



جامعة الملك عبد الله  
للعلوم والتقنية  
King Abdullah University of  
Science and Technology

## End-to-End Video Compressive Sensing Using Anderson-Accelerated Unrolled Networks

Item Type	Conference Paper
Authors	Li, Yuqi; Qi, Miao; Gulve, Rahul; Wei, Mian; Genov, Roman; Kutulakos, Kiriakos N.; Heidrich, Wolfgang
Citation	Li, Y., Qi, M., Gulve, R., Wei, M., Genov, R., Kutulakos, K. N., & Heidrich, W. (2020). End-to-End Video Compressive Sensing Using Anderson-Accelerated Unrolled Networks. 2020 IEEE International Conference on Computational Photography (ICCP). doi:10.1109/iccp48838.2020.9105237
Eprint version	Post-print
DOI	<a href="https://doi.org/10.1109/ICCP48838.2020.9105237">10.1109/ICCP48838.2020.9105237</a>
Publisher	Institute of Electrical and Electronics Engineers (IEEE)
Rights	Archived with thanks to IEEE
Download date	05/08/2022 07:13:27
Link to Item	<a href="http://hdl.handle.net/10754/663826">http://hdl.handle.net/10754/663826</a>

# End-to-End Video Compressive Sensing Using Anderson-Accelerated Unrolled Networks

Paper ID 69

**Abstract**—Compressive imaging systems with spatial-temporal encoding can be used to capture and reconstruct fast-moving objects. The imaging quality highly depends on the choice of encoding masks and reconstruction methods. In this paper, we present a new network architecture to jointly design the encoding masks and the reconstruction method for compressive high-frame-rate imaging. Unlike previous works, the proposed method takes full advantage of a denoising prior to provide a promising frame reconstruction. The network is also flexible enough to optimize full-resolution masks and efficient at reconstructing frames. To this end, we develop a new dense network architecture that embeds Anderson acceleration, known from numerical optimization, directly into the neural network architecture.

Our experiments show the optimized masks and the dense accelerated network respectively achieve 1.5 dB and 1 dB improvements in PSNR without adding training parameters. The proposed method outperforms other state-of-the-art methods both in simulations and on real hardware. In addition, we set up a coded two-bucket camera for compressive high-frame-rate imaging, which is robust to imaging noise and provides promising results when recovering nearly 1,000 frames per second.

**Index Terms**—high-frame-rate imaging, deep neural network, computational camera

## 1 INTRODUCTION

As a well-developed technique, compressive sensing (CS) is widely applied in reconstructing images with low sampling rates [1], [2]. In particular, a variety of mask-based CS cameras have been demonstrated for capturing high-dimensional image data (e.g., spectra, video, etc.) using a two-dimensional camera with encoding capacity. Compared to conventional cameras employing brute-force sampling strategies, such CS cameras have significant advantages in acquisition efficiency, storage consumption, and potentially cost [3], [4].

High-frame-rate imaging is concerned with recording videos at rates in excess of hundreds of frames per second. However, with bandwidth being a limiting factor, conventional cameras record either a very low spatial resolution with a relatively high frame rate, or at relatively high spatial resolution with a low frame rate. Using mask-based compressive sensing, it becomes feasible to capture high-frame-rate and high-spatial-resolution videos with an efficient spatio-temporal encoding. This approach is a good fit for recently developed image sensors with high-speed per-pixel programmable exposure control [5]. The exposure control can be viewed as an encoding of the captured frames with a set of binary temporal masks. With such cameras, it is possible to encode multiple *subframes* into a captured image and decode them later using frame reconstruction methods (Fig. 1).

Much research has focused on the improvement of the reconstruction techniques, usually by employing optimization-based approaches (see Section 2 for more detail). Less work has concentrated on the derivation of good encoding masks: it can be shown that optimal mask selection in CS is NP-complete, but random (Bernoulli or

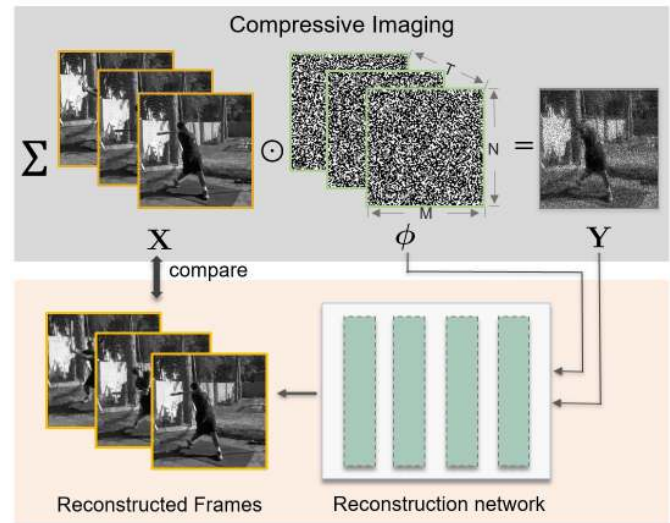


Fig. 1. Illustration of the encoding and reconstruction within the compressive high-frame-rate imaging system. In the system,  $T$  subframes with resolution  $M \times N$  are encoded with masks  $\phi$ . The reconstruction network reconstructs the frames from the measurement  $Y$  and the known mask  $\phi$ .

Gaussian) patterns are satisfactory with high probability [6]. However, the encoding and decoding components of the imaging system are highly interdependent. Based on this observation, we focus on the joint end-to-end design of encoding masks and reconstruction methods for improving both encoding efficiency and reconstruction accuracy. We put forward a compact end-to-end neural network that can handle the mask optimization for the whole image with fewer training parameters. We also show that this network design corresponds to Anderson acceleration, a well-known

• This paper is under review for ICCP 2020 and the PAMI special issue on computational photography. Do not distribute.

acceleration technique in numerical optimization [7].

Both simulations and experiments on real hardware show that our network outperforms existing methods. In addition, we show that our masks can also improve the reconstruction quality of existing methods. Our contributions can be summarized as follows:

- We present the *first work* to jointly design full-resolution coding masks and reconstruction methods for compressive high-frame-rate imaging using an end-to-end network. Our approach outperforms state-of-art methods by 2.2dB in PSNR.
- We show that the acceleration of the gradient descent algorithm is equivalent to adding dense skip connections to iterative optimization-unrolling neural networks. This speeds up training convergence and helps to design a compact and efficient *network architecture*.
- Experiments on both simulation and *real hardware* demonstrate the effectiveness of our reconstruction method and the designed masks. The two-bucket design of our camera shows improved noise suppression and can provide promising results in reconstructing video of frame rates up to almost 1,000 frames per second.

## 2 RELATED WORK

Many approaches have been developed to solve the ill-conditioned inverse problem in CS. The existing methods can be divided into model-based optimization methods, deep discriminative learning methods, and unrolled iterative optimization methods.

### *Model-based methods.*

Model-based methods utilize designed image priors for regularization, which can reduce the number of possible solutions and remove artifacts in frame reconstruction. For example, the Total Variation (TV) prior [4], [8] can simultaneously preserve edges while smooth away noise in flat regions; optical flow [9] can estimate the motion of moving objects and helps to eliminate ghosting effects; Gaussian mixture models [10] and dictionary learning methods [11], [12] take into account image statistics and reconstruct frames using learned atoms; non-local low-rank priors [13] consider correlation between small patches in the frames for denoising. Such model-based methods are straightforward to adapt to different sensing matrices without retraining, and the sensing matrix can be optimized based on the analysis of mutual coherence in dictionary-learning based methods [14]. However, such model-based methods have their respective drawbacks, and none of them is suitable for all scenes. In addition, these methods can be computationally expensive, especially compared to learning-based methods.

### *Learning-based methods.*

In recent years, deep discriminative learning methods have shown drastic improvements in image reconstruction quality. Some deep neural networks (DNNs) have been proposed for compressive imaging as well. Convolutional neural networks [15], [16], [17] and fully-connected networks [18],

[19] were developed to reconstruct small image patches. However, none of the convolutional networks are capable of simultaneously designing masks and optimizing parameters in the network. Compared to model-based methods, these DNN-based methods are efficient but difficult to adapt to different sensing matrices. These networks usually use random code masks, such as Gaussian or Bernoulli random masks [20], and thus cannot achieve optimal reconstruction quality. On the other hand, fully connected networks, suffer from a large search space, and can in practice only optimize a small repeated mask by preserving the essential connections. While repeated masks significantly reduce the scale of the optimization problem, they may also introduce structured artifacts during reconstruction.

### *Unrolling iterative optimization methods.*

More recently, a class of networks constructed by unrolling iterative optimization methods has started to be used in image reconstruction (e.g. LISTA [21]ADMM-net [22], LDAMP [23], IRCNN [24], ISTA-Net [25]). Such network architectures combine the advantages of both model-based methods and deep discriminative learning methods, and provide an efficient and flexible plug-and-play framework to solve inverse problems. Previous works have utilized the multistage iterative network for image restoration [26] and illumination optimization [27]. In this paper, we claim that such networks are effective in jointly optimizing the sensing matrix and reconstruction method if the elements of the sensing matrix are treated as trainable parameters in the network. Crucially, we also show how to improve the design of such unrolled networks to embed Anderson acceleration directly into the network architecture. This improvement will be applicable and useful far beyond our specific application scenario.

### *Computational video cameras.*

Many different prototype designs for computational video cameras have been proposed. Raskar *et al.* modified a conventional DSLR camera and added a control unit for high-speed control of the exposure pattern over the full frame. The camera can then be used for deblurring [28] and video compressive sensing [29]. Liu *et al.* used an LCoS to implement a single exposure mask and applied dictionary learning to reconstruct the scene [30]. To achieve a high-speed encoding, Bub *et al.* used a DMD for high-frame-rate imaging [31]. Llull *et al.* changed from active to passive codes to reduce the power consumption [4]. In their design, the static mask is spatially shifted over time, which provides a very limited design space for the spatio-temporal encoding.

Recently, several image sensor designs have been proposed that can implement the CS mask directly on the sensor. Luo *et al.* [32] invented a CMOS sensor that allows for active control of the exposure pattern in each pixel, and applied this design for image deblurring. Zhang *et al.* [33] a CMOS sensor for both high-speed and high-dynamic-range imaging [34]. However, since there is no charge bucket connect with PD, every pixel can only expose once during a frame. Sonoda *et al.* [35] built a sensor with quasi pixel-wise programmable control, but pixels on the sensor can only be controlled in blocks. Therefore, their camera cannot

generate arbitrary mask patterns [36]. Sarhangnejad *et al.* [37] implemented a coded-exposure-pixel camera with two-bucket pixels that has 180 subframes per second. In this camera every pixel is programmable and can be exposed many times during a single frame. Wei *et al.* use this system for a one-shot photometric stereo and develop an image formation model for computational video cameras [5].

### 3 METHOD

Our goal is to jointly learn both the full-resolution masks for encoding and the reconstruction method for decoding that together minimize subframe reconstruction error. We achieve this by training an end-to-end network that consists of  $K$  stages with dense skip connections and a mask layer, as shown in Fig. 4. Given a video sequence, the mask layer modulates each subframe using the learned mask and integrates all subframes into a single captured image; the  $K$  stages constructed via unrolling the optimization iterations for reconstruction can decode the captured images into multiple subframes.

In the following, we first present the encoding and decoding parts of our neural network architecture along with training details. Then we describe a set of simulations for comparing the proposed method with other existing methods. Lastly, we implement our approach on a real camera and evaluate the effectiveness of our network.

#### Image formation.

The image formation model for our compressive video capture system is shown in Fig.1, and can be formulated as:

$$\mathbf{Y} = \sum_{i=1}^T \phi^{(i)} \odot \mathbf{X}^{(i)} + \mathbf{N}, \quad (1)$$

where  $\phi^{(i)} \in \mathbb{R}^{M \times N}$  denotes the  $i$ -th binary encoding mask,  $\mathbf{X}^{(i)} \in \mathbb{R}^{M \times N}$  represents the  $i$ -th subframe we need to reconstruct,  $\odot$  denotes the element-wise product,  $\mathbf{N} \in \mathbb{R}^{M \times N}$  denotes the imaging noise, and  $\mathbf{Y}$  is the  $M \times N$  captured image. The system has a compression ratio of  $1/T$ , i.e.  $T$  successive subframes are encoded into a single captured image.

Eq. 1 can be transformed into the following equation:

$$\mathbf{y} = \Phi \mathbf{x} + \mathbf{n}, \quad (2)$$

where  $\Phi \in \mathbb{R}^{MN \times TMN}$  is the sensing matrix with diagonal blocks consisting of the masks  $\phi$ :

$$\Phi = [\text{diag}(\text{Vec}(\phi^{(1)})), \dots, \text{diag}(\text{Vec}(\phi^{(T)}))], \quad (3)$$

$\mathbf{x}$  represents the  $TMN \times 1$  vectorized subframes of  $\mathbf{X}$ ,  $\mathbf{y}$  is the  $MN \times 1$  vectorized captured image of  $\mathbf{Y}$ , and  $\mathbf{n}$  denotes the vectorized noise of  $\mathbf{N}$ .

#### 3.1 Mask generation

A layer containing only bias values is constructed to generate the encoding masks  $\phi$ . Since different pixels in the subframes are encoded independently, the operation  $\Phi \mathbf{x}$  can be realized by an element-wise multiplication of  $\phi$  and  $\mathbf{X}$  and a summation of the multiplication results; the operation  $\Phi^T \mathbf{y}$  can be realized by a repeat copy operation

of  $\mathbf{Y}$  and an element-wise multiplication, as shown in Fig.3. The two operations are beneficial for efficient calculation, as well as reduced storage requirements. Since the masks used in high-frame-rate imaging are binary, we need to add a constraint that the outputs of the mask layer must be either 0 or 1 during propagation. Inspired by the Binaryconnect method [38], this can be achieved by a simple but efficient deterministic binarization operation:

$$\hat{b} = \begin{cases} 1, & \text{when } b > 0, \\ 0, & \text{else.} \end{cases}, \quad (4)$$

where  $\hat{b}$  is the binarized value of the mask layer, and  $b$  is the real value. The sign function binarizes the values straightforwardly, however it is only activated during the forward and backward propagations but not during the parameter update since it is necessary to maintain good precision weights during the updates.

#### 3.2 Subframe reconstruction

##### Unrolled network reconstruction.

To present the subframe reconstruction method, we first mathematically formulate the reconstruction procedure as an unconstrained problem, and then loop-unroll the optimization to construct our multi-stage network. Subframe reconstruction is an optimization problem

$$\arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \Phi \mathbf{x}\|^2 + \lambda J(\mathbf{x}), \quad (5)$$

where  $J(\mathbf{x})$  is the denoising prior for regularization weighted by parameter  $\lambda$ . The first data fidelity term guarantees a minimal re-sensing error while the regularization term ensures that the reconstructed frames satisfy the desired prior model. Different from designed priors in model-based method, denoising prior depicts intrinsic statistics of images and results in better image reconstruction. ensures that the reconstructed subframes satisfy the desired prior model. Different from the hand-designed priors of model-based methods, the deep image prior captures the intrinsic statistics of images and results in better image reconstructions.

By introducing an auxiliary variable  $\mathbf{v}$ , Eq. 5 can be reformulated as a constrained optimization problem:

$$(\mathbf{x}, \mathbf{v}) = \arg \min_{\mathbf{x}, \mathbf{v}} \frac{1}{2} \|\mathbf{y} - \Phi \mathbf{x}\|^2 + \lambda J(\mathbf{v}), \text{ st. } \mathbf{x} = \mathbf{v}. \quad (6)$$

Inspired by previous image restoration works [24], we adopt the half-quadratic splitting method to convert the constrained optimization problem into an unconstrained one:

$$(\mathbf{x}, \mathbf{v}) = \arg \min_{\mathbf{x}, \mathbf{v}} \frac{1}{2} \|\mathbf{y} - \Phi \mathbf{x}\|^2 + \frac{\tau}{2} \|\mathbf{x} - \mathbf{v}\|^2 + \lambda J(\mathbf{v}), \quad (7)$$

where  $\tau$  is a weight term. Then, Eq. 7 can be solved by alternatively optimizing the two sub-problems with respect to  $\mathbf{z}$  and  $\mathbf{x}$ , respectively:

$$\begin{cases} \mathbf{x}^{i+1} &= \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \Phi \mathbf{x}\|^2 + \frac{\tau}{2} \|\mathbf{x} - \mathbf{v}^i\|^2 \\ \mathbf{v}^{i+1} &= \arg \min_{\mathbf{v}} \frac{\tau}{2} \|\mathbf{x}^{i+1} - \mathbf{v}\|^2 + \lambda J(\mathbf{v}) \end{cases} \quad (8)$$

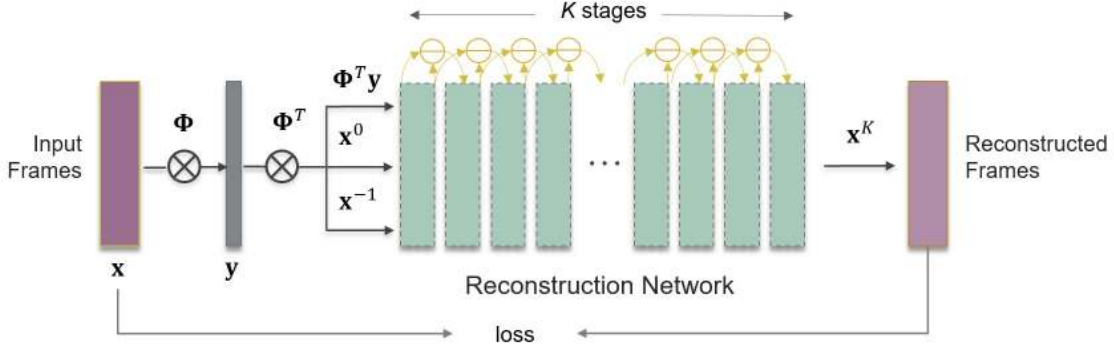


Fig. 2. Our deep network architecture. The overall network consists of a mask layer for generating masks and K stages for reconstruction. Note that the skip connections of residuals among stages make the network denser and more compact. (Here show is the case where the number of skip connections of each stage is  $m = 1$ .)

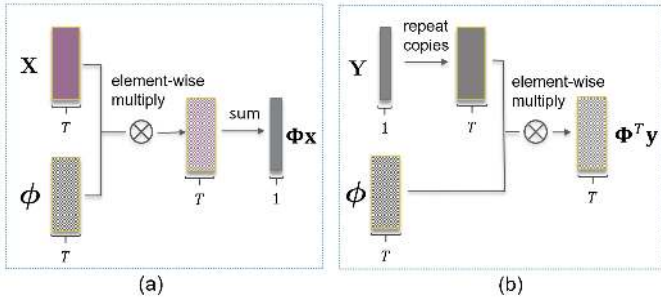


Fig. 3. Two matrix-vector multiplication operations: (a)  $\Phi \mathbf{x}$  and (b)  $\Phi^T \mathbf{y}$ .

By analyzing Eq. 8, it is evident that the optimization of  $\mathbf{x}$  in the first line is a quadratic problem, while optimization of  $\mathbf{v}$  in the second line is actually a denoising problem. To solve the first problem, we can calculate the closed-form solution

$$\mathbf{x}^{i+1} = (\Phi^T \Phi + \tau \mathbf{I})^{-1} (\Phi^T \mathbf{y} + \tau \mathbf{v}^i). \quad (9)$$

However, the matrix inversion is time consuming. More importantly, such inverse models consisting of the trainable sensing matrix  $\Phi$  are harder to train, compared to a forward model of  $\Phi$ . Previous work [26] suggests that using gradient descent algorithms to obtain an inexact solution in each step can also effectively and efficiently optimize the problem. In the general gradient descent method, the update step of  $\mathbf{x}$  can be performed as:

$$\begin{aligned} \mathbf{x}^{i+1} &= \mathbf{x}^i - \alpha^i g(\mathbf{x}^i) \\ &= \mathbf{x}^i - \alpha^i (\Phi^T \Phi \mathbf{x}^i - \Phi^T \mathbf{y} + \tau(\mathbf{x}^i - \mathbf{v}^i)) \end{aligned} \quad (10)$$

where  $g(\cdot)$  is the gradient function of  $\mathbf{x}$ , and  $\alpha^i$  is the length of the gradient descent step.

### Anderson acceleration

Many efforts have been devoted to developing acceleration methods for the gradient descent algorithm [39]. For example, the widely used Momentum acceleration method takes into account the previous gradients in the update step at each iteration [40]; Anderson acceleration uses the residuals of previous  $m$  iterations to adjust the current iteration point [7]. We claim that acceleration methods not only speed

up convergence but can also inform the *network's architecture*. Specifically, we use the general acceleration form:

$$\mathbf{x}^{i+1} = \mathbf{x}^i - \sum_{j=1}^{m'} w_j^i \mathbf{d}^{i-j} - \alpha_i g(\mathbf{x}^i - \sum_{j=1}^{m'} w_j^i \mathbf{d}^{i-j}), \quad (11)$$

where  $\mathbf{d}^{i-j}$  is the descent direction in the  $j$ -th iteration prior to iteration  $i$ , and  $w_j^i$  is the weight of the descent direction in iteration  $i$ . We choose  $m' = \min(m, i)$  to ensure that  $i - m'$  is a non-negative integer in the early layers.

Note that the form of Eq. 11 is exactly that of Anderson acceleration [7], [41], except that the parameters of Anderson acceleration are manually estimated while ours are learned from the network. Specifically, when  $m = 1$ , our acceleration becomes Nesterov's accelerated gradient method [42].

Since the norm of the residual in each iteration can be absorbed by its weights  $w_j^i$ , without loss of generality, we directly let

$$\mathbf{d}^i = \mathbf{x}^i - \mathbf{x}^{i-1}. \quad (12)$$

Combining Eq. 11 and the definition of  $g(\cdot)$  in Eq. 10, the update step of  $\mathbf{x}$  can be rewritten as:

$$\mathbf{x}^{i+1} = [(1-\beta^i)\mathbf{I} - \alpha^i \Phi^T \Phi] (\mathbf{x}^i - \sum_{j=1}^{m'} w_j^i \mathbf{d}^{i-j}) + \alpha^i \Phi^T \mathbf{y} + \beta^i \mathbf{v}^i, \quad (13)$$

where  $\alpha^i \tau$  is denoted as  $\beta^i$ . We show the detailed operations and connections in and between stages in Fig.4 (a). Compared to general unrolling networks, the skip connections between stages in our model make the network denser and more compact, and transform it from a Resnet to a Densenet.

The denoising network we used to solve the second sub-problem in Eq 8 consists of two cascaded residual blocks. The architecture of the denoising network is as shown in Fig. 4 (b). The number of used residual blocks is chosen empirically. Previous work [43] gave some convergence analysis and also showed that two residual blocks provide the best results for learning the proximal operator. Note that we can also apply non-local attention [44] and a multi-scale architecture [45], [46]. But to ensure the decoding network has a limited parameter count to prevent overfitting, each residual block in the denoising network contains only five convolutional layers, and all layers generate feature maps with  $3 \times 3$  kernels.



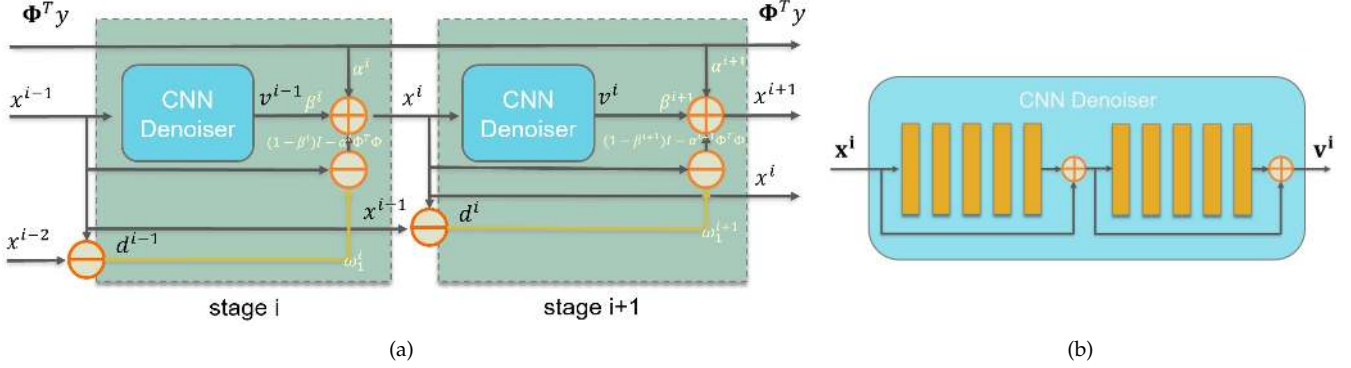


Fig. 4. (a) Illustration of two stages in our network. (here we show the case  $m = 1$ ) (b) The architecture of our denoising network.

**Algorithm 1** Accelerated subframe reconstruction

**Input:** Sensing matrix  $\Phi$ , captured image  $\mathbf{y}$ , number  $m$   
**Output:** Reconstructed subframes  $\mathbf{x}$

- 1: Initialize  $\mathbf{x}^0 = \Phi^T \mathbf{y}$ ,  $\mathbf{x}^{-1} = \mathbf{x}^0$  ( $i = 1, \dots, m$ ),  $\mathbf{d}^0 = 0$
- 2: **for**  $i = 1, 2, \dots, K$  **do**
- 3:    $\mathbf{v}^{i-1} = D(\mathbf{x}^{i-1})$
- 4:    $m' = \min(m, i)$
- 5:    $\mathbf{z}^i = \mathbf{x}^{i-1} - \sum_{j=1}^{m'} w_j^i \mathbf{d}^{i-j}$
- 6:    $\mathbf{x}^i = [(1 - \beta^i) \mathbf{I} - \alpha^i \Phi^T \Phi] \mathbf{z}^i + \alpha^i \Phi^T \mathbf{y} + \beta^i \mathbf{v}^{i-1}$
- 7:    $\mathbf{d}^i = \mathbf{x}^i - \mathbf{x}^{i-1}$
- 8: **end for**

**3.3 Training**

We constructed an end-end network by unrolling the algorithm shown in Algorithm 1. The proposed model mainly consists of a mask layer and a K-stage reconstruction network using convolutional layers. The input subframes  $\mathbf{x}$  are encoded using a trainable mask layer  $\phi$ . We multiply the transpose of the mask  $\Phi^T$  and the captured images  $\mathbf{y}$  to generate an initial guess  $\mathbf{x}^0 = \Phi^T \mathbf{y}$ . We then feed the initial image into the reconstruction. All layers use ReLU as their activation function, except the output layer, which uses a sigmoid. We choose the mean square error (MSE) as our loss function, expressed as

$$\mathcal{L}(\phi, w; \alpha; \beta; \theta) = \frac{1}{k} \sum_{i=1}^k \|f(\mathbf{x}; \phi; w; \alpha; \beta; \theta) - \mathbf{x}\|^2, \quad (14)$$

where  $k$  is the number of the training samples,  $\theta$  are the denoising network weights,  $\phi$  are the mask layer weights, and  $(w; \alpha; \beta)$  are the optimization parameters. We trained the proposed network to learn these parameters simultaneously. The parameters of each stage are set to be different, and the  $\alpha$  are set to be channel-wise.

The model was trained on an Intel Xeon E5 workstation with an NVIDIA GeForce RTX 2080 Ti GPU and 512 GB main memory. Our network is implemented using Keras 2.2.5 and trained using the Adam optimizer [47]. The initial learning rate is set to  $10^{-4}$  and decayed by a factor of 10 at the 20th iteration. We train the model for 80 iterations with a batch size of 1, which takes about two days to complete.

**4 SIMULATIONS**

In this section, we conduct numerical simulations to show the effectiveness of our proposed network and compare our method with other state-of-the-art compressive reconstruction methods.

**Datasets and Training.** The data we used for the simulations are two popular databases: the SumMe database from <https://gyglim.github.io/me/vsum/index.html> [48] and the "Sports Videos in the Wild" database from <http://cvlab.cse.msu.edu/project-svw.html> [49]. We randomly cropped and selected 3,000 video sequences of size  $256 \times 256 \times 32$  to train our network, and selected 800 video sequences of the same size for testing.

TABLE 1  
Ablation Studies.

Methods	Noiseless		Noisy ( $\sigma = 0.01$ )	
	PSNR	SSIM	PSNR	SSIM
Unopt [26]	30.68	0.896	28.52	0.861
Opt	32.35	0.921	30.52	0.897
Opt + SC ( $m=1$ )	33.18	0.930	31.24	0.905
Opt + SC ( $m=2$ )	33.30	0.932	31.43	0.908
Opt + SC ( $m=3$ )	33.32	0.932	31.46	0.909

**Ablation studies.** To clearly understand the effect of each component as well as choosing an appropriate  $m$  in our end-to-end network, we carried out five ablation simulations. We present our observations and quantitative results in Table 1. For all the simulations in the ablation study, we used the architecture shown in Fig. 4 with 39 stages for frame reconstruction, and calculated the average PSNR and SSIM of the reconstructed results in the presence and absence of noise. The baseline for comparison is model Unopt, a multistage network without mask optimization and dense skip connections, which is the same network architecture as in previous work [26]. Compared to this baseline, our method leads to a significant improvement in reconstruction quality as well as to a reduction of the number of training epochs needed for the same accuracy.

*Optimized vs. fixed mask:* For the Unopt model, we used a randomly shifted Bernoulli binary masks as shown in Fig. 5(a) while in other Opt models we used optimized

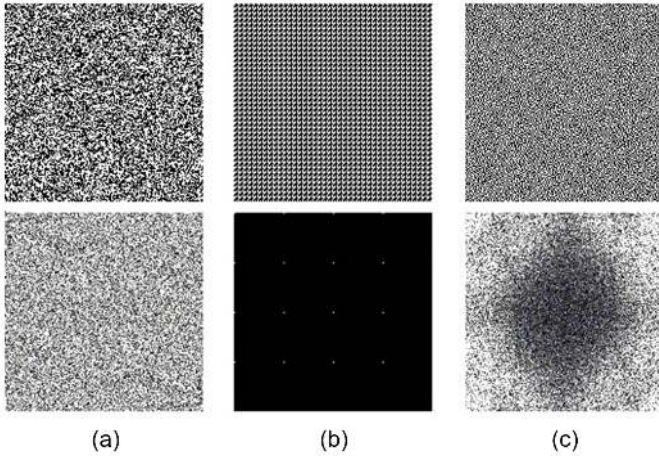


Fig. 5. The comparison of the first binary pattern(upper) and their spectrum distribution(bottom) of the three used masks sequences. (a) Bernoulli pattern used in [50] and [8]. (b) Optimized repeated pattern of [18]. (c) Our optimized pattern. Note that the patterns were cropped into  $160 \times 160$  for visualization.

364 masks as shown in Fig. 5(c). PSNRs can be improved by  
 365 nearly 1dB when replacing the random masks by the opt-  
 366 imized masks. It is worth noting that the loss of Unopt  
 367 is relatively low in the initial few epochs since random  
 368 Bernoulli masks are suitable for compressive reconstruction  
 369 [51]. However, Opt models catch up with and surpass  
 370 the Unopt model as the number of epochs increases, as  
 371 shown in Fig. 7. The results indicate that our network has  
 372 learned more efficient masks after several epochs of training.

373 *Skip connections (SC) vs. no skip-connections:* We tested  
 374 the effect of skip connections in our network. It is obvious that  
 375 skip connections can enhance reconstruction quality and  
 376 accelerate the convergence of training loss. The PSNRs are  
 377 improved by nearly 1dB when three skip connections for  
 378 a single stage ( $m = 3$ ) are applied. However, denser skip  
 379 connections require more memory, so we need to choose  
 380 an appropriate  $m$  for the best trade-off between memory  
 381 consumption and reconstruction accuracy. As shown in  
 382 Table 1, the model with  $m = 3$  outperforms the one with  
 383  $m = 2$ , but only by a small margin in both PSNR and SSIM.  
 384 Therefore, we choose  $m = 2$  as an empirical setting for our  
 385 reconstruction network.

386 **Comparison methods.** We compared the proposed  
 387 method with two representative DNN-based methods:  
 388 DeepMask [18] and Deep Tensor ADMM-Net (DTAN) [50];  
 389 and two state-of-the-art traditional methods: GAP-TV [8]  
 390 and GMM [10]. Following previous literature, we used  
 391 masks to modulated every eighth consecutive frame. Thus  
 392 we reconstructed 32 subframes from 4 measurements in the  
 393 simulations. To be specific, DeepMask is the only existing  
 394 method which can jointly optimize masks and reconstruction  
 395 method; it learns  $4 \times 4 \times 8$  repeated masks for encoding  
 396 and reconstructs frames via a fully-connected network. The  
 397 other three methods use a  $256 \times 256 \times 8$  shifting Bernoulli  
 398 binary masks. The masks of different methods and their  
 399 frequency spectra are shown in Fig.5. It can be observed  
 400 that our masks perform as a ‘high-pass filter’ that blocks  
 401 low-frequency spatial content.

402 **Quantitative results.** The PSNR and SSIM results of

TABLE 2  
The comparison of reconstruction quality of the five methods with  $T=8$  subframes.

Methods	Noiseless		Noisy( $\sigma = 0.01$ )	
	PSNR	SSIM	PSNR	SSIM
GAP-TV [8] + random	29.82	0.857	27.99	0.835
GAP-TV + optimized	30.72	0.884	29.04	0.843
GMM [10] + random	27.24	0.797	27.00	0.774
GMM + optimized	27.35	0.807	27.10	0.785
DTAN [50] + random	26.08	0.803	25.12	0.799
DTAN + optimized	27.28	0.816	26.45	0.813
DeepMask [18]	31.05	0.905	29.28	0.882
Ours	<b>33.32</b>	<b>0.932</b>	<b>31.43</b>	<b>0.908</b>

TABLE 3  
The comparison of reconstruction quality of the four methods with  $T=32$  subframes.

Methods	Noiseless		Noisy( $\sigma = 0.01$ )	
	PSNR	SSIM	PSNR	SSIM
GAP-TV [8] + random	23.44	0.725	23.15	0.700
GMM [10] + random	22.19	0.589	22.16	0.583
DeepMask [18]	27.58	0.814	25.46	0.792
Ours	<b>28.01</b>	<b>0.840</b>	<b>26.15</b>	<b>0.810</b>

403 different methods with different masks are shown in Table  
 404 2. As an optimization method, GAP-TV is effective and  
 405 efficient in reconstructing subframes, but the reconstruction  
 406 quality is not competitive compared to ours due to the  
 407 used handcrafted priors. The GMM approach reconstructs  
 408 frames patch-by-patch, and also cannot produce competi-  
 409 tive results. To our surprise, DTAN performs worst among  
 410 these methods, although it works well on its ‘NBA’ dataset.  
 411 This might be because the non-local low-rank prior fails  
 412 in reconstructing spatial high-frequency content. Due to  
 413 the joint design of masks and reconstruction, the average  
 414 PSNR and SSIM of DeepMask exceed 31dB and 0.9, re-  
 415 spectively. However, we finds serious structured artifacts in  
 416 the reconstructed images of DeepMask (see Fig.6) caused  
 417 by the use of repeated masks. Our method outperforms  
 418 state-of-art methods by more than 2.2dB in PSNR and more  
 419 than 0.03 in SSIM. This is further confirmed by visual  
 420 comparison of the reconstructed images in Fig. 6, where  
 421 we show ground truth and the reconstructed results of  
 422 four frames. Our method generates much more visually  
 423 pleasant images with more accurate detail information. We  
 424 also compared our method with GAP-TV, GMM, and Deep-  
 425 Mask with  $T=32$  subframes. In this simulation, 64 frames  
 426 are reconstructed from two encoded images. The results are  
 427 shown in Table 3. Compression ratios of 1:32 are very  
 428 challenging for compressive sensing algorithms in general, so  
 429 the results are worse than for 8 subframes, however our  
 430  
 431

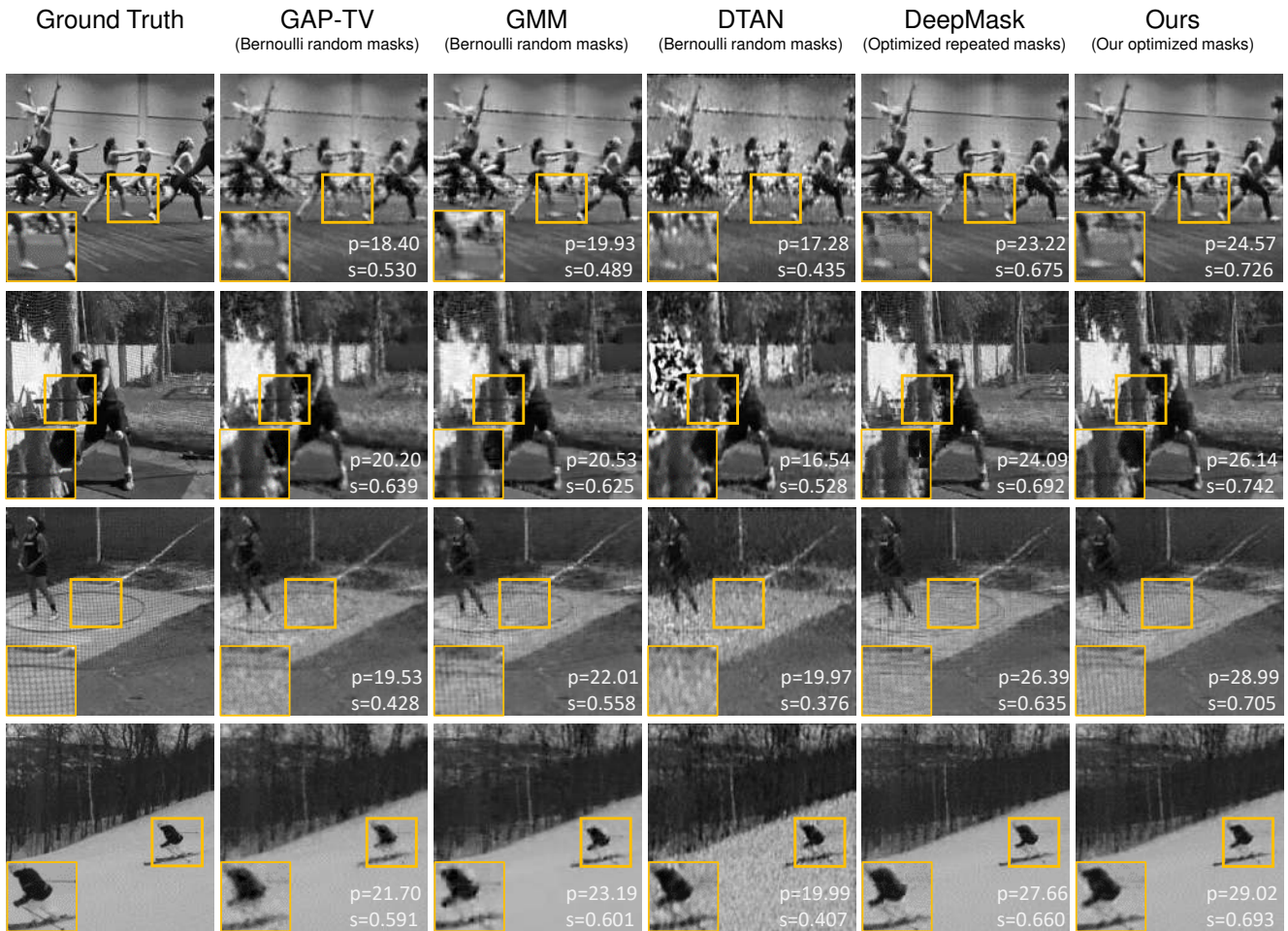


Fig. 6. The comparison of reconstructed frames and the statistics on the PSNR and SSIM. From top to bottom: ground truth;reconstructed results of GAP-TV, GMM, Deep Tensor Admm-net, DeepMask, and ours.

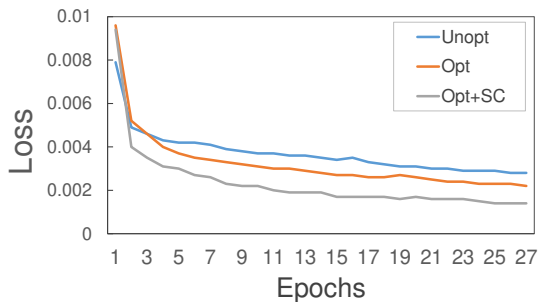


Fig. 7. Training loss vs number of epochs on the neural network models in ablation study.

approach still dominates the comparison methods.

**Mask evaluation.** We also evaluated our optimized mask by comparing it with random masks using the same reconstruction method. Since GAP-TV is a model-based optimization method which does not memorize data, we reconstructed frames using GAP-TV with random masks and our proposed masks respectively to present the behavior of the two masks. Fig. 8 shows the reconstructed results. The frames reconstructed from the image encoded by our masks are significantly better than those by random masks,

especially around the edges. We also observed the improvement brought by the optimized masks using other existing methods (e.g. [10] and [50]). Note that the improvement for the GMM method is not significant since it reconstructs frames patch-by-patch while our masks are optimized as a whole.

## 5 REAL EXPERIMENTS

Previous work on mask-based video compressive sensing uses either a static mask that is shifted over time, or a setup with some form of spatial light modulator, such as a DMD or LCOS, which can be controlled with high temporal resolution. However, the drawback of these methods is that they are difficult to align and rather bulky due to the need for re-imaging optics [52].

Fortunately, recent developments in image sensor technology allow us to directly implement the CS mask on the sensor itself. Specifically, there are now several prototypes of image sensors with per-pixel programmable exposure control [5], [37]. In this paper, we use the Coded two-Bucket (C2B) camera from Wei *et al.* [5]. In this camera, each pixel has two charge-collection sites (i.e. two buckets). The exposure control signal for each pixel can select which of the two buckets integrates incident light at any given point in

432  
433  
434  
435  
436  
437  
438  
439  
440  
441

442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464



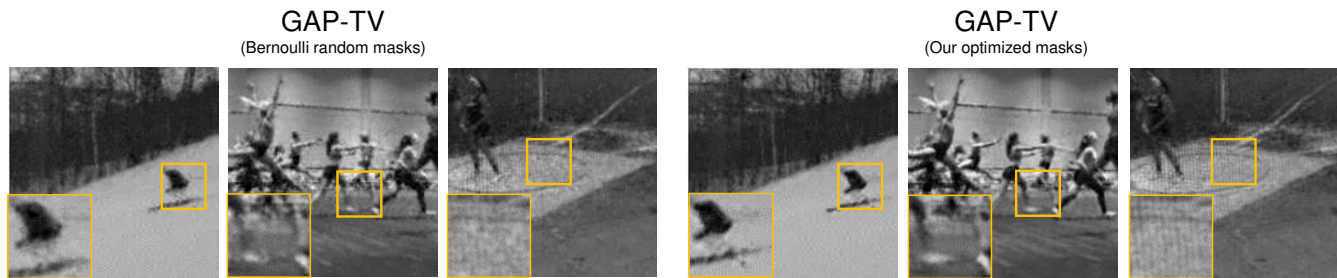


Fig. 8. The comparison of reconstructed results using GAP-TV method with different encoding masks.

465 time. The major advantage of this design is that it makes  
 466 use of all incident photons and simultaneously encodes  
 467 subframes with a pair of complementary masks. Using this  
 468 camera, subframes are reconstructed from the pair of cap-  
 469 tured complementary images. The spatial resolution of the  
 470 camera is  $312 \times 320$ , and the frame rate can reach 30 frames  
 471 per second with over 100 different masks per frame. In our  
 472 experiments we use only up to 32 masks per frame since a  
 473 compression ration of 1:32 is already extremely challenging  
 474 for all compressive sensing approaches.

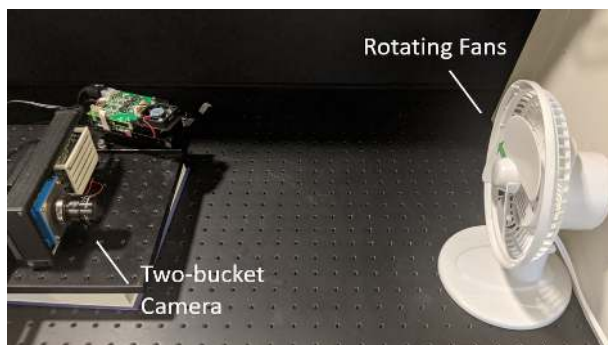


Fig. 9. The setup of our experiments.

475 We captured several dynamic scenes using the camera to  
 476 compare the reconstruction quality of four different meth-  
 477 ods: GAP-TV [8], GMM [10], DeepMask [18], and ours. The  
 478 setup for our experiments is shown in Fig. 9. Unlike the  
 479 simulation, here, the number of subframes we used is 32 to  
 480 explore the limits of the four methods; thus, a high-frame-  
 481 rate ( $32 \times 30 = 960$ ) imaging can be achieved. In the exper-  
 482 iment, the first two methods used  $312 \times 320 \times 32$  random  
 483 masks, DeepMask used optimized repeated  $4 \times 4 \times 32$  masks,  
 484 and our method used  $312 \times 320 \times 32$  optimized masks. We  
 485 reconstructed 64 subframes from two successively captured  
 486 images. Fig. 10 shows two examples of the reconstructed  
 487 results. It can be seen that the GAP-TV method created  
 488 watercolor-like artifacts due to the drawbacks of the hand-  
 489 crafted prior; GMM and DeepMask introduced significant  
 490 structured artifacts in the patch-by-patch reconstruction.  
 491 The proposed methods, on the other hand, can produce  
 492 better results with fewer artifacts, clearer contents, and  
 493 higher contrast compared with the other three methods  
 494 (please zoom in for details).

495 We also investigated the improvement brought by the  
 496 two bucket mechanism of the camera. With the two-bucket  
 497 mechanism each subframe is encoded by a pair of com-

plementary masks, so that the number of measurements  
 is doubled when compared to the one-bucket mechanism.  
 To demonstrate the improvements due to the two-bucket  
 design, we captured a fan with varying rotation speeds and  
 reconstructed 64 subframes from two one-bucket images  
 and four two-bucket images respectively. The results are  
 shown in Fig.11. It can be seen that the reconstructed results  
 from two-bucket images are significantly better than those  
 from one-bucket images. We can also observe that the ad-  
 vantages of our method over the state of the art are even  
 more compelling in real experiments than in simulation.  
 That is because our method depends on a deep image prior  
 rather than handcrafted priors and thus can better handle  
 complicated video content found in real scenes.

## 6 CONCLUSION AND FUTURE WORKS

We have presented a new end-to-end learned method and  
 prototype system for video reconstruction from mask-based  
 compressive sensing cameras. Unlike existing approaches,  
 the proposed method is suited for optimizing full-resolution  
 masks, and can reconstruct subframes efficiently. The re-  
 construction quality of the proposed method significantly  
 outperforms that of previous methods due to the utilized  
 deep image prior. We implemented a two-bucket camera for  
 high-frame-rate imaging; the frame rate can reach close to  
 1,000 frames with superior image quality compared to other  
 CS video approaches.

In addition to providing a superior solution to the  
 compressive sensing video reconstruction problem, we also  
 make a fundamental improvement to loop-unrolled neural  
 network architectures for image reconstruction problems  
 in general: we demonstrate that dense skip connections  
 can implement Anderson acceleration directly in the neural  
 network to make it compact and efficient. The proposed  
 dense network is not limited to CS problems, but can be  
 applied to solve other inverse problems directly.

We believe that the frames in the near future can be  
 predicted from previously reconstructed frames. Therefore,  
 in future work, we plan to explore more efficient frame  
 reconstruction and adaptively optimize masks in real-time  
 for even better results.

## REFERENCES

- [1] E. J. Candès, J. Romberg, and T. Tao, "Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information," *IEEE Transactions on information theory*, vol. 52, no. 2, pp. 489–509, 2006.

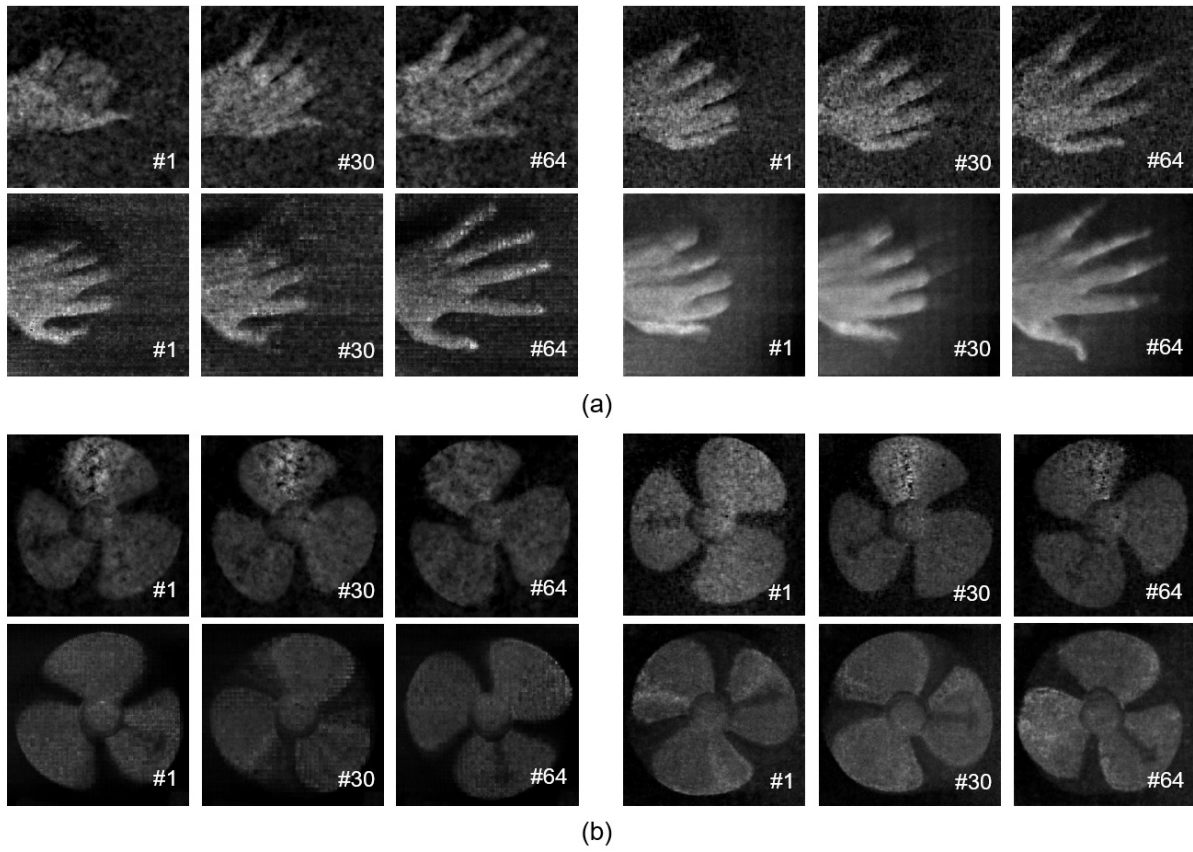


Fig. 10. The reconstructed results of (a)an opening hand and (b) a rotating fans using four methods. Left-top: GAP-TV; Right-top: GMM; Left-bottom: DeepMask; Right-bottom: our method. Here shows the 1st, 30th, and the 64th subframes reconstructed from two one-bucket images. The rotating speed of the fans is 2.5 rounds pre second. Note that the reconstructed subframes are scaled by the maximum intensity for visualization.

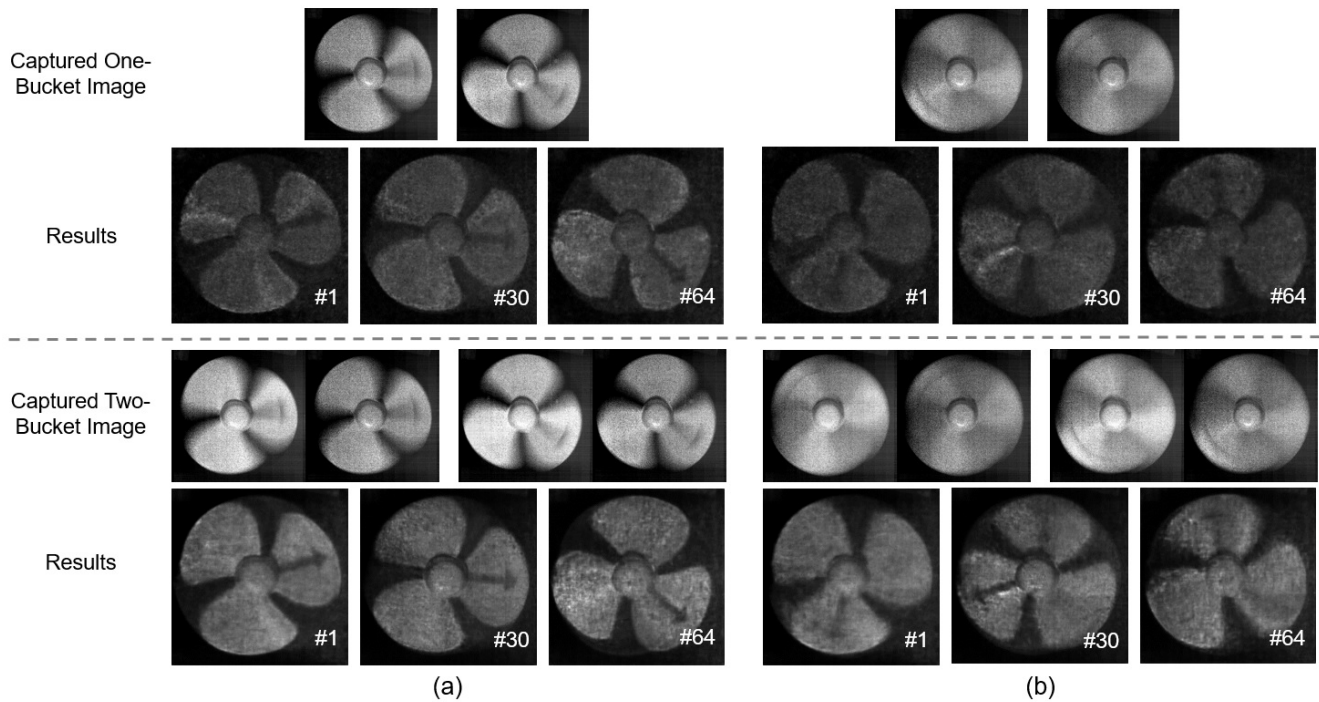


Fig. 11. The 1st, 30th, and 64th subframes of a rotating fans reconstructed from two one-bucket encoded images and four two-bucket encoded images. The fans are captured under the rotating speeds (a) 2.5 rounds and (b) 7 rounds per second. Note that our method can reconstruct clear results from the two-bucket encoded images with heavy motion-blur.

- 543 [2] M. Lustig, D. L. Donoho, J. M. Santos, and J. M. Pauly, "Com-  
544 pressed sensing mri," *IEEE signal processing magazine*, vol. 25, no. 2,  
545 p. 72, 2008.
- 546 [3] A. Wagadarikar, R. John, R. Willett, and D. Brady, "Single disperser  
547 design for coded aperture snapshot spectral imaging," *Appl. Opt.*,  
548 vol. 47, no. 10, pp. B44–B51, Apr 2008.
- 549 [4] P. Llull, X. Liao, X. Yuan, J. Yang, D. Kittle, L. Carin,  
550 G. Sapiro, and D. J. Brady, "Coded aperture compressive  
551 temporal imaging," *Optics Express*, vol. 21, no. 9, p. 10526,  
552 May 2013. [Online]. Available: [https://www.osapublishing.org/  
553 abstract.cfm?URI=oe-21-9-10526](https://www.osapublishing.org/abstract.cfm?URI=oe-21-9-10526)
- 554 [5] M. Wei, N. Sarhangnejad, Z. Xia, N. Gusev, N. Katic, R. Genov,  
555 and K. N. Kutulakos, "Coded two-bucket cameras for computer  
556 vision," in *Computer Vision ECCV 2018*, V. Ferrari, M. Hebert,  
557 C. Sminchisescu, and Y. Weiss, Eds. Springer International  
558 Publishing, 2018, vol. 11207, pp. 55–73. [Online]. Available:  
559 [http://link.springer.com/10.1007/978-3-030-01219-9\\_4](http://link.springer.com/10.1007/978-3-030-01219-9_4)
- 560 [6] G. Zhang, S. Jiao, X. Xu, and L. Wang, "Compressed sensing  
561 and reconstruction with bernoulli matrices," in *The 2010 IEEE  
562 International Conference on Information and Automation*, June 2010,  
563 pp. 455–460.
- 564 [7] H. F. Walker and P. Ni, "Anderson acceleration for fixed-point  
565 iterations," *SIAM Journal on Numerical Analysis*, vol. 49, no. 4, pp.  
566 1715–1735, 2011.
- 567 [8] X. Yuan, "Generalized alternating projection based total variation  
568 minimization for compressive sensing," in *2016 IEEE International  
569 Conference on Image Processing (ICIP)*. IEEE, 2016, pp. 2539–2543.
- 570 [9] D. Reddy, A. Veeraraghavan, and R. Chellappa, "P2c2: Pro-  
571 grammable pixel compressive camera for high speed imaging,"  
572 in *CVPR 2011*, June 2011, pp. 329–336.
- 573 [10] J. Yang, X. Yuan, X. Liao, P. Llull, D. J. Brady, G. Sapiro, and  
574 L. Carin, "Video compressive sensing using gaussian mixture  
575 models," *IEEE Transactions on Image Processing*, vol. 23, no. 11, pp.  
576 4863–4878, Nov 2014.
- 577 [11] Y. Hitomi, J. Gu, M. Gupta, T. Mitsunaga, and S. K. Nayar, "Video  
578 from a single coded exposure photograph using a learned over-  
579 complete dictionary," in *2011 International Conference on Computer  
580 Vision*, Nov 2011, pp. 287–294.
- 581 [12] D. Liu, J. Gu, Y. Hitomi, M. Gupta, T. Mitsunaga, and S. K. Nayar,  
582 "Efficient space-time sampling with pixel-wise coded exposure  
583 for high-speed imaging," *IEEE Transactions on Pattern Analysis and  
584 Machine Intelligence*, vol. 36, no. 2, pp. 248–260, Feb 2014.
- 585 [13] Y. Liu, X. Yuan, J. Suo, D. Brady, and Q. Dai, "Rank minimization  
586 for snapshot compressive imaging," *IEEE transactions on pattern  
587 analysis and machine intelligence*, 2018.
- 588 [14] R. Obermeier and J. A. Martinez-Lorenzo, "Sensing matrix design  
589 via mutual coherence minimization for electromagnetic compres-  
590 sive imaging applications," *IEEE Transactions on Computational  
591 Imaging*, vol. 3, no. 2, pp. 217–229, June 2017.
- 592 [15] S. Lohit, K. Kulkarni, R. Kerviche, P. Turaga, and A. Ashok,  
593 "Convolutional neural networks for noniterative reconstruction of  
594 compressively sensed images," *IEEE Transactions on Computational  
595 Imaging*, vol. 4, no. 3, pp. 326–340, Sep. 2018.
- 596 [16] H. Yao, F. Dai, S. Zhang, Y. Zhang, Q. Tian, and C. Xu, "Dr2-  
597 net: Deep residual reconstruction network for image compressive  
598 sensing," *Neurocomputing*, vol. 359, pp. 483 – 493, 2019.
- 599 [17] K. Kulkarni, S. Lohit, P. Turaga, R. Kerviche, and A. Ashok,  
600 "ReconNet: Non-Iterative Reconstruction of Images from  
601 Compressively Sensed Measurements," in *2016 IEEE Conference  
602 on Computer Vision and Pattern Recognition (CVPR)*. Las Vegas,  
603 NV, USA: IEEE, Jun. 2016, pp. 449–458. [Online]. Available:  
604 <http://ieeexplore.ieee.org/document/7780424/>
- 605 [18] M. Iliadis, L. Spinoulas, and A. Katsaggelos, "Deep fully-  
606 connected networks for video compressive sensing," *Digital Signal  
607 Processing: A Review Journal*, vol. 72, pp. 9–18, 1 2018.
- 608 [19] M. Yoshida, A. Torii, M. Okutomi, K. Endo, Y. Sugiyama, R.-i.  
609 Taniguchi, and H. Nagahara, "Joint optimization for compressive  
610 video sensing and reconstruction under hardware constraints," in  
611 *Computer Vision – ECCV 2018*, V. Ferrari, M. Hebert, C. Sminchis-  
612 escu, and Y. Weiss, Eds. Cham: Springer International Publishing,  
613 2018, pp. 649–663.
- 614 [20] J. Ma, X.-Y. Liu, Z. Shou, and X. Yuan, "Deep Tensor ADMM-  
615 Net for Snapshot Compressive Imaging," in *IEEE Conference on  
616 Computer Vision and Pattern Recognition (CVPR)*, 2019, p. 10.
- 617 [21] K. Gregor and Y. LeCun, "Learning fast approximations of sparse  
618 coding," in *Proceedings of the 27th International Conference on In-  
619 ternational Conference on Machine Learning*. Omnipress, 2010, pp.  
620 399–406.
- 621 [22] J. Sun, H. Li, Z. Xu *et al.*, "Deep admn-net for compressive sensing  
622 mri," in *Advances in neural information processing systems*, 2016, pp.  
623 10–18.
- 624 [23] C. Metzler, A. Mousavi, and R. Baraniuk, "Learned d-amp: Prin-  
625 ciple neural network based compressive image recovery," in  
626 *Advances in Neural Information Processing Systems*, 2017, pp. 1772–  
627 1783.
- 628 [24] K. Zhang, W. Zuo, S. Gu, and L. Zhang, "Learning deep cnn  
629 denoiser prior for image restoration," in *Proceedings of the IEEE  
630 conference on computer vision and pattern recognition*, 2017, pp. 3929–  
631 3938.
- 632 [25] J. Zhang and B. Ghanem, "Ista-net: Interpretable optimization-  
633 inspired deep network for image compressive sensing," in *Pro-  
634 ceedings of the IEEE Conference on Computer Vision and Pattern  
635 Recognition*, 2018, pp. 1828–1837.
- 636 [26] W. Dong, P. Wang, W. Yin, G. Shi, F. Wu, and X. Lu, "Denoising  
637 prior driven deep neural network for image restoration," *IEEE  
638 transactions on pattern analysis and machine intelligence*, vol. 41,  
639 no. 10, pp. 2305–2318, 2018.
- 640 [27] M. Kellman, E. Bostan, N. Repina, and L. Waller, "Physics-based  
641 learned design: Optimized coded-illumination for quantitative  
642 phase imaging," *IEEE Transactions on Computational Imaging*, 2019.
- 643 [28] R. Raskar, A. Agrawal, and J. Tumblin, "Coded exposure pho-  
644 tography: motion deblurring using fluttered shutter," in *ACM  
645 transactions on graphics (TOG)*, vol. 25, no. 3. ACM, 2006, pp.  
646 795–804.
- 647 [29] A. Veeraraghavan, D. Reddy, and R. Raskar, "Coded strobing pho-  
648 tography: Compressive sensing of high speed periodic videos,"  
649 *IEEE Transactions on Pattern Analysis and Machine Intelligence*,  
650 vol. 33, no. 4, pp. 671–686, 2010.
- 651 [30] D. Liu, J. Gu, Y. Hitomi, M. Gupta, T. Mitsunaga, and S. K. Nayar,  
652 "Efficient Space-Time Sampling with Pixel-Wise Coded Exposure  
653 for High-Speed Imaging," *IEEE Transactions on Pattern Analysis and  
654 Machine Intelligence*, vol. 36, no. 2, pp. 248–260, Feb. 2014.
- 655 [31] G. Bub, M. Tecza, M. Helmes, P. Lee, and P. Kohl, "Temporal  
656 pixel multiplexing for simultaneous high-speed, high-resolution  
657 imaging," *Nature methods*, vol. 7, no. 3, p. 209, 2010.
- 658 [32] Y. Luo, D. Ho, and S. Mirabbasi, "Exposure-programmable cmos  
659 pixel with selective charge storage and code memory for computa-  
660 tional imaging," *IEEE Transactions on Circuits and Systems I: Regular  
661 Papers*, vol. 65, no. 5, pp. 1555–1566, 2017.
- 662 [33] J. Zhang, T. Xiong, T. Tran, S. Chin, and R. Etienne-Cummings,  
663 "Compact all-cmos spatiotemporal compressive sensing video  
664 camera with pixel-wise coded exposure," *Optics express*, vol. 24,  
665 no. 8, pp. 9013–9024, 2016.
- 666 [34] J. P. Newman, X. Wang, C. S. Thakur, J. Rattray, R. Etienne-  
667 Cummings, M. A. Wilson *et al.*, "A closed-loop all-electronic pixel-  
668 wise adaptive imaging system for high dynamic range video,"  
669 *arXiv preprint arXiv:1906.10045*, 2019.
- 670 [35] T. Sonoda, H. Nagahara, K. Endo, Y. Sugiyama, and R.-i.  
671 Taniguchi, "High-speed imaging using cmos image sensor with  
672 quasi pixel-wise exposure," in *2016 IEEE International Conference  
673 on Computational Photography (ICCP)*. IEEE, 2016, pp. 1–11.
- 674 [36] M. Yoshida, A. Torii, M. Okutomi, K. Endo, Y. Sugiyama, R.-i.  
675 Taniguchi, and H. Nagahara, "Joint optimization for compressive  
676 video sensing and reconstruction under hardware constraints," in  
677 *Proceedings of the European Conference on Computer Vision (ECCV)*,  
678 2018, pp. 634–649.
- 679 [37] N. Sarhangnejad, N. Katic, Z. Xia, M. Wei, N. Gusev, G. Dutta,  
680 R. Gulve, H. Haim, M. M. Garcia, D. Stoppa *et al.*, "5.5 dual-tap  
681 pipelined-code-memory coded-exposure-pixel cmos image sensor  
682 for multi-exposure single-frame computational imaging," in *2019  
683 IEEE International Solid-State Circuits Conference-ISSCC*. IEEE,  
684 2019, pp. 102–104.
- 685 [38] M. Courbariaux, Y. Bengio, and J.-P. David, "Binaryconnect: Train-  
686 ing deep neural networks with binary weights during propaga-  
687 tions," in *Advances in Neural Information Processing Systems 28*,  
688 C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett,  
689 Eds. Curran Associates, Inc., 2015, pp. 3123–3131.
- 690 [39] S. Ruder, "An overview of gradient descent optimization algo-  
691 rithms," *arXiv preprint arXiv:1609.04747*, 2016.
- 692 [40] N. Qian, "On the momentum term in gradient descent learning  
693 algorithms," *Neural Networks*, vol. 12, no. 1, pp. 145 – 151, 1999.
- 694 [41] J. Zhang, Y. Peng, W. Ouyang, and B. Deng, "Accelerating admn

- 695 for efficient simulation and optimization," in *Siggraph Asia*, Nov.  
696 2019.
- 697 [42] Y. Nesterov, "A method for solving a convex programming prob-  
698 lem with convergence rate  $o(1/k^2)$ ," *Soviet Mathematics Doklady*,  
699 no. 27, pp. 372–367, 1983.
- 700 [43] M. Mardani, Q. Sun, D. Donoho, V. Pappyan, H. Monajemi,  
701 S. Vasanawala, and J. Pauly, "Neural proximal gradient descent for  
702 compressive imaging," in *Advances in Neural Information Processing*  
703 *Systems*, 2018, pp. 9573–9583.
- 704 [44] Y. Zhang, K. Li, K. Li, B. Zhong, and Y. Fu, "Residual non-  
705 local attention networks for image restoration," *arXiv preprint*  
706 *arXiv:1903.10082*, 2019.
- 707 [45] S. Nah, T. Hyun Kim, and K. Mu Lee, "Deep multi-scale con-  
708 volutional neural network for dynamic scene deblurring," in  
709 *Proceedings of the IEEE Conference on Computer Vision and Pattern*  
710 *Recognition*, 2017, pp. 3883–3891.
- 711 [46] T. Rott Shaham, T. Dekel, and T. Michaeli, "Singan: Learning a  
712 generative model from a single natural image," in *Computer Vision*  
713 *(ICCV), IEEE International Conference on*, 2019.
- 714 [47] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimiza-  
715 tion," *arXiv preprint arXiv:1412.6980*, 2014.
- 716 [48] M. Gygli, H. Grabner, H. Riemenschneider, and L. Van Gool,  
717 "Creating summaries from user videos," in *ECCV*, 2014.
- 718 [49] S. M. Safdarnejad, X. Liu, L. Udpa, B. Andrus, J. Wood, and  
719 D. Craven, "Sports videos in the wild (svw): A video dataset for  
720 sports analysis," in *Proc. International Conference on Automatic Face*  
721 *and Gesture Recognition*, Ljubljana, Slovenia, May 2015.
- 722 [50] J. Ma, X.-Y. Liu, Z. Shou, and X. Yuan, "Deep tensor admm-net  
723 for snapshot compressive imaging," in *Proceedings of the IEEE*  
724 *International Conference on Computer Vision*, 2019, pp. 10 223–10 232.
- 725 [51] G. Zhang, S. Jiao, X. Xu, and L. Wang, "Compressed sensing  
726 and reconstruction with bernoulli matrices," in *The 2010 IEEE*  
727 *International Conference on Information and Automation*. IEEE, 2010,  
728 pp. 455–460.
- 729 [52] Q. Sun, X. Dun, Y. Peng, and W. Heidrich, "Depth and transient  
730 imaging with compressive spad array cameras," in *Proceedings*  
731 *of the IEEE Conference on Computer Vision and Pattern Recognition*,  
732 2018, pp. 273–282.