

Low-Cost Compressive Sensing for Color Video and Depth

Xin Yuan, Patrick Lull, Xuejun Liao, Jianbo Yang, David J. Brady, Guillermo Sapiro and Lawrence Carin
Department of Electrical and Computer Engineering, Duke University, Durham, NC, 27708, USA
{xin.yuan, prl12, xjliao, jy118, dbrady, guillermo.sapiro, lcarin}@duke.edu

Abstract

A simple and inexpensive (low-power and low-bandwidth) modification is made to a conventional off-the-shelf color video camera, from which we recover multiple color frames for each of the original measured frames, and each of the recovered frames can be focused at a different depth. The recovery of multiple frames for each measured frame is made possible via high-speed coding, manifested via translation of a single coded aperture; the inexpensive translation is constituted by mounting the binary code on a piezoelectric device. To simultaneously recover depth information, a liquid lens is modulated at high speed, via a variable voltage. Consequently, during the aforementioned coding process, the liquid lens allows the camera to sweep the focus through multiple depths. In addition to designing and implementing the camera, fast recovery is achieved by an anytime algorithm exploiting the group-sparsity of wavelet/DCT coefficients.

1. Introduction

A variety of techniques have been developed to extract information from a single image. For example, the depth-from-focus method [11, 24, 25] allows estimation of a 3D scene (depth-dependent focus) from a single 2D image. The mosaic and demosaic technique [2] allows recovery of color information from a gray-scale image. Recently, inspired by compressive sensing (CS) [6], video has been extracted from a single image [8, 9, 14, 16, 23]. In this setting the measured data are acquired at a low frame rate, with coding at a faster rate, and high-frame-rate video is computationally recovered subsequently.

In this paper we develop a new method that borrows and extends ideas from this previous work. Specifically, like [8, 9, 14, 16] we perform high-frequency coding of video collected at a low frame rate, with CS-based inversion. Our coding strategy differs from previous work in that we use a *single* code that is inexpensively translated via a piezoelectric device. We recover color via a hardware mosaicing and computational demosaicing procedure like in conventional cameras. The newest aspect of the proposed approach is that we use a lens with voltage-dependent index of refraction

(a liquid lens), and by varying the voltage at high rate, the recovered high-rate video corresponds to capturing data at varying focus points (depths). For each of the measured frames, we recover multiple color frames, and these multiple frames capture variable focus depths.

We here show example results that summarize the three key aspects of the approach: mosaicing for color, high-speed coding for video, and fast time-dependent focus for depth, with the data measured at a low frame rate. We first consider mosaicing and high-speed coding, with the focus held constant. In Figure 1 two compressive measurements (real data from our camera) are shown at left. These are two consecutive frames collected at frame rate 30 frames/sec, using an off-the-shelf Guppy Pro camera [1], with a high-speed coding element, as summarized in Figure 2. At right in Figure 1, are shown $n_t = 22$ recovered frames from each of the compressive measurements: 22 color video frames recovered from a single *monochromatic* coded image. Each measurement in Figure 1 employs a high-speed code (here a shifting mask) to modulate the light during the integration time-period Δ_t ; see Figure 2(a).

In Figure 3 we now consider results in which the focus (observed depth) is varied at a rate that is fast relative to the rate of the video camera collecting the data (now the measurements employ mosaicing, high-speed coding, and variation of the focal depth). The variable focus is manifested by varying the voltage on a liquid lens (see Figure 4). The coded data, with subsequent CS inversion, allows recovery of multiple frames for each measured frame. Since the focus has been adjusted at a fast rate, these high-frequency recovered frames also capture multiple depths.

In Figure 3, the top figure depicts the camera and scene, composed of multiple objects at varying depths. At the bottom in this figure, we depict one of the measured gray-scale frames (here measured at a frame rate of 30 frames/sec), and at bottom right is depicted the recovered color video from this single frame. Note that because of the high-speed time-varying focus, we effectively recover multiple depths, defined by the focus for which a given region of the image is sharpest.

In this paper, we describe in detail how the camera that

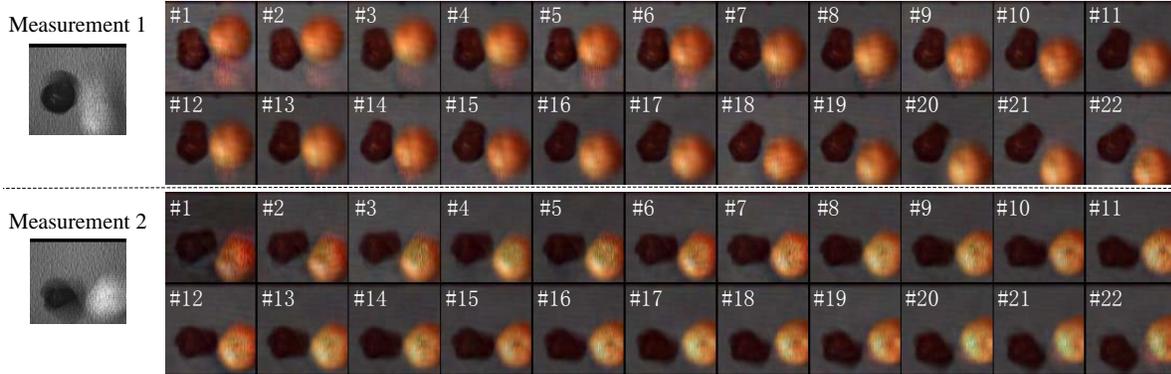


Figure 1. Reconstruction of real measurement with frame size 512×512 . Two plastic fruits dropping, touch a surface at the bottom, and then rebound. Left-most are two contiguous measurements, and the right part shows 44 corresponding reconstructed video frames ($n_t = 22$ recovered video frames from each compressive measurement). We can see clear motion and color information from the reconstructed frames (video online at [18]). Note in this dataset we did not change the focal plane during the capture of the measurement.

took these measurements was implemented, with a summary provided in Figure 2 for the coding and mosaicing components, and in Figure 4 for the camera setup. We also discuss a new CS inversion algorithm, that is endowed with a guarantee that upon every iteration the residual between the true video and the estimate is assured to be reduced (with technical requirements, that are discussed). This is an “any-time” algorithm, in that the estimated underlying video may be approximated at any time based upon computations thus far, and the quality of the inversion is guaranteed to improve with additional computations.

The contributions of this paper are: *i*) development of a new *low-cost and low-power* CS camera, that for each low-frame-rate measured image allows recovery of multiple color frames focused at different depths (at high-frame-rate); and *ii*) application of an anytime CS inversion algorithm to data measured by the camera, providing fast recovery of *high-speed motion, color and depth* information from the scene.

2. Hardware Design

The proposed imaging system is built with an off-the-shelf camera, specifically a Guppy Pro camera [1], by adding a *liquid lens* to change the focal plane, and by moving a *single mask* to modulate the high speed frames during one integration time-period, both extremely low-power additions (in contrast with common alternatives in the literature, as detailed in Section 2.3). The main challenge is to synchronize temporal modulation (coding) with time variation of the focus location (for capture of variable depth). Figure 4 depicts the setup of our camera, and Figure 5(b) shows the synchronized control of the mask and liquid lens.

2.1. Coding strategy

The focal plane of the liquid lens used in the camera is controlled by the voltage of the input signal, which also controls the piezoelectric translator to shift the mask. The control signal (a triangular wave) is generat-

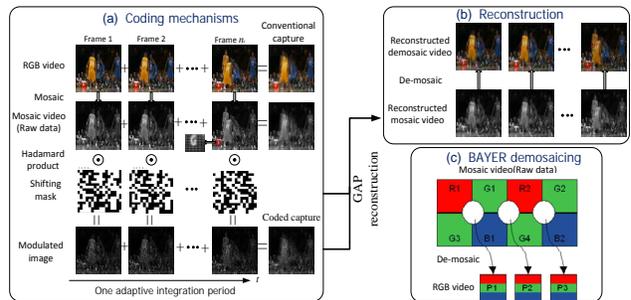


Figure 2. Illustration of the proposed method. (a) First row shows n_t color (RGB) frames of the original high-speed video; second row shows each color frame rearranged into a Bayer-filter mosaic; third row depicts the (horizontally) moving mask used to modulate the high-speed frames (black is zero, white is one); fast translation manifested by a piezoelectric translator. Fourth row shows the modulated frames, whose sum gives a single coded capture. (b) Recovered RGB frames arranged into a Bayer-filter mosaic (second row), which is de-mosaiced to give the color frames (first row). (c) The demosaicing process [1].

ed by a function generator and then we use power divider to distribute the signal (Figure 5) to the liquid lens and the piezoelectric translator to achieve the synchronization. The shifting (through the piezoelectric translator) of the same mask (Figure 2(a)) is utilized to modulate the high-speed frames. This modulation enjoys advantages of low-power ($\sim 0.2W$), low-cost and low-bandwidth implementation, compared for example with the modulation by liquid-crystal-on-silicon (LCoS) in [8, 16] (power $> 3W$, and high-bandwidth electronic switching/coding). In the experiments, we show that the proposed efficient coding mechanism yields similar results to that of the relatively expensive LCoS-type coding. During the calibration, we approximate the continuous moving of the mask by discrete patterns (Figure 5(a)).

We record temporally (and depth) compressed measurements for RGB colors on a Bayer-filter mosaic, where the

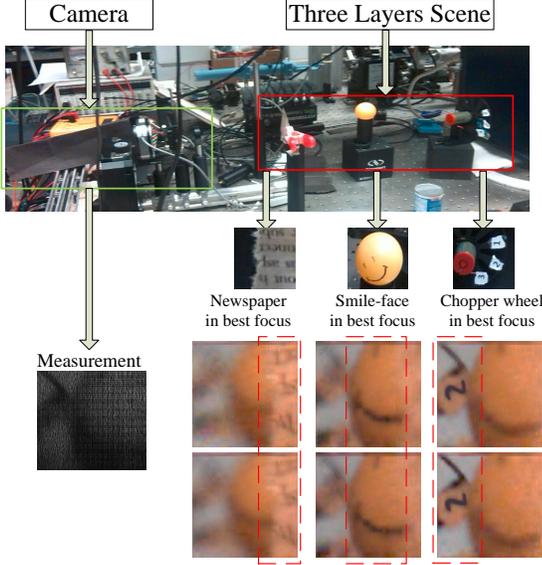


Figure 3. Experimental scene (top), three objects are placed at three depths. 6 selected recovered frames (bottom right) reconstructed all from the real single measurement (bottom left) are shown as examples. The focal plane varies from the newspaper (near) to the chopper-wheel (far) during the integration time period. Note how chopper-wheel goes from blur (two left columns) to sharp (refer to the number “2”). Note the motion of the *moving* chopper-wheel (video online at [18]).



Figure 4. Setup of the camera.

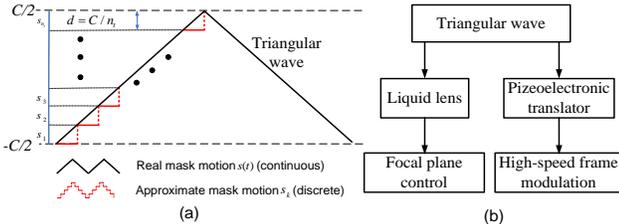


Figure 5. Control signal: (a) triangular wave and (b) distributed control.

three colors are arranged in the pattern shown in the right bottom of Figure 2. The single coded image is partitioned into four components, one for R and B and two for G (each

is 1/4 the size of the original spatial image). The CS recovery (video from a single measurement) is performed separately on these four mosaiced components, and then the demosaicing process in Figure 2(b) is employed to manifest the final color video. One may also jointly perform CS inversion on all 4 components, with the hope of sharing information on the importance of (here wavelet and DCT) components; this was also done, and the results were very similar to processing R, B, G1 and G2 separately.

2.2. Measurement model

Let $f(x, y, z, t)$ denote the continuous/analog spatio-temporal volume of the video being measured, with (x, y, z) symbolizing the 3D space coordinate and t denoting time. Note the depth of the scene is defined by z . Additionally, let object-space and image-space coordinates be respectively designated with unprimed and primed variables. Given an $n_x \times n_y$ -pixel sensor with pixel size Δ and integration time Δ_t , space-time compressive measurements $g(x', y', t') \in \mathbb{R}^2$ are formed on the detector, with $t' < t$. The digital data used to represent the scene consists of discrete samples of the continuous transformation

$$g(x', y', t') = \int \int \int f(x, y, z, t) T(x-r(t), y-s(t)) \times \text{rect}\left(\frac{x-x'}{\Delta}, \frac{y-y'}{\Delta}\right) \text{rect}\left(\frac{t-t'}{\Delta_t}\right) dy dx dt dz, (1)$$

where $T(x-r(t), y-s(t))$ represents a random binary transmission pattern that translates with periodic transverse (x, y) motion parameterized by $(r(t), s(t))$. The spatial and temporal pixel sampling functions, $\text{rect}\left(\frac{x-x'}{\Delta}, \frac{y-y'}{\Delta}\right)$ and $\text{rect}\left(\frac{t-t'}{\Delta_t}\right)$, bandlimit the incident optical datastream, which is equivalently represented as $f = c * b$, with $*$ denoting the convolution operator. Here, $c(x, y, z, t)$ denotes an instantaneous all-in-focus representation of the spatiotemporal scene and $b(x, y, z(t))$ is a time and depth varying (blur) kernel imparted by the liquid lens on the focal volume.

The discrete formulation of the model can be simplified. The frame at each depth (n_t depths in Figure 5(a)) is approximated as being modulated by a single unique code (approximated by the shifting mask, Figure 5(a)); in reality the code/mask is always moving continuously. At each depth, the (physical/continuous) frame can be denoted by $\tilde{f}(x, y, t)$ (z is now manifested by t), and after digitization, we represent it as $\mathbf{X}_k \in \mathbb{R}^{n_x \times n_y}$, $\forall k = 1, \dots, n_t$. Denoting the coding pattern of the mask by $\mathbf{H}_k \in \mathbb{R}^{n_x \times n_y}$, the measurement $\mathbf{Y} \in \mathbb{R}^{n_x \times n_y}$ is $y_{i,j} = \sum_{k=1}^{n_t} x_{i,j,k} h_{i,j,k} + e_{i,j}$, $\forall i = 1, \dots, n_x$; $j = 1, \dots, n_y$; or

$$\mathbf{Y} = \sum_k \mathbf{X}_k \odot \mathbf{H}_k + \mathbf{E}, (2)$$

where \mathbf{E} denotes the noise and \odot symbolizes the Hadamard (element-wise) product. By vectorizing each frame and then concatenating them, we have $\mathbf{x} =$

$$\begin{aligned}
& [(\text{vec}(\mathbf{X}_1))^T, \dots, (\text{vec}(\mathbf{X}_{n_t}))^T]^T \in \mathbb{R}^{n_x n_y n_t}, \text{ and (2) can} \\
& \text{be written as:} \\
& \mathbf{y} = [\text{diag}(\text{vec}(\mathbf{H}_1)), \dots, \text{diag}(\text{vec}(\mathbf{H}_{n_t}))] \mathbf{x} + \mathbf{e} \\
& = \mathbf{\Psi} \mathbf{x} + \mathbf{e}, \tag{3}
\end{aligned}$$

where $\mathbf{\Psi} = [\text{diag}(\text{vec}(\mathbf{H}_1)), \dots, \text{diag}(\text{vec}(\mathbf{H}_{n_t}))]$. The imaging process may therefore be modeled as in the standard CS problem. The goal is to estimate \mathbf{x} , given \mathbf{y} and $\mathbf{\Psi}$. Before presenting how we address this, we further comment on related work, now in that we have introduced the proposed hardware.

2.3. Related work

Video compressive sensing has been investigated in [8, 9, 14, 16, 19], by capturing low frame-rate video to reconstruct high frame-rate video. The LCoS used in [8, 16] can modulate as fast as 3000 fps by pre-storing the exposure codes, but, because the coding pattern is continuously changed at each pixel throughout the exposure, it requires considerable energy consumption ($> 3W$) and bandwidth compared with the proposed modulation, in which a single mask is translated using a pizeoelectronic translator (requiring $\sim 0.2W$). Similar coding was used in [14]. However, we investigate color video here, and thus demosaicing is needed; because of the R, G and B channels, we need to properly align (in hardware of course) the mask more accurately compared with the monochromatic video in [14]. Therefore, [14] can be seen as a special case of the proposed camera. Furthermore, we also extract the depth information from the defocus phenomenon of the reconstructed frames, which has not been considered in the above papers.

Coded apertures have been used often in computational imaging for depth estimation [11, 12, 24]. However, these only consider still images. From the algorithms investigated therein, one can get the depth map from a still image. In [10] an imaging system was presented that enables one to control the depth of field by varying the position and/or orientation of the image detector, during the integration time of a single photograph. However, moving the detector costs more energy than controlling the liquid lens in the proposed design (almost no power consumption), and the camera developed in [10] can only provide a single all-in-focus image without the depth information. Furthermore, no motion information is considered in the above coded-aperture cameras, while here we consider video (allowing depth estimation on moving scenes).

3. Reconstruction algorithm

We reconstruct high-frame-rate video from low-frame-rate measurements via an *anytime* algorithm, the generalized alternating projection (or GAP) algorithm, first developed in [13] for other applications. GAP produces a sequence of partial solutions that *monotonically* converge to the *true signal* (thus, anytime). In [13], the authors did not mention how to select group weights and no real data or application was considered. The manner in which the GAP

algorithm is employed here, as well as the application considered, is significantly different from [13]. In the following, we first review the underlying GAP algorithm and then show how to improve it to get better results for the data considered here.

GAP is used to investigate the *group-sparsity* of wavelet/DCT coefficients of the video to be reconstructed. Let $\mathbf{T}_x \in \mathbb{R}^{n_x \times n_x}$, $\mathbf{T}_y \in \mathbb{R}^{n_y \times n_y}$, $\mathbf{T}_t \in \mathbb{R}^{n_t \times n_t}$ be orthonormal matrices defining bases such as wavelets or DCT [15]. Define $\mathbf{w} = (\mathbf{T}_t^T \otimes \mathbf{T}_y^T \otimes \mathbf{T}_x^T) \mathbf{x}$, and $\mathbf{\Phi} = \mathbf{\Psi} (\mathbf{T}_t \otimes \mathbf{T}_y \otimes \mathbf{T}_x)$, where \otimes denotes Kronecker product [15]. Then we can write (3) concisely as $\mathbf{y} = \mathbf{\Phi} \mathbf{w} + \mathbf{e}$, where $\mathbf{\Phi} \in \mathbb{R}^{n_x n_y \times n_x n_y n_t}$ with $\mathbf{\Phi} \mathbf{\Phi}^T = \text{diag}(\text{vec}(\sum_{k=1}^{n_t} \mathbf{H}_k \odot \mathbf{H}_k))$. For simplification, from now we ignore possible noise \mathbf{e} . Note that \mathbf{y} reflects one $n_x \times n_y$ compressively measured image, as denoted at left in Figure 1, and $\mathbf{x} = (\mathbf{T}_t \otimes \mathbf{T}_y \otimes \mathbf{T}_x) \mathbf{w}$ is the $n_x \times n_y \times n_t$ video we wish to recover (Figure 1 right, for $n_t = 22$).

3.1. GAP for CS inversion

Let $\mathbf{\Phi} \in \mathbb{R}^{n_1 \times n}$, $\mathbf{w} \in \mathbb{R}^n$, and $\mathbf{y} \in \mathbb{R}^{n_1}$, with $n_1 < n$. Assume $\mathbf{\Phi}$ has full row rank. Let $\mathcal{G} = \{\mathcal{G}_1, \mathcal{G}_2, \dots, \mathcal{G}_m\}$ be a set of nonempty mutually-disjoint and collectively exhaustive subsets of $\{1, 2, \dots, n\}$. Let $\beta = [\beta_1, \dots, \beta_m]^T$ be a column of constant positive weights with β_k associated with \mathcal{G}_k . We solve the weighted- $\ell_{2,1}$ minimization problem

$$\min_{\mathbf{w}} \|\mathbf{w}\|_{\ell_{2,1}^{\mathcal{G}\beta}}, \text{ subject to } \mathbf{\Phi} \mathbf{w} = \mathbf{y}, \tag{4}$$

with $\|\mathbf{w}\|_{\ell_{2,1}^{\mathcal{G}\beta}} = \sum_{k=1}^m \beta_k \|\mathbf{w}_{\mathcal{G}_k}\|_2$, where $\mathbf{w}_{\mathcal{G}_k}$ is a sub-vector of \mathbf{w} containing components indexed by \mathcal{G}_k , and $\|\cdot\|_2$ denotes ℓ_2 norm; $\|\cdot\|_{\ell_{2,1}^{\mathcal{G}\beta}}$ is referred as a weighted- $\ell_{2,1}$ norm or $\ell_{2,1}^{\mathcal{G}\beta}$ norm. The groups and weights are below related to the anticipated importance of wavelet/DCT coefficients; the larger β_k , the more importance is placed on the k th group of coefficients being sparse.

The problem in (4) can be equivalently rewritten as

$$\min_{\mathbf{w}, C} C \text{ subject to } \|\mathbf{w}\|_{\ell_{2,1}^{\mathcal{G}\beta}} \leq C \text{ and } \mathbf{\Phi} \mathbf{w} = \mathbf{y} \tag{5}$$

Denote $B_{2,1}^{\mathcal{G}\beta}(C) = \{\mathbf{w} : \|\mathbf{w}\|_{\ell_{2,1}^{\mathcal{G}\beta}} \leq C\}$ and $\mathcal{S}_{\mathbf{\Phi}, \mathbf{y}} = \{\mathbf{w} : \mathbf{\Phi} \mathbf{w} = \mathbf{y}\}$, where $\mathcal{S}_{\mathbf{\Phi}, \mathbf{y}}$ is a given linear manifold and $B_{2,1}^{\mathcal{G}\beta}(C)$ is a weighted- $\ell_{2,1}$ ball with radius C . Geometrically, the problem in (5) is to find the smallest weighted- $\ell_{2,1}$ ball that has a nonempty intersection with the given linear manifold; we refer to this ball as the *critical ball* and denote its radius as C^* . When the smallest intersection is a singleton, the solution to (5) is unique.

We solve (5) as a series of alternating projection problems,

$$\begin{aligned}
(\mathbf{w}^{(t)}, \boldsymbol{\theta}^{(t)}) &= \arg \min_{\mathbf{w}, \boldsymbol{\theta}} \|\mathbf{w} - \boldsymbol{\theta}\|_2^2, \\
&\text{subject to } \|\boldsymbol{\theta}\|_{\ell_{2,1}^{\mathcal{G}\beta}} \leq C^{(t)} \text{ and } \mathbf{\Phi} \mathbf{w} = \mathbf{y}, \tag{6}
\end{aligned}$$

where a special rule is used to update $C^{(t)}$ to ensure that $\lim_{t \rightarrow \infty} C^{(t)} = C^*$. For each $C^{(t)}$, we solve an equivalent problem

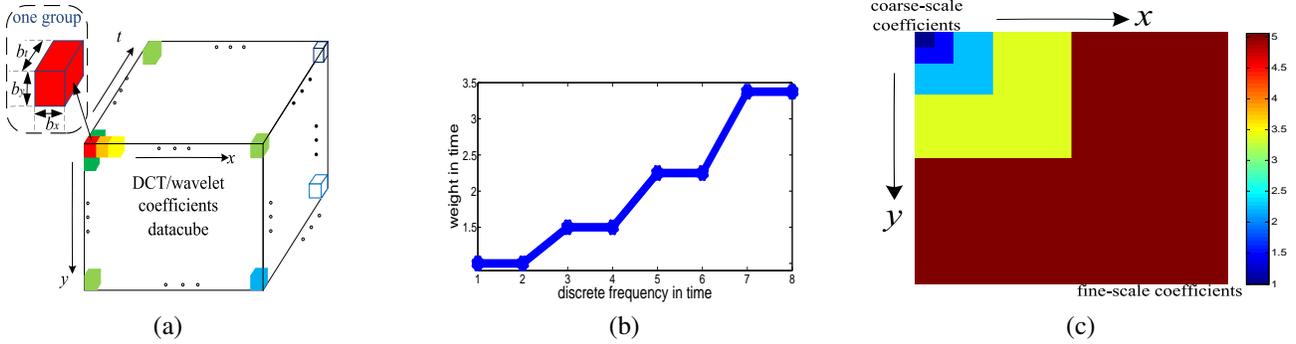


Figure 6. (a) Demonstration of groups of the wavelet/DCT coefficients of the video. The large 3D cube represents the wavelet/DCT coefficients and the small 3D cubes in different color denote different groups, with the size $\{b_x, b_y, b_t\}$ shown by one example group on the top-left. (b) DCT block structure time-weight. Each b_t (here $b_t = 2$) frames in time share the same weight. (c) Wavelet-tree structure scale-weight in space. The groups in the same wavelet level (shown in the same color) share the same weight.

$$\begin{aligned} (\mathbf{w}^{(t)}, \boldsymbol{\theta}^{(t)}) &= \arg \min_{\mathbf{w}, \boldsymbol{\theta}} \|\mathbf{w} - \boldsymbol{\theta}\|_2^2 + \lambda^{(t)} \|\boldsymbol{\theta}\|_{\ell_{2,1}^{\mathcal{G}\beta}} \\ \text{subject to } &\Phi \mathbf{w} = \mathbf{y}, \end{aligned} \quad (7)$$

where $\lambda^{(t)} \geq 0$ is the Lagrangian multiplier uniquely associated with $C^{(t)}$. Denote by λ^* the multiplier associated with C^* . It suffices to find a sequence $\{\lambda^{(t)}\}_{t \geq 1}$ such that $\lim_{t \rightarrow \infty} \lambda^{(t)} = \lambda^*$.

We solve (7) by alternately projection between \mathbf{w} and $\boldsymbol{\theta}$. Given one, the other is solved analytically: \mathbf{w} is an Euclidean projection of $\boldsymbol{\theta}$ on the linear manifold, while $\boldsymbol{\theta}$ is the result of applying group-wise shrinkage to \mathbf{w} . An attractive property of GAP is that, by using a special rule of updating $\lambda^{(t)}$, we only need to run a single iteration of (7) for each $\lambda^{(t)}$ to make $\{\lambda^{(t)}\}_{t \geq 1}$ converge to λ^* . In particular, GAP starts from $\boldsymbol{\theta}^{(0)} = \mathbf{0}$ and computes two sequences, $\{\boldsymbol{\theta}^{(t)}\}_{t \geq 1}$ and $\{\mathbf{w}^{(t)}\}_{t \geq 1}$:

$$\mathbf{w}^{(t)} = \boldsymbol{\theta}^{(t-1)} + \Phi^T (\Phi \Phi^T)^{-1} (\mathbf{y} - \Phi \boldsymbol{\theta}^{(t-1)}), \quad (8)$$

$$\boldsymbol{\theta}_{\mathcal{G}_k}^{(t)} = \mathbf{w}_{\mathcal{G}_k}^{(t)} \max \left\{ 1 - \frac{\lambda^{(t)} \beta_k}{\|\mathbf{w}_{\mathcal{G}_k}^{(t)}\|_2}, 0 \right\}, \forall k = 1, \dots, m \quad (9)$$

$$\text{where } \lambda^{(t)} = \left\| \mathbf{w}_{\mathcal{G}_{j_{m^*+1}^{(t)}}}^{(t)} \right\|_2 \beta_{j_{m^*+1}^{(t)}}^{-1}, \quad m^* < m \quad (10)$$

with $(j_1^{(t)}, j_2^{(t)}, \dots, j_m^{(t)})$ a permutation of $(1, 2, \dots, m)$

such that $\left\| \mathbf{w}_{\mathcal{G}_{j_1^{(t)}}}^{(t)} \right\|_2 \beta_{j_1^{(t)}}^{-1} \geq \left\| \mathbf{w}_{\mathcal{G}_{j_2^{(t)}}}^{(t)} \right\|_2 \beta_{j_2^{(t)}}^{-1} \geq \dots \geq \left\| \mathbf{w}_{\mathcal{G}_{j_m^{(t)}}}^{(t)} \right\|_2 \beta_{j_m^{(t)}}^{-1}$.

The algorithm (8)-(9) is referred as *generalized alternating projection* (or GAP) [13] to emphasize its difference from alternating projection (AP) in the conventional sense: conventional AP produces a sequence of projections between two *fixed* convex sets, while GAP produces a sequence of projections between two convex sets that undergo systematic changes over the iterations. In the GAP algorithm as shown in (8)-(9), the alternating projection is per-

formed between a fixed linear manifold $\mathcal{S}_{\Phi, \mathbf{y}}$ and a changing weighted- $\ell_{2,1}$ ball, i.e., $B_{2,1}^{\mathcal{G}\beta}(C^{(t)})$ whose radius $C^{(t)}$ is a function of the iteration number t .

By interpreting $\frac{1}{\lambda^{(t)}} \|\boldsymbol{\theta} - \mathbf{w}\|_2^2$ as a penalty to enforce $\boldsymbol{\theta} = \mathbf{w}$, one may view that iteration of (7) with t constituting a penalty method for solving the following constrained problem,

$$\min_{\boldsymbol{\theta}} \|\boldsymbol{\theta}\|_{\ell_{2,1}^{\mathcal{G}\beta}} \quad \text{subject to } \boldsymbol{\theta} = \mathbf{w} \text{ and } \Phi \mathbf{w} = \mathbf{y}, \quad (11)$$

which is an equivalent formulation of (4). The GAP algorithm in (8)-(10) is a special penalty method for solving (11), using (10) to adjust the penalty strength $\{\lambda^{(t)}\}$. Bregman iteration [21] can also solve (4) or (11). However, Bregman penalizes $\|\mathbf{y} - \Phi \mathbf{w}\|_2^2$, while GAP keeps $\mathbf{y} = \Phi \mathbf{w}$ as a constraint and fulfills it via the orthogonal projection in (8). Under a set of group-restricted isometry property (group-RIP) conditions, the use of (8)-(10) ensures monotonic decrease of the reconstruction error to zero and makes GAP an anytime algorithm [13]. By contrast, Bregman iteration and classic penalty method (which adjusts $\lambda^{(t)}$ differently) do not have the anytime property, nor do other popular algorithms such as TwIST and ADMM [5].

3.2. Extension of GAP for the proposed camera

The diagonalization of $\Phi \Phi^T$ is the key to fast GAP recovery of video. The inversion of $\Phi \Phi^T$ in (8) now just requires computation of the reciprocals of the diagonal elements, as a result of the *hardware implementation* of the proposed camera. Best results were found when \mathbf{T}_x and \mathbf{T}_y correspond to a wavelet basis (here the Daubechies-8 [15]), and \mathbf{T}_t corresponds to a DCT. The basis-function weights (β) are defined with respect to these bases, and the groups are manifested in the domain of these wavelet-DCT coefficients (see Figure 6(a) for a depiction of the groups). For \mathbf{T}_t the coefficients are arranged from low frequencies to high frequencies in Figure 6(b), and for \mathbf{T}_x and \mathbf{T}_y the 2D arrangement of coefficients is as is done typically with wavelets [15], and illustrated in Figure 6(c). Let $\{b_x, b_y, b_t\}$ represent the edge lengths of each group of coefficients (Figure 6(a)). In all experiments, $b_x = b_y = 2$ and $b_t = \lfloor \frac{n_t}{4} \rfloor$, where $\lfloor \cdot \rfloor$ denotes the closest integer of the

number inside $[\]$. The weight is constituted by the product of a time-weight associated with the DCT (Figure 6(b)) and a scale-weight associated with the wavelets (Figure 6(c)). We show the details of the weight design in the following.

Let $\{g_x, g_y, g_t\}$ index each group with $g_x = 1, \dots, \frac{n_x}{b_x}$, $g_y = 1, \dots, \frac{n_y}{b_y}$, $g_t = 1, \dots, \frac{n_t}{b_t}$. The time weights are defined by $\beta_t(g_t) = a^{(g_t-1)}$. For the scale weights, we exploit the wavelet-tree structure (Figure 6(c)), and enforce the groups in the ℓ th level (assuming L levels in total, and $\ell = 1, \dots, L$) of the wavelet coefficients sharing the same weight, determined by $\beta_{x,y}(g_x^\ell, g_y^\ell) = a^{(\ell-1)}$, $\forall g_x^\ell = \frac{n_x^{(\ell-1)}}{b_x} + 1, \dots, \frac{n_x^\ell}{b_x}$, $g_y^\ell = \frac{n_y^{(\ell-1)}}{b_y} + 1, \dots, \frac{n_y^\ell}{b_y}$, with $\{n_x^\ell, n_y^\ell, n_t^\ell\}$ denoting the end points of wavelet coefficients at the ℓ th level, and $n_x^0 = n_y^0 = 0$. The weight for each 3D group is $\beta(g_x^\ell, g_y^\ell, g_t) = \beta_{x,y}(g_x^\ell, g_y^\ell)\beta_t(g_t)$. Setting $a = 1.5$ was found to yield good results. After constructing the groups and weights as above, the GAP performance is improved significantly in the application here.

3.3. Temporal overlap in inversion

In Section 3.1, each coded CS measurement \mathbf{Y}_l is employed to recover n_t frames of video. This may lead to discontinuities in the video recovered from consecutive CS measurements. To mitigate this, we also consider the *joint* inversion of two consecutive CS measurements, \mathbf{Y}_l and \mathbf{Y}_{l+1} , from which $2n_t$ consecutive frames are recovered at once. Therefore the n_t frames associated with \mathbf{Y}_l are estimated jointly from $\mathbf{Y}_{l-1} \cup \mathbf{Y}_l$ and (separately) from $\mathbf{Y}_l \cup \mathbf{Y}_{l+1}$. The final recovered video within a particular contiguous set of n_t frames is taken as the average of the results inferred from $\mathbf{Y}_{l-1} \cup \mathbf{Y}_l$ and $\mathbf{Y}_l \cup \mathbf{Y}_{l+1}$. As demonstrated below, this tends to improve the quality of the recovered video (yields smoother results).

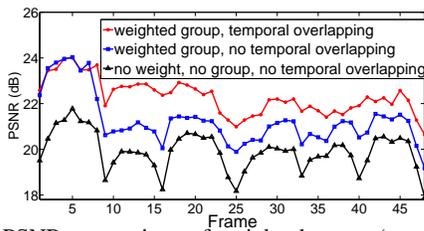


Figure 7. PSNR comparison of weighted group (average PSNR: 22.13dB), no weight, no group (average PSNR: 19.92dB), and temporal overlapping weighted group (average PSNR: 22.80dB).

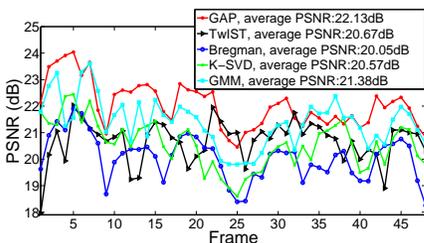


Figure 8. PSNR comparison of GAP, TwiST, linearized Bregman, K-SVD, and GMM algorithms with the simulation dataset.

4. Simulation results

To demonstrate the performance of the reconstruction algorithm, we start with simulated data (see Section 5 for real data from the proposed camera), and employ a color video sequence in which a basketball player performs a dunk [18]; this video is challenging due to the complicated motion of the basketball players and the varying lighting conditions and multiple depths scene; see the example video frames in Figure 2(a). We consider a binary mask, with 1/0 coding drawn at random Bernoulli(0.5); the code is shifted spatially via the coding mechanism in Figure 2(a), as in our physical camera. The video frames are 256×256 spatially, and we choose $n_t = 8$.

We first investigate the efficacy of weighted groups in the proposed GAP algorithm. Figure 7 demonstrates the improvement by the weighting and grouping of the wavelet/DCT coefficients; these parameters (weights and groups) were not optimized, and many related settings yield similar results – there is a future opportunity for optimization. In Figure 7 we also demonstrate the performance improvement manifested by temporal overlapping and averaging results from two consecutive measurements. Note that *without* temporal overlapping, the PSNR degrades for frames at the beginning (*e.g.*, 1, 9, 17, ...) and end (8, 16, 24, ...) of a given measurement, while the PSNR curve with temporal overlapping (red curve) is much smoother. The experiments with (real) measured data consider temporal overlapping when performing inversion.

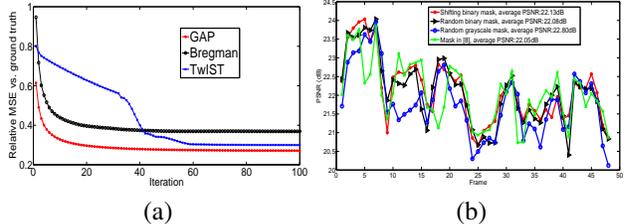


Figure 9. (a) Convergence comparison of GAP, TwiST, and linearized Bregman. (b) PSNR of reconstructed video frames with three different coding scenarios. “Shifting binary mask” means the coding mechanism proposed in our camera. “Random binary mask” denotes that for every frame, each element of the mask is randomly drawn from Bernoulli(0.5). “Random grayscale mask” symbolizes each element of the mask is randomly drawn from $[0, 1]$ for every frame.

4.1. Reconstruction performance and convergence

The developed GAP algorithm is compared with the following: *i*) two-step iterative shrinkage/thresholding (TwiST) [4] (with total variation norm), *ii*) K-SVD [3] with orthogonal matching pursuit (OMP) [17] used for inversion, *iii*) a Gaussian mixture model (GMM) based inversion algorithm [7, 20, 22], and *iv*) the linearized Bregman algorithm [21]. The ℓ_1 -norm of DCT or wavelet coefficients is adopted in linearized Bregman algorithm with the same transformation as GAP. GMM and K-SVD are patch-based algorithms and we used a separate dataset for training. A

batch of training videos were used (shown in the website [18]) to pre-train K-SVD and GMM, and we selected the best reconstruction results for presentation here. The PSNR curves of videos reconstructed by GAP and the four alternative algorithms are shown in Figure 8, demonstrating the GAP algorithm outperforms the other algorithms by 0.7-2dB higher in PSNR. In this simulation, the temporal overlap was not used. The readers can refer to the reconstructed videos on the website [18].

We further investigate the convergence properties of GAP, TwIST and linearized Bregman, as they are each solved iteratively. We run each algorithm a total of 100 iterations and compute the relative mean square error (MSE) of the estimate compared with the ground truth for each iteration. Relative MSE versus iteration number are plotted in Figure 9(a). It can be seen that GAP and linearized Bregman converge much faster than TwIST, while GAP provides the smallest relative MSE in every iteration. We also verify the *anytime* property of GAP by computing the first-order difference of the MSE for each iteration (compared with TwIST and Bregman); see Supplementary Material [18].

4.2. Comparison with other coding mechanisms

One advantage of the proposed coded aperture compressive camera (manifested by spatially shifting a *single* binary mask), compared with others [8, 16], is its low-power, low-bandwidth characteristics. However, the use of a single shifted binary mask to yield temporal coding may appear limiting, and therefore it is of interest to compare to other strategies. In Figure 9(b) we compare the proposed translated coding mechanism in our hardware system with even more general coding strategies than in [8, 16]. Specifically, we compare to a unique random binary mask at each time point, for which *each* of the n_t codes is a *distinct* i.i.d. draw $\text{Bernoulli}(0.5)$. We also consider the case for which each code/mask element, for each of the n_t codes, is drawn uniformly at random from $[0,1]$, reflecting the degree of code transmission (each of the n_t codes is distinct, and each is not restricted to be binary). Summarizing Figure 9(b), the simple shifted binary code in the proposed system yields similar results (even a little higher PSNR for some datasets) compared to these alternative coding strategies, at an order of magnitude less power. We found similar results when comparing to the coding strategies in [8, 16].

5. Experimental results: real data

The physical camera we have developed captures the measurements (coded capture) at $1/\Delta_t = 30$ fps, and the reconstructed video has 660 fps ($n_t = 22$, although different n_t may be considered in the inversion). One result from this camera was shown in Figure 1 (recovery of high-speed motion). From the reconstructed frames in Figure 1, one can clearly identify the spin of the red apple and the rebound of the yellow orange; the full video is at [18], along with many additional examples. At [18], we show compar-

isons to a diverse set of alternative algorithms, for example via [21].

All algorithms considered here were implemented in Matlab 2012(a), and the experiments are conducted on a PC with a CPU@3.30 GHz and 16GB RAM. For the real data ($512 \times 512 \times 22$), GAP, TwIST, and linearized Bregman use 50, 100, and 500 iterations, respectively (required to yield good results). Each iteration in these three algorithms are similar (around 0.8 seconds). Hence, TwIST and linearized Bregman cost much longer ($> 2\times$) time than GAP, but typically provide worse results, and do not have an anytime property. K-SVD and GMM may be made fast if parallel computing is used (processing the multiple patches in parallel with GPU or networks), but for serial computing on a PC like that considered here these methods are slower than GAP.

5.1. Recovery of depth and motion simultaneously

Figure 3 shows the reconstruction frames (6 out of 14 frames are shown for demonstration) recovered from one measurement. The three objects, “newspaper,” “smile-face” and “chopper-wheel” are in three different depths. The motion of the *rotating* chopper-wheel is also recovered from the reconstructed frames. The first column of the recovered frames (bottom right of Figure 3) shows the newspaper is in focus; the second column shows the smile-face is in focus, and finally, the chopper-wheel is in focus at column three. We have built a depth-frames relationship table with calibration. The newspaper is best in focus in frame 3, which corresponds to 14cm away from the objective lens (truth 15cm). The smile-face is best in focus in frame 8, corresponding to 40cm (truth 38cm). The chopper-wheel is best in focus in frame 12, which corresponds to 64cm (truth 65cm). It can be seen the depths of these objects are estimated correctly.

5.2. Recovery of high-speed motion

As one additional example of high-speed motion recovery (assuming all the objects are in focus, *i.e.*, here without the liquid lens), Figure 10 shows the results of a purple hammer quickly hitting a red apple. In this dataset, 44 frames are reconstructed from 2 measurements (left part) showing the entire process of hitting (full video is at [18]), and 5 example frames out of 44 are plotted in the right part of Figure 10. To demonstrate the *anytime* property of GAP, we show the results after 2, 10, 20 and 50 iterations. Note that good results are manifested with as few as 10 iterations, with convergence after about 20.

5.3. Recovery of depth

When the scene is not moving, we can get space-depth-color information from the reconstructed data. An example is shown in Figure 11. We can see that the smile-face is first out-of-focus, then in-focus and finally out-of-focus again.

6. Conclusions

This paper proposes a means of recovering depth, time, and color information from a single coded image, via de-

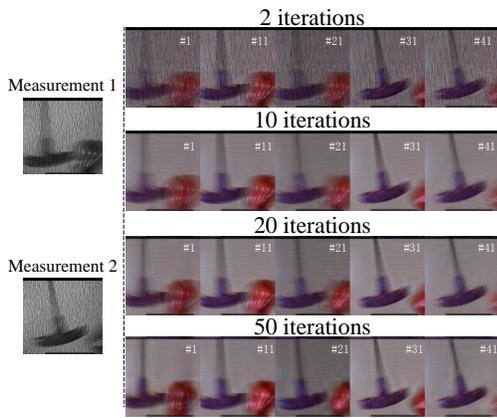


Figure 10. Reconstruction from two real measurements (left part). Five selected frames out of the 44 reconstructed video frames ($n_t = 22$ recovered video frames from each compressive measurement) after 2, 10, 20, and 50 iterations of GAP are shown in the right part.

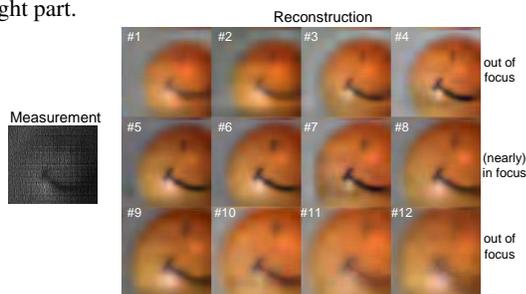


Figure 11. Reconstruction from a real single measurement (left part). 12 selected frames are shown in the right part to demonstrate the process of the smile-face from out-of-focus to in focus and then out-of-focus again.

velopment of a new color CS camera for high-speed depth-video reconstruction. In the presented computational time comparisons, GAP was run in MATLAB, since that was the language in which all of the comparison algorithms had available code, and therefore provided a good comparison point for *relative* speed. In the context of *absolute* speed, we have implemented GAP in C++ on a GPU, and the *total* time for reconstructing a $512 \times 512 \times 22$ video from a single CS measurement is less than 0.5 seconds, opening the door for real-time fast (3D) video capturing and reconstruction.

Acknowledgements

This work was supported by DARPA under the KECOM program, and by the DOE, NGA, ONR, NSF and ARO.

References

- [1] AVT Guppy PRO, technical manual. *Allied Vision Technologies, v 4.0*, 2012.
- [2] J. Adams, K. Parulski, and K. Spaulding. Color processing in digital cameras. *IEEE Micro*, 1998.
- [3] M. Aharon, M. Elad, and A. Bruckstein. K-SVD: An algorithm for designing overcomplete dictionaries for sparse representation. *IEEE T-SP*, 2006.
- [4] J. Bioucas-Dias and M. Figueiredo. A new TwIST: Two-step iterative shrinkage/thresholding algorithms for image restoration. *IEEE T-IP*, 2007.

- [5] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *FTML*, 2011.
- [6] E. Candès, J. Romberg, and T. Tao. Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information. *IEEE T-IT*, 2006.
- [7] M. Chen, J. Silva, J. Paisley, C. Wang, D. Dunson, and L. Carin. Compressive sensing on manifolds using a non-parametric mixture of factor analyzers: Algorithm and performance bounds. *IEEE T-SP*, 2010.
- [8] Y. Hitomi, J. Gu, M. Gupta, T. Mitsunaga, and S. K. Nayar. Video from a single coded exposure photograph using a learned over-complete dictionary. *ICCV*, 2011.
- [9] J. Holloway, A. C. Sankaranarayanan, A. Veeraraghavan, and S. Tambe. Flutter shutter video camera for compressive sensing of videos. In *ICCP*, 2012.
- [10] S. Kuthirummal, H. Nagahara, C. Zhou, and S. K. Nayar. Flexible depth of field photography. *IEEE T-PAMI*, 2010.
- [11] A. Levin. Analyzing depth from coded aperture sets. *ECCV*.
- [12] A. Levin, R. Fergus, F. Durand, and W. T. Freeman. Image and depth from a conventional camera with a coded aperture. *ACM SIGGRAPH*, 2007.
- [13] X. Liao, H. Li, and L. Carin. Generalized alternating projection for weighted $\ell_{2,1}$ minimization with applications to model-based compressive sensing. *SIIMS*, 2014.
- [14] P. Llull, X. Liao, X. Yuan, J. Yang, D. Kittle, L. Carin, G. Sapiro, and D. J. Brady. Coded aperture compressive temporal imaging. *Optics Express*, 2013.
- [15] S. Mallat. *A Wavelet Tour of Signal Processing: The Sparse Way*. Academic Press, 2008.
- [16] D. Reddy, A. Veeraraghavan, and R. Chellappa. P2C2: Programmable pixel compressive camera for high speed imaging. *CVPR*, 2011.
- [17] J. A. Tropp and A. C. Gilbert. Signal recovery from random measurements via orthogonal matching pursuit. *IEEE T-IT*, 2007.
- [18] URL. <https://sites.google.com/site/2014cvpr>.
- [19] A. Veeraraghavan, D. Reddy, and R. Raskar. Coded strobing photography: Compressive sensing of high speed periodic videos. *IEEE T-PAMI*, 2011.
- [20] J. Yang, X. Yuan, X. Liao, P. Llull, D. Brady, G. Sapiro, and L. Carin. Video compressive sensing using gaussian mixture models. *IEEE-TIP*, 2014.
- [21] W. Yin, S. Osher, D. Goldfarb, and J. Darbon. Bregman iterative algorithms for ℓ_1 -minimization with applications to compressed sensing. *SIIMS*, 2008.
- [22] G. Yu, G. Sapiro, and S. Mallat. Solving inverse problems with piecewise linear estimators: From Gaussian mixture models to structured sparsity. *IEEE T-IP*, 2012.
- [23] X. Yuan, J. Yang, P. Llull, X. Liao, G. Sapiro, D. J. Brady, and L. Carin. Adaptive temporal compressive sensing for video. *ICIP*, 2013.
- [24] C. Zhou, S. Lin, and S. K. Nayar. Coded aperture pairs for depth from defocus and defocus deblurring. In *ICCV*, 2009.
- [25] C. Zhou, S. Lin, and S. K. Nayar. Coded aperture pairs for depth from defocus. *IJCV*, 2011.