation of the modeling of the parametric signal $p(t)$.

**Fig. 3**. (**A**) Binary patterned scene illumination. (**B**) Scene response to all-ones scene illumination. (**C**) Diagrammatic explanation of the high-resolution SLM (small $\Delta$) approximation. (**D**) Modeling of the parametric signal $\mathbf{U}^p(t)$ as a weighted sum of equally-spaced Diracs. Note that $\mathbf{U}^p(t)$ has the same time envelope as the signal $\mathbf{P}(t, T_0, T_\phi \delta\phi, T_\theta \delta\theta)$.

to the distance $2\,|OQ|$. Thus the signal incident on the photodetector in response to impulse illumination of $Q$ is $q(t) = a\,\delta(t - 2\,|OQ|)$, where $a$ is the total attenuation (transmissivity) of the unit-intensity pulse. Denoting the photodetector impulse response by $h(t)$, the electrical output $r_q(t)$ of the photodetector is

$$r_q(t) \;=\; h(t) * a\,\delta(t - 2\,|OQ|) \;=\; a\,h(t - 2\,|OQ|).$$

The response to a fully-transparent SLM pattern (see Fig. 2) is obtained by integrating $r_q(t)$ over $\theta \in [\theta_1, \theta_2]$ and $\phi \in [\phi_1, \phi_2]$. For notational simplicity, let $\gamma(\phi_1, \theta_1) = \sqrt{\sec^2 \phi_1 + \tan^2 \theta_1}$. Then under an appropriate small-angle approximation, we obtain

$$r(t) \;=\; \frac{a}{T_\phi T_\theta}\, h(t) * \mathbf{P}(t, T_0, T_\phi \delta\phi, T_\theta \delta\theta)$$

where $\mathbf{P}(t, T_0, T_\phi \delta\phi, T_\theta \delta\theta)$ is defined graphically in Fig. 2-C with

$$T_\phi = \frac{2d_\perp}{\gamma(\phi_1, \theta_1)} \tan \phi_1 \sec^2 \phi_1, \quad T_\theta = \frac{2d_\perp}{\gamma(\phi_1, \theta_1)} \tan \theta_1 \sec^2 \theta_1,$$

and $T_0 = 2d_\perp \gamma(\phi_1, \theta_1)$.

*Parameter Recovery.* Estimating the function $\mathbf{P}(t, T_0, T_\phi \delta\phi, T_\theta \delta\theta)$ by processing the digital samples $r[k]$ of function $r(t)$ enables estimation of various parameters including $T_{\min}$ and $T_{\max}$. The detector impulse response $h(t)$ is generally modeled as a bandlimited low-pass filter. Thus, the general deconvolution problem of obtaining $\mathbf{P}(t, T_0, T_\phi \delta\phi, T_\theta \delta\theta)$ from samples $r[k]$ is ill-posed and highly sensitive to noise. However, our modeling shows that the light transport function $\mathbf{P}(t, T_0, T_\phi \delta\phi, T_\theta \delta\theta)$ is piecewise linear. This knowledge makes the recovery of $\mathbf{P}(t, T_0, T_\phi \delta\phi, T_\theta \delta\theta)$ a *parametric deconvolution* problem that we solve using the parametric signal processing framework described in [11].

The analysis up to this point is independent of the tilt $\alpha$ and orientation of the rectangular plane with respect to the global coordinate system $(X, Y, Z)$. This is intuitive because the illumination and sensing are omnidirectional.

## 4. RESPONSE TO BINARY SLM PATTERN

As discussed in Section 3, the SLM pixels discretize the FOV into small squares of size $\Delta \times \Delta$. We index both the SLM pixels and
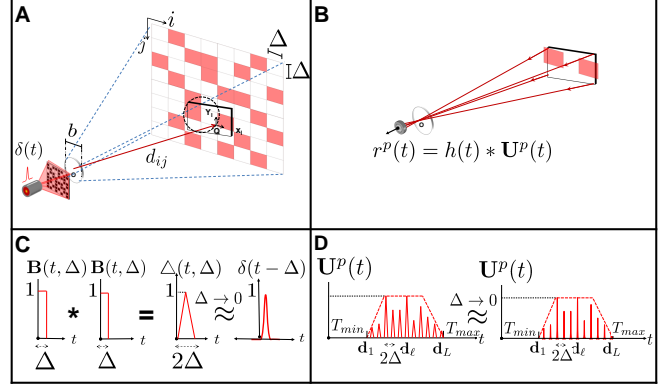
the corresponding scene points by $(i, j)$, and we illuminate the $M$ binary SLM patterns by $p$. Consider the scene shown in Fig. 3 consisting of only a single rectangular planar facet and illuminated with a binary SLM pattern given by a collection of values $\{c_{ij}^p : i = 1, \ldots, N,\ j = 1, \ldots, N\}$ where each $c_{ij}^p$ is either 0 or 1. Let $d_{ij}$ denote the depth in the direction of illumination of SLM pixel $(i, j)$. The depths for all $(i, j)$ associated with the rectangular facet form the *depth map*, $\mathbf{D} = \{\mathbf{D}_{ij} : i = 1, \ldots, N,\ j = 1, \ldots, N\}$, where $\mathbf{D}_{ij} = d_{ij}$ if rays along SLM illumination voxel $(i, j)$ intersect the rectangular facet and $\mathbf{D}_{ij} = 0$ otherwise. Also define the binary valued *index map*, $\mathbf{I} = \{\mathbf{I}_{ij} : i = 1, \ldots, N,\ j = 1, \ldots, N\}$, where $\mathbf{I}_{ij} = 1$ if $\mathbf{D}_{ij} \neq 0$.

Analysis of the scene response is detailed in [10]. Digital samples of the received signal $r^p[k]$ allow us to recover the depth map $\mathbf{D}$ as follows. The set of distance values, $\{d_{ij}\}$, contains repetitions; i.e., several $(i, j)$ positions may have the same depth value $d_{ij}$. All these points will lie on a circular arc on the rectangular facet as shown in Fig. 3-A. Each $d_{ij}$ belongs to the set of equally-spaced distinct depth values $\{\mathbf{d}_1, \mathbf{d}_2, \ldots, \mathbf{d}_L\}$ where

$$L = \frac{T_{\max} - T_{\min}}{2\Delta}, \quad \mathbf{d}_1 = T_{\min}, \quad \mathbf{d}_\ell = \mathbf{d}_1 + 2\Delta\ell \quad \ell = 1, \ldots, L.$$

In the case of all-ones SLM illumination discussed in Section 3 we obtain the continuous signal $\mathbf{P}(t, T_0, T_\phi \delta\phi, T_\theta \delta\theta)$, but in the patterned illumination case, we obtain a signal $\mathbf{U}^p(t)$ satisfying

$$\lim_{\Delta \to 0} \mathbf{U}^p(t) \;=\; \sum_{\ell=1}^{L} \left( \sum_{i=1}^{N} \sum_{j=1}^{N} c_{ij}^p I_{ij}^\ell \right) \delta(t - 2\mathbf{d}_\ell), \qquad (1)$$

where $I_{ij}^\ell = 1$ if $d_{ij} = \mathbf{d}_\ell$ and rays along SLM illumination voxel $(i, j)$ intersect with the rectangular facet and $I_{ij}^\ell = 0$ otherwise. Clearly $\mathbf{I}_{ij} = \sum_{\ell=1}^{L} I_{ij}^\ell$ and $\mathbf{D}_{ij} = \sum_{\ell=1}^{L} \mathbf{d}_\ell I_{ij}^\ell$. With this new notation, the depth map $\mathbf{D}$ associated with the rectangular facet is the weighted sum of the index maps $\{I^\ell : \ell = 1, \ldots, L\}$.

Let $K$ denote the total number of samples collected by the ADC so we have $\{r^p[k] : k = 1, \ldots, K\}$ for each $p$, and let $R^p[k]$ denote

the discrete Fourier transform of $r^p[k]$. Similarly define $H^p[k]$ for the impulse response samples $h^p[k]$. Then

$$\frac{R^p[k]}{H[k]} = \frac{af}{4} \sum_{\ell=1}^{L} \left( \sum_{i=1}^{N} \sum_{j=1}^{N} c_{ij}^p I_{ij}^\ell \right) e^{-\mathbf{j}(4\pi f \mathbf{d}_\ell)k}, \quad k = 1, \ldots, K. \tag{2}$$

For notational simplicity let

$$y_\ell^p = \sum_{i=1}^{M} \sum_{j=1}^{N} c_{ij}^p I_{ij}^\ell, \qquad \ell = 1, \ldots, L. \tag{3}$$

The constants $a$ and $f$ are computed using calibration and are computationally compensated using normalization. Since the values $\{\mathbf{d}_1, \mathbf{d}_2, ..., \mathbf{d}_L\}$ are known, Eq. (2) can be represented as a system of linear equations that can be compactly written as

$$\mathbf{R^P}/\mathbf{H} = \mathbf{V}\mathbf{y}^p \tag{4}$$

where the division is elementwise and $\mathbf{V}$ is a Vandermonde matrix. We need $K \geq L$ so that we can uniquely solve the linear system in (4). Furthermore, a larger value of $K$ allows us to mitigate the effect of noise by producing least square estimates of $\mathbf{y^P}$.

As we illuminate the scene with a total of $M \ll N^2$ different binary SLM illumination patterns and process the samples of the received light signal with (4) for each choice of SLM illumination pattern, we obtain the transformed data $\{y_\ell^p : \ell = 1, \ldots, L, p = 1, \ldots, M\}$ which is related to the index images $\{I^\ell : \ell = 1, \ldots, L\}$ via a system of linear equations:

$$\mathbf{y}_{L \times M} = \left[ I^1 \ldots I^\ell \ldots I^L \right]_{L \times N^2} \mathbf{C}_{N^2 \times M}. \tag{5}$$

This system of equations is under-constrained since there are $L \times N^2$ unknowns (corresponding to the unknown values of $\left[ I^1 \ldots I^\ell \ldots I^L \right]$) and only $L \times M$ available transformed data observations $\mathbf{y}$.

*Algorithms for Depth Map Reconstruction.* To enable depth map reconstruction even though we have much fewer observations than unknowns, we exploit the structure of scene depth. The depth values $d_{ij}$ correspond to the distances from $O$ to points that are constrained to lie on a rectangular facet, and the distances $d_{ij}$ are also linearly spaced between $\mathbf{d}_1$ and $\mathbf{d}_L$. The planar constraint and linear variation imply that the depth map $\mathbf{D}$ is *sparse* in the second-finite difference domain. Thus, it is possible to recover $\mathbf{D}$ from $\mathbf{y}$ by solving the following constrained $\ell_1$-regularized optimization problem:

$$\underset{\mathbf{D}}{\text{minimize}} \quad \left\| \mathbf{y} - \left[ I^1 \ldots I^\ell \ldots I^L \right] \mathbf{C} \right\|_{\mathrm{F}} + \| \mathbf{\Phi} \mathbf{D} \|_1$$

subject to $\quad \sum_{\ell=1}^{L} I_{ij}^\ell = 1, \forall i,j \qquad \sum_{\ell=1}^{L} \mathbf{d}_\ell I^\ell = \mathbf{D}, \quad$ and

$$I_{ij}^\ell \in \{0,1\}, \quad \ell = 1, \ldots, L, \quad i = 1, \ldots, N, \quad j = 1, \ldots, N,$$

where $\| . \|_{\mathrm{F}}$ is the Frobenius norm and $\mathbf{\Phi}$ is the second-order finite difference operator matrix.

While this optimization problem already contains a convex relaxation in its use of $\| \mathbf{\Phi} \mathbf{D} \|_1$, it is nevertheless computationally intractable because of the integrality constraints $I_{ij}^\ell \in \{0, 1\}$. Using a further relaxation of $I_{ij}^\ell \in [0, 1]$ yields the following tractable formulation:

$$\underset{\mathbf{D}}{\text{minimize}} \quad \left\| \mathbf{y} - \left[ I^1 \ldots I^\ell \ldots I^L \right] \mathbf{C} \right\|_{\mathrm{F}} + \| \mathbf{\Phi} \mathbf{D} \|_1$$

subject to $\quad \sum_{\ell=1}^{L} I^\ell = 1, \forall i,j \qquad \sum_{\ell=1}^{L} \mathbf{d}_\ell I^\ell = \mathbf{D}, \quad$ and

$$I_{ij}^\ell \in (0,1) \quad \ell = 1, \ldots, L, \quad i = 1, \ldots, N, \quad j = 1, \ldots, N.$$

We solved this convex optimization problem using CVX, a package for specifying and solving convex programs [12].

## 5. CONCLUSION

We have presented a method for acquiring 2D depth maps of piecewise-planar scenes using samples measured by a single photodetector in response to spatiotemporally-modulated scene illuminations. Our acquisition architecture consists of a non-scanning 2D SLM and a single photodetector. We have demonstrated (see [10]) that it is possible to acquire scene depth at both high range resolution and spatial resolution with significantly-reduced device complexity and hardware cost as compared to state-of-the-art LIDAR systems and TOF cameras. These gains are achieved by developing a depth acquisition framework based on parametric signal processing that exploits the sparsity of natural scene structures.

## 6. REFERENCES

[1] K. Carlsson, P. E. Danielsson, R. Lenz, A. Liljeborg, L. Majlöf, and N. Åslund, "Three-dimensional microscopy using a confocal laser scanning microscope," *Opt. Lett.*, vol. 10, no. 2, pp. 53–55, 1985.

[2] A. Wehr and U. Lohr, "Airborne laser scanning—an introduction and overview," *ISPRC J. Photogrammetry & Remote Sensing*, vol. 54, no. 2–3, pp. 68–82, Jul. 1999.

[3] B. Schwarz, "LIDAR: Mapping the world in 3D," *Nature Photonics*, vol. 4, no. 7, pp. 429–430, Jul. 2010.

[4] S. Foix, G. Alenyà, and C. Torras, "Lock-in time-of-flight (ToF) cameras: A survey," *IEEE Sensors J.*, vol. 11, no. 9, pp. 1917–1926, Sep. 2011.

[5] S. Hussmann, T. Ringbeck, and B. Hagebeuker, "A performance review of 3D TOF vision systems in comparison to stereo vision systems," in *Stereo Vision*, A. Bhatti, Ed. InTech, 2008, pp. 103–120.

[6] M. F. Duarte, M. A. Davenport, D. Takhar, J. N. Laska, T. Sun, K. Kelly, and R. G. Baraniuk, "Single-pixel imaging via compressive sampling," *IEEE Signal Process. Mag.*, vol. 25, no. 2, pp. 83–91, Mar. 2008.

[7] E. J. Candès and M. B. Wakin, "An introduction to compressive sampling," *IEEE Signal Process. Mag.*, vol. 25, no. 2, pp. 21–30, Mar. 2008.

[8] G. Howland, P. Zerom, R. W. Boyd, and J. C. Howell, "Compressive sensing LIDAR for 3D imaging," in *Conf. Lasers and Electro-Optics.* OSA Technical Digest Series, 2011, p. CMG3.

[9] T. Blu, P.-L. Dragotti, M. Vetterli, P. Marziliano, and L. Coulot, "Sparse sampling of signal innovations," *IEEE Signal Process. Mag.*, vol. 25, no. 2, pp. 31–40, Mar. 2008.

[10] A. Kirmani, A. Colaço, F. N. C. Wong, and V. K. Goyal, "Exploiting sparsity in time-of-flight range acquisition using a single time-resolved sensor," *Opt. Expr.*, 2011, to appear.

[11] P. L. Dragotti, M. Vetterli, and T. Blu, "Sampling moments and reconstructing signals of finite rate of innovation: Shannon meets Strang-Fix," *IEEE Trans. Signal Process.*, vol. 55, no. 5, pp. 1741–1757, May 2007.

[12] M. Grant and S. Boyd, "CVX: Matlab software for disciplined convex programming, version 1.21," http://cvxr.com/cvx, Apr. 2011.