

DeepBinaryMask: Learning a Binary Mask for Video Compressive Sensing

Michael Iliadis, *Member, IEEE*, Leonidas Spinoulas, *Member, IEEE*, and Aggelos K. Katsaggelos, *Fellow, IEEE*

Abstract—In this paper, we propose a novel encoder-decoder neural network model referred to as DeepBinaryMask for video compressive sensing. In video compressive sensing one frame is acquired using a set of coded masks (sensing matrix) from which a number of video frames, equal to the number of coded masks, is reconstructed. The proposed framework is an end-to-end model where the sensing matrix is trained along with the video reconstruction. The encoder maps a video block to compressive measurements by learning the binary elements of the sensing matrix. The decoder is trained to map the measurements from a video patch back to a video block via several hidden layers of a Multi-Layer Perceptron network. The predicted video blocks are stacked together to recover the unknown video sequence. The reconstruction performance is found to improve when using the trained sensing mask from the network as compared to other mask designs such as random, across a wide variety of compressive sensing reconstruction algorithms. Finally, our analysis and discussion offers insights into understanding the characteristics of the trained mask designs that lead to the improved reconstruction quality.

Index Terms—Deep Learning, Compressive Sensing, Mask Optimization, Binary Mask, Video Reconstruction.

I. INTRODUCTION

In signal processing, Compressive Sensing (CS) is a popular problem which has been incorporated in various applications [1], [2]. In principle, CS theory suggests that a signal can be perfectly reconstructed using a small number of random incoherent linear projections by finding solutions to underdetermined linear systems. The underdetermined linear system in CS is defined by,

$$\mathbf{y} = \Phi \mathbf{x}, \quad (1)$$

where Φ is the $M_f \times N_f$ measurement or sensing matrix with $M_f \ll N_f$. We denote the vectorized versions of the unknown signal and compressive measurements as $\mathbf{x} : N_f \times 1$ and $\mathbf{y} : M_f \times 1$, respectively. Thus, having more unknowns than equations, to guarantee a single solution in system (1) sparsity on the signal is enforced. Many signals, such as natural images, are sparse in well-known bases (e.g., Wavelet). Therefore, most reconstruction approaches employ a regularization term $F(\cdot)$ which promotes sparsity of the unknown signal \mathbf{x} on some chosen transform domain. Thus, the following minimization problem is sought after,

$$\hat{\mathbf{a}} = \underset{\mathbf{a}}{\operatorname{argmin}} F(\mathbf{a}) \quad \text{s.t.} \quad \mathbf{y} = \Phi D \mathbf{a}, \quad (2)$$

M. Iliadis, L. Spinoulas and A. K. Katsaggelos are with the Department of Electrical Engineering and Computer Science, Northwestern University, Evanston, IL 60208-3118 USA (e-mail: miliad@northwestern.edu; leonisp@u.northwestern.edu; agk@eecs.northwestern.edu).

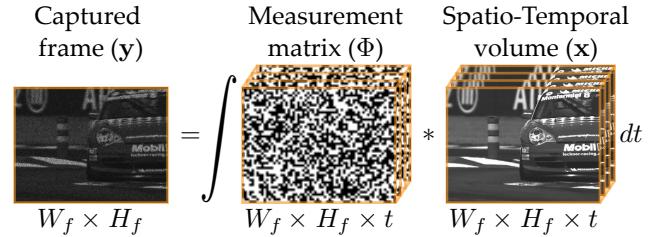


Fig. 1. Temporal compressive sensing measurement model.

where D is a chosen sparse representation transform resulting in a sparse \mathbf{a} , such that $\mathbf{x} = D\mathbf{a}$. For example, in the case $F = \|\mathbf{a}\|_0$, the problem in Eq. (2) is translated to an ℓ_0 minimization problem, which can be solved with standard numerical methods such as Orthogonal Matching Pursuit (OMP) and Basis Pursuit (BP).

Multiple algorithms have been proposed for reconstructing still images using CS by solving the problem in (2). The problem of video compressive sensing (VCS) refers to the recovery of an unknown spatio-temporal volume from the limited compressive measurements. There are two different approaches in VCS, namely spatial and temporal. Spatial VCS architectures perform spatial multiplexing per measurement based on the well-known single-pixel-camera [3] and enable video recovery by expediting the capturing process [4], [5], [6]. In temporal VCS, multiplexing occurs along the time dimension. Figure 1 demonstrates this process, where a spatio-temporal signal of size $W_f \times H_f \times t = N_f$ is modulated by t binary random masks during the exposure time of a single capture and produces a coded frame of size $W_f \times H_f = M_f$. The acquisition model in (1) applies to the temporal VCS case as well. However, the construction of Φ is different in this case. In particular, it is sparse and is given by,

$$\Phi = [\operatorname{diag}(\phi_1), \dots, \operatorname{diag}(\phi_t)] : M_f \times N_f, \quad (3)$$

where each vectorized sampling mask is expressed as ϕ_1, \dots, ϕ_t and $\operatorname{diag}(\cdot)$ creates a diagonal matrix from its vector argument. It is noted here that the spatio-temporal volume is lexicographically ordered into the vector \mathbf{x} by considering first the spatial and then the temporal dimensions.

Performance guarantees for sparse reconstruction methods, i.e., OMP, indicate that matrix Φ must be an incoherent unit norm tight frame [7]. Incoherence is a property that characterizes the degree of similarity between the columns of Φ (or ΦD). Therefore, the choice of matrix Φ is crucial for the reconstructed image and video quality irrespectively of the choice of $F(\cdot)$. For signals that can be represented sparsely

in some basis, various popular matrices in the literature are known to perform particularly well (e.g., Gaussian). However, in VCS the design of Φ as part of the acquisition hardware (e.g., camera) introduces certain limitations. For practical implementations, binary random matrices (e.g., Bernoulli) are better suited while they perform favorably to Gaussian random matrices [8].

The problem of optimizing the Φ matrix has been analyzed by several researchers [9], [10], [7], [11]. Unfortunately, optimization approaches typically rely on minimizing the coherence between the sampling matrix Φ and the sparsifying basis (ΦD), which mostly applies to spatial compressive sensing where dense matrices are used. Instead, the masks used for temporal VCS systems, as the one described herein, result in a sparse binary matrix with entries across diagonals, as presented by Eq. (3), and therefore existing results are not applicable.

In this work, we optimize the sensing matrix Φ for temporal VCS and transform it into a form that is more suitable for reconstruction using deep neural networks. The proposed neural network architecture, which is referred to as *DeepBinaryMask*, consists of two components that act as a pair of an encoder and a decoder. The encoder maps a video block to compressive measurements by learning binary weights (which correspond to the entries on the diagonals of the measurement matrix). The decoder maps the measurements back to a video block, via several hidden layers of a Multi-Layer Perceptron (MLP) network, utilizing real-valued weights. Both networks are trained jointly. We show that the mask trained from data using neural networks provides significantly improved recovery performance as compared to a non-trained sensing mask.

A. Contributions

- **Learning binary weights and reconstruction simultaneously:** Since existing approaches of optimizing the Φ matrix for *spatial* CS are not applicable for *temporal* VCS, we consider using deep learning-based recovery to train the Φ matrix and optimize mask parameters via back-propagation [12]. On this front, we propose a novel encoder-decoder neural network for temporal VCS in which the encoder learns binary weights that form the sensing mask and the decoder learns to reconstruct the video sequence given the encoded measurements. Our learning approach is performed on 3D video blocks.
- **Learning a general mask:** We show that the reconstruction performance is improved when using the optimized trained mask over a random one. Performance improvements are reported not only when the reconstruction method is the neural network decoder but also when other popular reconstruction methods are employed (e.g., based on ℓ_1 optimization).
- **Mask analysis:** We present a reconstruction performance analysis of the trained sensing mask/matrix for different mask initializations (e.g., initial number of nonzero elements). Furthermore, we conduct experiments using different random seeds when initializing the binary random masks to confirm performance stability.

II. MOTIVATION AND RELATED WORK

Recent advances in Deep Neural Networks (DNNs) [13] have demonstrated state-of-the-art performance in several computer vision and image processing tasks, such as image recognition [14] and object detection [15]. In this section we briefly discuss previous works in designing optimal masks for VCS and then we survey recent studies in image recovery problems using DNNs. Finally, we describe advances in DNNs utilizing binary weights, a key ingredient of our proposed method.

Designing optimal masks. Most of the previously proposed optimized mask patterns for temporal VCS rely on some heuristic constraints and trial-and-error patterns. A thresholded Gaussian matrix was employed in [16], [17] and [18] as it was assumed that it results in a sensing matrix that most closely resembles a dense Gaussian matrix. A normalized mask such that the total amount of light collected at each pixel is constrained to be constant was proposed by [19]. It was found in [20] that these normalized patterns produce improved reconstruction performance. In [20] a hybrid normalized and Gaussian thresholded mask was utilized which was found to outperform the masks proposed in [19] and [17].

Differently from these works, our proposed approach is data-driven and does not impose any mask constraints but instead generates mask patterns learnt from the training data. To the best of our knowledge this is the first study that investigates the construction of an optimized binary temporal VCS mask through DNNs.

Concurrent with our work, the method in [21] is developed for learning sensor's color multiplexing patterns for image demosaicking. The sensor's mask is jointly learnt with reconstruction, however, the task of demosaicking differs from ours. To that extend, the learning process is also fundamentally different; The task in [21] is to learn a mask to use one of the discrete set of color filters at each pixel location. To learn such masks a good choice is to apply a softmax function to the weights during training (therefore, estimating class probabilities for each color channel) and create an one-hot vector during testing to indicate color channel selection. In our case, the compressive video measurements are not discrete and the theory of CS suggests that to be able to recover the unknown signal a weighted linear combination of samples is required. Thus, multiple ones and zeros may be realized in each location during the acquisition of the frames and thresholded binary weights (thus, not class probabilities) are estimated during training and testing for accurate VCS reconstruction.

DNNs for image recovery. The capabilities of deep architectures have been investigated in image recovery problems such as deconvolution [22], [23], [24], denoising [25], [26], [27], [28], [29], inpainting [30], and super-resolution [31], [32], [33], [34]. Deep architectures have also been proposed for CS of still images. In [35], stacked denoising auto-encoders (SDAs) were employed to

learn a mapping between the CS measurements and image blocks. A similar approach was also utilized in [36], [37] but instead of SDAs, convolutional neural networks (CNNs) were used in [36] and residual networks (ResNets) in [37].

A closely related study is our previous work in [38] which focuses on learning to map directly temporal VCS measurements to video frames using deep fully-connected networks when the measurement matrix is fixed. We showed that the deep learning framework enables the recovery of video frames from temporal compressive measurements in a few seconds at significantly improved reconstruction quality compared to different, optimization based, schemes.

Binary neural networks. Recently, several approaches have been proposed on the development of neural networks with binary weights [39], [40], [41], [42] for image recognition applications. The main objective of such an approach is to simplify computations in neural networks, thus making them more efficient while requiring reduced storage. Efficiency is achieved by approximating the standard real-valued DNNs with binary weights. In BinaryConnect [40] the authors proposed to binarize the weights for all layers during the forward and backward propagations while keeping the real-valued weights during the parameter update. The real-valued updates were found to be necessary for the application of stochastic gradient descent (SGD). Performance with various classification tasks demonstrated that binary neural networks compare favorably with real-valued weight networks. In [42], the authors introduced a weight binarization scheme where both a binary filter and a scaling factor are estimated. Such scheme was proven more effective compared to the BinaryConnect.

Motivation for using DNNs to learn mask parameters. Motivated by the success of DNNs in CS reconstruction and binary DNNs in classification, we investigate in this paper the problem of learning an optimized binary sensing matrix using DNNs for temporal VCS.

The work presented in this paper is different from the studies in image recovery using DNNs and from the binary neural networks. First, this work is different from our work in [38] since our focus is on learning an optimized sensing mask along with the video reconstruction. In [38] the scope was to recover video frames directly from the temporal measurements (i.e., the mask is pre-defined). Furthermore, our objective in this paper is to learn binary masks that will encode video frames on VCS cameras for video reconstruction which is different from that in binary neural network studies, which is efficiency for image recognition problems.

III. DEEPBINARYMASK

In this work, we propose a novel neural network architecture that learns to *encode* a three dimensional (3D) video block to compressive two-dimensional (2D) measurements by learning the binary weights of Φ and to *decode* the measurements back to a video block, as illustrated in Figure 2. Let us now describe in detail the encoder and decoder.

A. Encoder

In order for our learning approach to be practical, reconstruction has to be performed on 3D video blocks [36], [38]. Thus, each video block must be sampled with a block-based measurement matrix which should be the same for all blocks. Furthermore, such a measurement matrix should be realizable in hardware. We follow the pattern in [38] and we consider a Φ which consists of repeated identical building blocks of size $w_p \times h_p \times t = N_p$ corresponding to the matrix Φ_p of size $M_p \times N_p$, where $M_p = w_p \times h_p$. In other words Φ_p has the structure shown in Eq. (3), in which M_f and N_f have been respectively replaced by M_p and N_p . An implementation of such a matrix on existing systems employing Digital Micromirror Devices (DMDs), Spatial Light Modulators (SLMs) or Liquid Crystal on Silicon (LCoS) [6], [18], [19], [16], [5] can easily be performed. At the same time, a repeated mask can be printed and shifted appropriately to produce the same effect in systems utilizing translating masks [20], [17].

Let us consider a set of N training 3D video blocks, each of size $w_p \times h_p \times t$. They are lexicographically ordered by considering first the spatial and then the temporal dimensions to form vectors \mathbf{x}_i , each of size $N_p \times 1$. The encoder is defined as the mapping $g(\cdot)$ that transforms each \mathbf{x}_i to a measurement \mathbf{y}_i of size $M_p \times 1$, which represents the lexicographically ordered $w_p \times h_p$ image patch, followed by a non-linearity given as,

$$\mathbf{y}_i = g(\mathbf{x}_i; \theta_e) = \sigma_e(\Phi_p \mathbf{x}_i), \quad (4)$$

where $\theta_e = \{\Phi_p\}$ is the parameter set and function $\sigma_e(\cdot)$ represents the non-linearity. We use the subscript “e” to denote quantities pertaining to the encoder, in order to distinguish them from the decoder quantities to be introduced later.

The formulation in (4) would have been straightforward to handle if matrix Φ_p were dense and consisting of real-values. However, as mentioned earlier, in the case of temporal VCS, matrix Φ is binary (due to implementation considerations) and sparse following the structure defined in (3). For ease of presentation let us now also define a matrix $B \in \{0, 1\}^{t \times M_p}$ containing the binary weights as,

$$B = [\mathbf{b}_1, \dots, \mathbf{b}_{M_p}] = \begin{bmatrix} b_{1,1} & \dots & b_{1,M_p} \\ \vdots & \ddots & \vdots \\ b_{t,1} & \dots & b_{t,M_p} \end{bmatrix}. \quad (5)$$

It is related to the measurement matrix Φ_p as,

$$\Phi_p = \begin{bmatrix} b_{1,1} & \mathbf{0} & \mathbf{0} & b_{t,1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \ddots & \mathbf{0} & \dots & \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & b_{1,M_p} & & \mathbf{0} & \mathbf{0} & b_{t,M_p} \end{bmatrix}. \quad (6)$$

In order to realize such a structure in a neural network and be able to train it we transform the encoder into a network that involves the following steps:

- 1) The first step consists of M_p binary parallel layers. To describe this step we need to introduce a new column

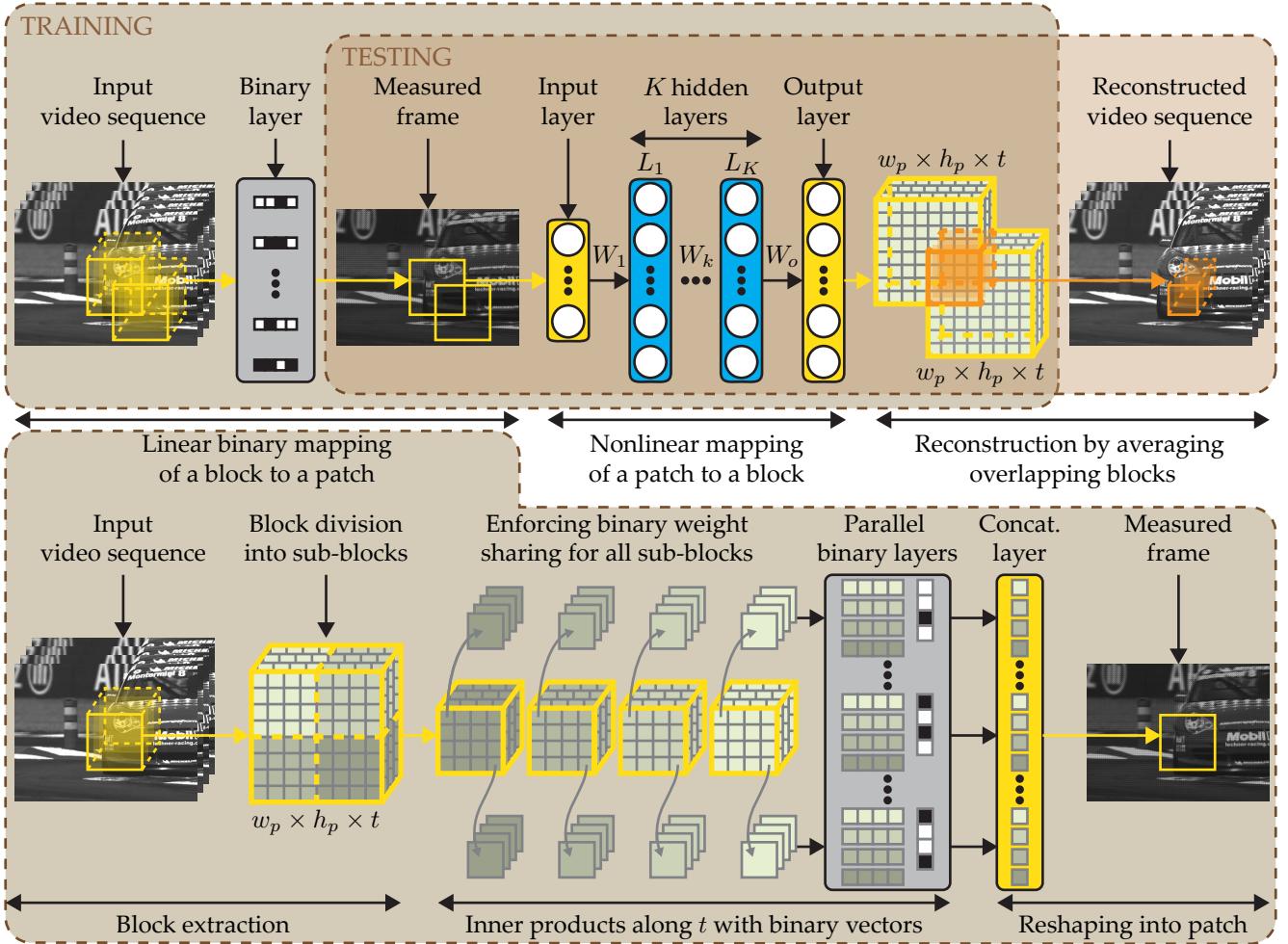


Fig. 2. Illustration of the proposed encoder-decoder neural network for video compressive sensing. The bottom part demonstrates the *encoder* network that is responsible for learning the binary mask and outputs CS measurements. The upper part, labeled as “TESTING” illustrates the *decoder* network which takes as input CS measurements and outputs a video sequence.

$(t \times 1)$ vector $\mathbf{x}_{i,j}$, which consists of all the temporal elements at a given spatial location j , that is,

$$\mathbf{x}_{i,j} = \begin{bmatrix} \mathbf{x}_i(j) \\ \mathbf{x}_i(M_p + j) \\ \mathbf{x}_i(2M_p + j) \\ \vdots \\ \mathbf{x}_i((t-1)M_p + j) \end{bmatrix}, \quad (7)$$

where $\mathbf{x}_i(j)$ denotes the j -th element of vector \mathbf{x}_i . Then in parallel the following inner products are computed,

$$e(\mathbf{x}_{i,j}) = \mathbf{b}_j^T \mathbf{x}_{i,j}, \quad \text{for } j = 1, \dots, M_p. \quad (8)$$

- 2) The second step consists of a concatenation layer which concatenates the outputs of the parallel layers in order to construct a single measurement vector that is,

$$\mathbf{y}_i = g(\mathbf{x}_i; \theta_e) = \text{concat}\left(e(\mathbf{x}_{i,1}), \dots, e(\mathbf{x}_{i,M_p})\right), \quad (9)$$

with a parameter set $\theta_e = \{\mathbf{b}_1, \dots, \mathbf{b}_{M_p}\}$, as defined by Eqs. (5) and (6). Note, that a non-linearity such as the rectified linear unit (ReLU) [43] defined as, $\sigma(z) = \max(0, z)$, is implicitly applied here after the

concatenation since the output is always positive. This is due to the fact that the weights are binary with values 0 and 1 and the video inputs have non-negative values.

The above two steps follow the model presented in Figure 1 but translated to a neural network, where the set θ_e consists of the elements of the trained projection matrix. The two steps of the encoder are illustrated at the bottom part of Figure 2. Note that the figure refers to the encoding of overlapping blocks, as we describe next.

Overlapping blocks and weight sharing. The $t \times M_p$ binary weight matrix B we have considered so far corresponds to non-overlapping video blocks. In order to realize overlapping blocks which usually aid in improving reconstruction quality we can utilize repeating blocks of dimensions $\frac{w_p}{2} \times \frac{h_p}{2} \times t$, which we call sub-blocks as shown in Figure 2. Thus, for the final trained matrix Φ_p each $\frac{w_p}{2} \times \frac{h_p}{2} \times t$ sub-block is the same allowing reconstruction of overlapping blocks of size $w_p \times h_p \times t$ with spatial overlap of $\frac{w_p}{2} \times \frac{h_p}{2} = w_s \times h_s$, as presented in Figure 3. In such a case the parameter set θ_e is also different. Instead of learning M_p binary weight vectors we learn $M_p/4$, where each weight vector is *shared* four

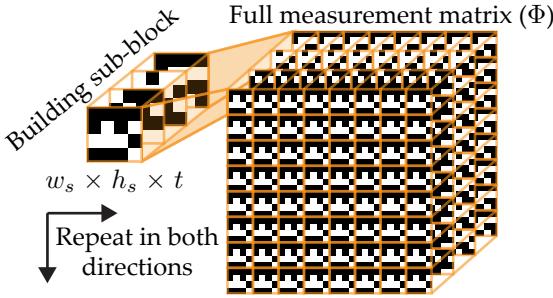


Fig. 3. Construction of the full measurement matrix by repeating a three dimensional random binary array (building sub-block) in the horizontal and vertical directions.

times for each of the corresponding pixel positions of the input. For example, in the case when $w_p \times h_p = 8 \times 8$ there will be four identical $4 \times 4 \times t$ sub-block projection matrices. Notice in Figure 2 that the values of the input block at the corresponding pixel locations at each of the sub-blocks are multiplied by the same binary vector. Thus, for this example, we only need to estimate 16 binary weight vectors and each one is shared by four different inputs. For instance, in order to calculate $e(\mathbf{x}_{i,1}), e(\mathbf{x}_{i,5}), e(\mathbf{x}_{i,33}), e(\mathbf{x}_{i,37})$ the weight vector \mathbf{b}_1 is used.

Binary weights. Let us now proceed to describe how to estimate the binary weights. We follow the BinaryConnect method [40] to constrain the weights of the encoder to be equal to either 0 or 1 during propagation. The binarization scheme to transform the real-valued weights to binary values is based on the sign function, that is,

$$b_b = \begin{cases} 1 & \text{if } b_r \geq 0, \\ 0 & \text{otherwise,} \end{cases} \quad (10)$$

where b_b and b_r are the binarized and real-valued weights of B , respectively. Following the training process in [40] we binarize the weights of the encoder only during the forward and backward propagations. The update of the parameter set θ_e is performed using the real-valued weights. As explained in [40], keeping the real-valued weights during the updates is necessary for training the networks using SGD. In addition, we enforced the real-valued weights to lie within the $[-1, 1]$ interval at each training iteration. The weight clipping was chosen since otherwise the weights may become infinitely large having no impact during binarization.

Weight initialization. The network weight initialization of the encoder corresponds in our case to the mask initialization. Typically, in VCS the mask is generated randomly. Similarly here, we start with a randomly generated mask by a Bernoulli distribution. However, since real-valued weights are also required by the network to perform their updates we consider the following initialization scheme,

$$b_b \sim \text{Bern}(p),$$

$$b_r = \begin{cases} \sim \text{Unif}(0, 1/\sqrt{t}) & \text{if } b_b = 1, \\ \sim \text{Unif}(-1/\sqrt{t}, 0) & \text{otherwise,} \end{cases} \quad (11)$$

where $\text{Bern}(\cdot)$ and $\text{Unif}(\cdot)$ denote the Bernoulli and Uniform distributions, respectively, p is the probability of the weight to be initialized with 1 and notation (\cdot, \cdot) refers to the lower and upper bounds for the values of the distribution. The bounds of the Uniform distribution follow the scheme introduced in [44]. The initialization scheme proposed above allows us to fully understand the benefits of learning as compared to non-learning the mask along with reconstructing the video. This is due to the fact that in the case of choosing the non-learning mode we keep the initial b_b weights, drawn from the Bernoulli distribution, fixed.

B. Decoder

The resulting hidden measurement \mathbf{y}_i produced by the encoder is then mapped back to a reconstructed $N_p \times 1$ vector through the decoder $f(\mathbf{y}_i; \theta)$, which when unstacked results in the $w_p \times h_p \times t$ dimensional video block, as illustrated in the upper part of Figure 2. Thus, the decoder of the proposed method is another network which is trained to reconstruct the video output sequence given \mathbf{y}_i . We consider an MLP architecture to learn a nonlinear function $f(\cdot)$ that maps a measured frame patch \mathbf{y}_i via several hidden layers to a video block \mathbf{x}_i as in [38].

The output of the k^{th} hidden layer L_k , $k = 1, \dots, K$ is defined as,

$$h_k(\mathbf{y}_i) = \sigma_d(W_k h_{k-1}(\mathbf{y}_i) + \mathbf{c}_k), \quad \text{with } h_0(\mathbf{y}_i) = \mathbf{y}_i, \quad (12)$$

where W_k is the output weight matrix, and $\mathbf{c}_k \in \mathbb{R}^{N_p}$ the bias vector. $W_1 \in \mathbb{R}^{N_p \times M_p}$ connects \mathbf{y}_i , the output of the encoder, to the first hidden layer of the decoder, while for the remaining hidden layers, $\{W_2, \dots, W_K\} \in \mathbb{R}^{N_p \times N_p}$. The last hidden layer is connected to the output layer via $\mathbf{c}_o \in \mathbb{R}^{N_p}$ and $W_o \in \mathbb{R}^{N_p \times N_p}$ without nonlinearity. The nonlinear function $\sigma_d(\cdot)$ is the ReLU and the weights of each layer are initialized to random values uniformly distributed in $(-1/\sqrt{N_p}, 1/\sqrt{N_p})$ [44].

C. Choice of network architecture for decoder

One could argue that a CNN architecture could be used, instead of an MLP one. There are several reasons why the MLP architecture for the decoder is a reasonable choice for the temporal video compressive sensing problem which have been explained in [38]. First, unlike other imaging problems (e.g., deconvolution) the measuring process in video CS cannot be modeled as a convolution since spatially neighboring pixels do not contribute to each measured pixel in a 2D patch. Second, since we want to move from 2D CS measurements, to 3D video blocks, a fully-connected layer should be employed as a first layer to increase the dimensionality of the unknown variables to be estimated. Clearly, one could employ CNN layers after the first layer and still be able to recover a reasonable video reconstruction. However, the small size of patches and video blocks used for reconstruction, in order to make training feasible, would not allow for convolutions to be effective. Indeed, we experimented with architectures that contained subsequent convolutional layers and consistently obtained worse or similar performance to the one reported

using a full MLP architecture. Since our focus in this work is to primarily investigate and compare the performance of the trained versus the non-trained sensing matrix we adopt the decoder design in [38].

D. Training the encoder-decoder network

The two components of the proposed MLP encoder-decoder are jointly trained by learning all the weights and biases of the model. Using spatial overlap $\frac{w_p}{2} \times \frac{h_p}{2}$ the set of all parameters is denoted by $\theta = \{\mathbf{b}_1, \dots, \mathbf{b}_{M_p/4}; W_1, \dots, W_K; W_o; \mathbf{c}_1, \dots, \mathbf{c}_K; \mathbf{c}_o\}$ and is updated by the backpropagation algorithm [45] minimizing the quadratic error between the set of the encoded mapped measurements $f(\mathbf{y}_i; \theta)$ and the corresponding video blocks \mathbf{x}_i . The loss function is the Mean Squared Error (MSE) which is given by,

$$L(\theta) = \frac{1}{N} \sum_{i=1}^N \|f(\mathbf{y}_i; \theta) - \mathbf{x}_i\|_2^2. \quad (13)$$

The MSE was used in this work since our goal is to optimize the Peak Signal to Noise Ratio (PSNR) which is directly related to the MSE.

Training procedure. The overall training procedure can be summarized by the following steps:

- 1) Forward propagation is performed by using weights B after binarization in the encoder and real-valued weights W in the decoder.
- 2) Then, backpropagation is performed to compute the gradients with respect to layer's activation knowing B and W .
- 3) Parameter updates are computed using the real-valued weights for both encoder and decoder.

Note that one other difference between our work and [40] is that our encoder-decoder neural network does not utilize binary weights in all layers; instead it utilizes binary weights at the encoder and standard real-valued weights at the decoder.

Implementation details. Our encoder-decoder neural network is trained for 480 epochs using a mini-batch size of 200. We used SGD with the momentum set equal to 0.9. We further used ℓ_2 norm gradient clipping to keep the gradients in a certain range. Gradient clipping is a widely used technique in recurrent neural networks to avoid exploding gradients [46]. The threshold of gradient clipping was set equal to 0.1.

One hyper-parameter that was found to affect the performance in our approach is the learning rate. Based on experimentation we chose a starting learning rate for the encoder that was 10 times larger than that for the decoder. This was found to be important as we wanted the weights of the encoder to have their sign changed during the training iterations. In addition, the learning rate was divided by 2 at every 10 epochs in the encoder and by 10 after 400 epochs in the decoder.

All hyper-parameters were selected after cross-validation using a validation test set.

Test inference. Once the encoder-decoder neural network is trained we use the trained sensing matrix $B \rightarrow \Phi$ to calculate the compressive measurements \mathbf{y} . Then, given \mathbf{y} we can use any VCS algorithm (in addition to the decoder network) to reconstruct the video blocks.

IV. EXPERIMENTAL RESULTS

In this section we present quantitative and qualitative reconstruction results to demonstrate the effectiveness of the proposed projection mask in temporal VCS. The performance of our trained masks is investigated using various reconstruction algorithms and initial mask parameters. Our analysis offers insights into understanding how the different initial parameters of the mask affect reconstruction performance. The metrics used for reconstruction evaluation were the PSNR and SSIM (Structural SIMilarity).

A. Training data collection and test set

In order to train our encoder-decoder architecture we collected a diverse set of training samples using 400 high-definition videos from YouTube, depicting natural scenes. The video sequences contain more than 10^5 frames which were converted to grayscale. We randomly extracted 1 million video blocks of size $w_p \times h_p \times t$ to train our encoder-decoder neural network while keeping the amount of blocks extracted per video proportional to its duration.

Our test set consists of 14 video sequences that were used in [19] which are provided by the authors. We also included in the test set the “Basketball” video sequence used in [47]. All test video sequences are unrelated to the training set.

B. Mask patterns and decoding layers

Our experimental investigation is motivated by the following two questions: 1) “How does performance of trained and non-trained masks compare using different reconstruction algorithms?” and 2) “Does the training procedure result in a unique sensing matrix Φ irrespectively of the initialization parameters?”

In order to answer these two questions we simulated noiseless compressive video measurements by realizing four different $\frac{w_p}{2} \times \frac{h_p}{2} \times t$ mask patterns. We denote by “RandomMask- p ” the mask that is initialized with Bern(p), as in Eq. (11), and is not learnt (that is, the elements of the encoder are fixed). We also denote by “DeepMask- p ” the learnt mask trained by our proposed encoder-decoder network described in section III and which is initialized by Bern(p) as in Eq (11). Thus, we consider the following four mask patterns:

- RandomMask-20 and DeepMask-20, with $p = 20\%$.
- RandomMask-40 and DeepMask-40, with $p = 40\%$.
- RandomMask-60 and DeepMask-60, with $p = 60\%$.
- RandomMask-80 and DeepMask-80, with $p = 80\%$.

For the remainder of this paper, we describe the selection of block sizes of $w_p \times h_p \times t = 8 \times 8 \times 16$, such that $N_p = 1024$ and $M_p = 64$. Therefore, the compression ratio is 1/16. Although larger block sizes can be used in our framework, block sizes of 8×8 have provided good

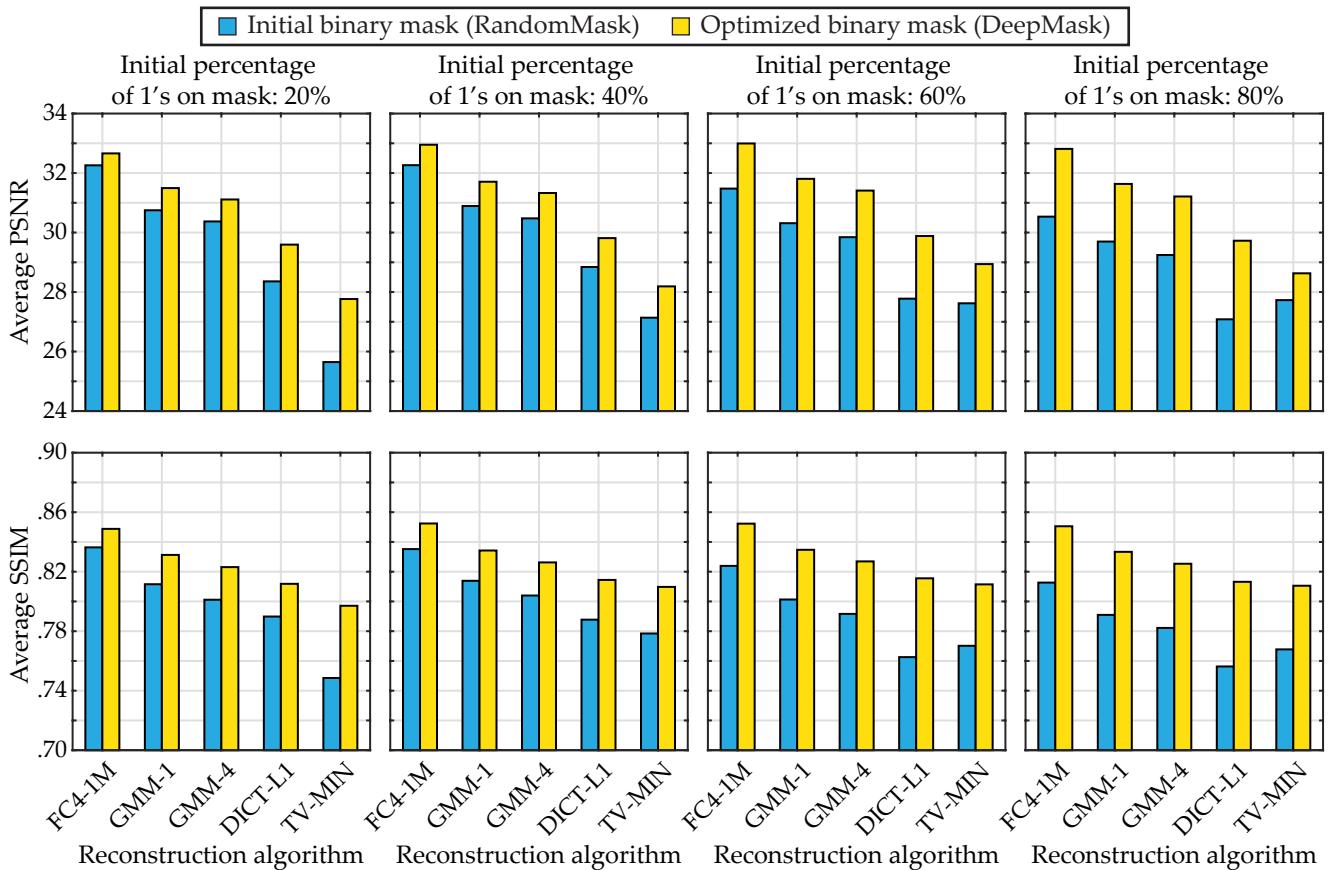


Fig. 4. Average PSNR and SSIM over all test video sequences for several reconstruction methods using the RandomMasks and DeepMasks. The test set consists of 14 video sequences and the reported PSNR and SSIM corresponds to the average values for the reconstruction of the first 32 frames of each sequence. The PSNR metric is measured in dB while the SSIM is unitless.

reconstruction quality with sharper reconstructed frames in learning approaches used for CS video reconstruction [19], [48]. It was shown that choosing larger block sizes led to worse reconstruction quality. In addition, as mentioned earlier, the small size of patches and video blocks used for reconstruction make training feasible. Thus, we adopted this same setting in our approach. Furthermore, for each of the eight Φ mask types above, each $\frac{w_p}{2} \times \frac{h_p}{2} \times t = 4 \times 4 \times 16$ block is the same allowing reconstruction for overlapping blocks of size $8 \times 8 \times 16$ with spatial overlap of 4×4 . Note that the same random seed was utilized for all patterns.

Further, we used $K = 4$ hidden layers for the decoder architecture of Figure 2. We found out experimentally that for the number of training data used (1 million video blocks) 4 layers provided the best performance. A similar observation was reported in [38] where the addition of extra layers for this number of training data did not lead to performance improvement.

In the following section we present the compressive sensing reconstruction algorithms used to test the eight mask types.

C. Reconstruction Algorithms

Since our main goal is to compare the performance between a trained sensing mask over a non-trained one in an implementation agnostic to mask patterns, we tested a number

of different reconstruction algorithms. Candidate reconstruction algorithms were selected for their utility in solving the underdetermined system in the video compressive sensing setting. We evaluated the following optimization algorithms as potential solvers:

- 1) **DICT-L1:** In (1), we have described an underdetermined system where data are noise-free. However, real data are typically noisy and dealing with small dense noise is required. In order to deal with such noise we transform the problem in (2) into the LASSO (Least-Absolute Shrinkage and Selection Operator) problem for $F = \lambda \|\mathbf{a}\|_1$ given as,

$$\hat{\mathbf{a}} = \underset{\mathbf{a}}{\operatorname{argmin}} \|\mathbf{y} - \Phi D \mathbf{a}\|_2^2 + \lambda \|\mathbf{a}\|_1, \quad (14)$$

where $\lambda > 0$ is the regularization parameter whose value is related to the noise tolerance. For this problem we chose to use an overcomplete dictionary D as a sparsifying basis. The dictionary consists of 20,000 atoms trained on a subset of 200,000 video blocks from our training database and reconstruction is performed block-wise on overlapping sets of 7×7 patches of pixels. For the optimization problem in (14), λ was set equal to 0.005.

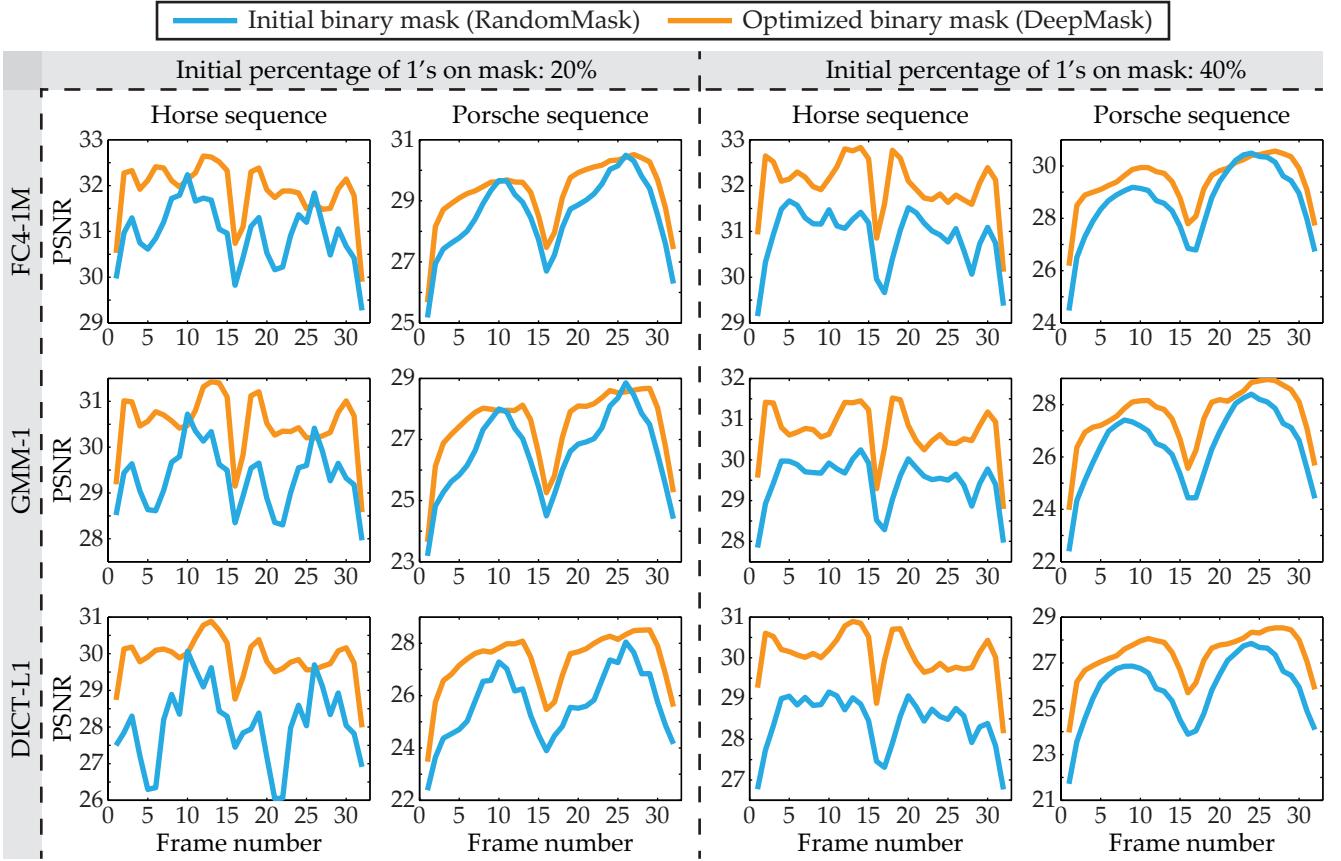


Fig. 5. PSNR (in dB) comparison for the first 32 frames of 2 video sequences among the proposed method FC4-1M and the previous methods GMM-1 [48] and DICT-L1 [19]. Notice that the vertical scale changes among the various plots.

- 2) **TV-MIN:** A popular CS reconstruction method utilizes for $F(\cdot)$ the total variation (TV) norm defined as,

$$TV(\mathbf{z}) = \sum_{i,j,n} \left((z(i+1, j, n) - z(i, j, n))^2 + (z(i, j+1, n) - z(i, j, n))^2 \right)^{1/2}, \quad (15)$$

where \mathbf{z} is the stacked version of the 3D array $z(i, j, n)$, where (i, j, n) are respectively the two spatial and one temporal coordinates. Thus, the TV minimization problem is given as,

$$\hat{\mathbf{x}} = \operatorname{argmin}_{\mathbf{x}} \|\mathbf{y} - \Phi \mathbf{x}\|_2^2 + \lambda TV(\mathbf{x}). \quad (16)$$

In order to solve (16) we used the two-step iterative shrinkage/thresholding (TwIST) algorithm [49] with $\lambda = 0.01$.

- 3) **GMM-TP:** Another reconstruction algorithm we considered in our experiments is a Gaussian mixture model (GMM)-based algorithm [48] learned from Training Patches (TP) referred as GMM-TP. We followed the settings proposed by the authors and used our training data (randomly selecting 20,000 samples) to train the underlying GMM parameters. In our experiments we refer to this method by GMM-4 and GMM-1 to denote

reconstruction of overlapping blocks with spatial overlap of 4×4 and 1×1 pixels, respectively.

- 4) **FC4-1M:** Finally, another reconstruction method we considered is the decoder neural network introduced in subsection III-B. The decoder is a $K = 4$ MLP trained on 1 million samples similarly to [38]. In this case, a collection of overlapping patches of size 8×8 is extracted by each coded measurement of size $W_f \times H_f$ and subsequently reconstructed into video blocks of size $8 \times 8 \times 16$. Overlapping areas of the recovered video blocks are then averaged to obtain the final video reconstruction results as shown in the upper part of Figure 2. The step of the overlapping patches was set to 4×4 due to the special construction of the utilized measurement matrix, as discussed in subsection III-A.

For each algorithm, λ values were determined based on the best performance among different settings. All code implementations are publicly available provided by the authors while the deep network architectures were implemented in Torch7 [50], a Lua library that allowed us to develop an optimized GPU code.

D. Reconstruction Results

For each reconstruction algorithm described above, we tested the eight mask types presented in subsection IV-B.



Fig. 6. Qualitative reconstruction performance. The figure shows reconstruction of a single frame for 2 test video sequences when using different reconstruction algorithms and different mask initializations. The PSNR metric is measured in dB while the SSIM is unitless.

Quantitative and qualitative results. Figure 4 shows average reconstruction quality for each mask and algorithm combination, using the PSNR and SSIM metrics. The presented metrics refer to average performance for the reconstruction of the first 32 frames of each test video sequence, using 2 consecutive captured coded frames for each of the eight masks for every algorithm. First, we note that the DeepMasks perform consistently better compared to the RandomMasks across all reconstruction algorithms and initial percentage of nonzeros. In particular, we observe an improvement around 1-2 dB, in terms of PSNR between the trained and non-trained masks across all initial percentages and algorithms. Furthermore, we observe that the decoder FC4-1M demonstrates the highest PSNR and SSIM values among all algorithms.

Figure 5 compares the PSNR for each of the 32 frames of 2 video sequences (“Horse” and “Porsche” sequences) using our FC4-1M algorithm and the previous methods GMM-1 [48] and DICT-L1 [19] between the RandomMasks and DeepMasks. The varying PSNR performance across the frames of a 16 frame block is consistent for all algorithms and is reminiscent of the reconstruction tendency observed in other video CS papers in the literature [20], [17], [47], [48]. Please notice that the scale on the vertical axes of the plot varies.

Finally, Figure 6 compares the reconstruction quality between the optimized and non-optimized masks for a single frame from the 2 video sequences under different algorithms and different mask initializations. Consistent with the PSNR and SSIM results, it is clear from the visual evaluations that reconstruction quality improves when using the optimized masks, especially for the case when the initial percentage of nonzeros is $p = 60$. In particular, background details of the “Horse” sequence are more visible and letters in the “Porsche” sequence appear sharper. At the same time, the proposed end-to-end optimization using the deep network (FC4-1M) provides the highest visual quality among the candidate compared algorithms. In general, it is observed that the optimized masks reconstruct the frames with less blurring and sharper edges than the non-optimized masks.

E. Training analysis

We start our analysis by examining the real-valued weight histogram of the encoder (DeepMask-40) upon convergence in Figure 7. First, we observe that negative values are more frequent than positive ones, which suggests that zero elements of the mask (after binarization) are more important than the nonzero ones. More importantly, we observe that a number of weights are around zero, hesitating between becoming negative or positive, a phenomenon that was also reported in [40].

Average test MSE per epoch calculated on a validation test set for DeepMask-40 and RandomMask-40 is shown in Figure 8. It is shown that the test error curve of the RandomMask-40 is smooth while the DeepMask-40 is noisy. This is due to the fact that many binary weights switch between 1 and 0 frequently, especially during the first epochs of training when the learning rate has high values, thus constantly changing the way the VCS measurements are performed. Furthermore,

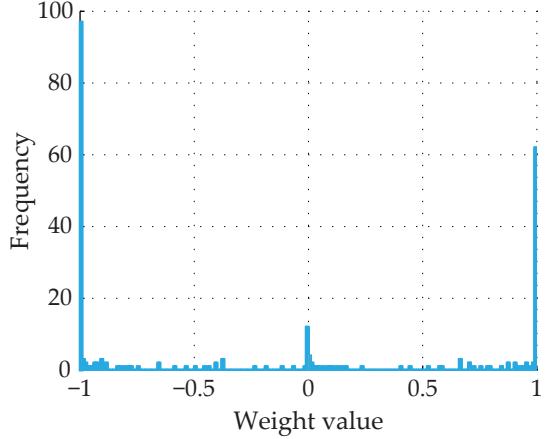


Fig. 7. Histogram of the real-valued weights produced by the encoder neural network for DeepMask-40. We report similar observations with [40] as we found out that most of the weights have the tendency to become deterministic (-1 and 1) and reduce the training error while some stay around zero.

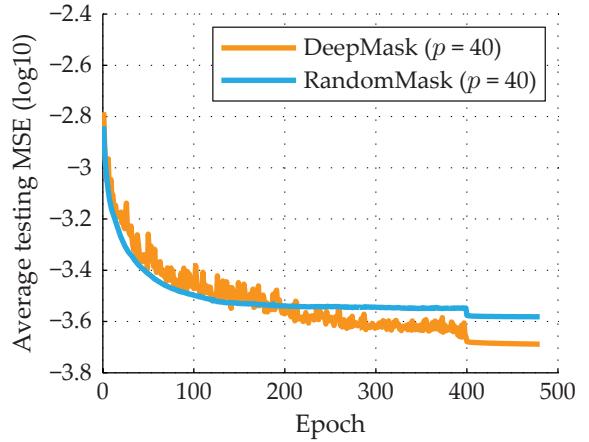


Fig. 8. Test error curves between the RandomMask-40 decoder and DeepMask-40 encoder-decoder calculated on a validation test set. The latter provided lower test error upon convergence as is optimized end-to-end.

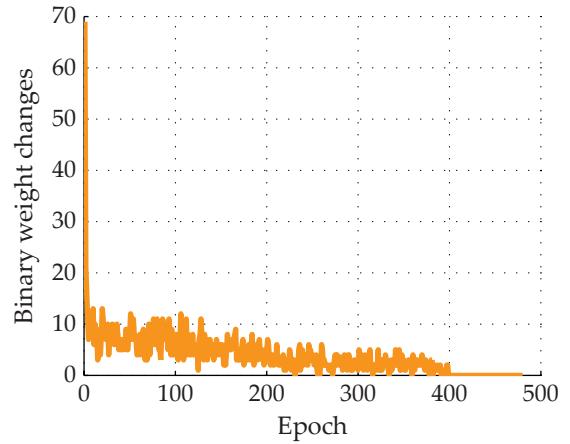


Fig. 9. Number of binary weight changes per epoch of DeepMask-40 encoder-decoder. A large number of weights change in the first few epochs; this number decreases and finally becomes zero in the last few epochs.

even at the later stages of training, many real-valued elements of the encoder remain around zero, as observed in Figure 7.

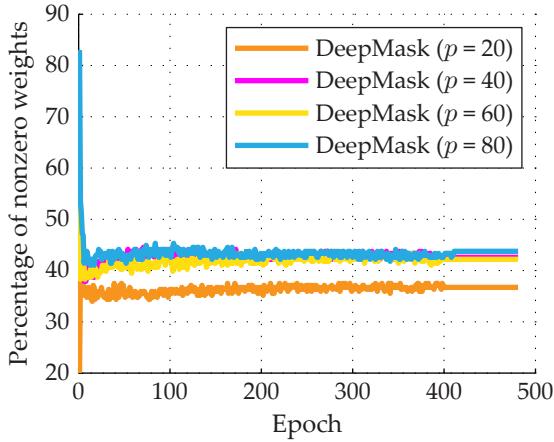


Fig. 10. Percentage of nonzero binary values for DeepMasks. Irrespectively to the initial nonzero percentage, DeepMasks converge to a point around 40% nonzeros.

Therefore, during the binarization process some of the encoder's binary weights change from zero to one and vice versa even with very small learning rates. However, as the encoder's learning rate becomes really small the curve becomes smooth. A better optimized learning rate decay schedule of the encoder would have probably provided a smoother curve and perhaps a higher performance. We leave this as a task for future work as further investigation into this may be needed. Finally, as showed in the reconstruction results, DeepMask performs consistently better than RandomMask which also explains the lower test MSE produced by the former during training.

Lastly, in Figure 9 we show the number of binary weights that change from zero to one and vice versa per epoch. We observe that a large number of weights change in the first few epochs and this number decreases as the number of epochs increases.

V. DISCUSSION

Having obtained better reconstruction performance using the DeepMasks across a wide range of reconstruction algorithms our next step is to analyze the masks produced by the networks and highlight a few crucial points. Note that the analysis provided below is based on the specific dataset and MLP architecture used in this work. We start our analysis by posing the following question.

Does DeepMask produce a unique sensing matrix Φ ? To answer this question we examine the differences between the produced DeepMasks with respect to their percentage of nonzero elements and to their support.

First, in Figure 10 we show the percentage of nonzero binary weights per epoch for the different DeepMasks. This figure allows us to examine uniqueness with respect to the percentage of nonzero elements produced by each DeepMask. We observe that the masks with p equal 40, 60 and 80 converge to a percentage around 40%. The $p = 20$ mask though, converges to a bit lower percentage (around 38% nonzero binary elements).

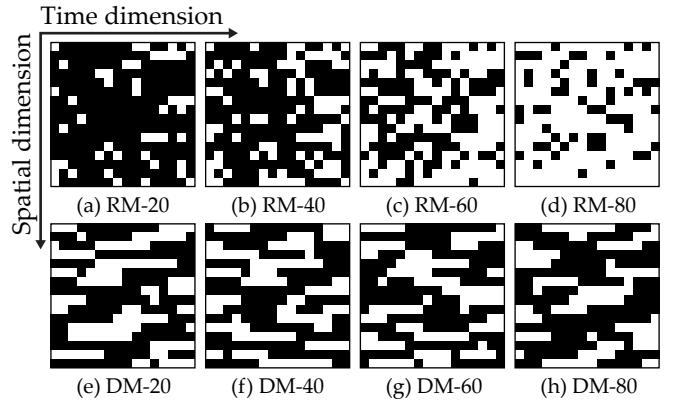


Fig. 11. The eight mask patterns produced in this work of size 16×16 . The initial RandomMasks (RM) are presented on the top row while on the bottom row we present the optimized DeepMasks (DM).

Next, Figure 11 demonstrates the four mask patterns produced in this work. The first row illustrates the RandomMasks and the second row presents the four DeepMasks produced by the proposed encoder-decoder neural network. All $4 \times 4 \times 16$ masks are reshaped into a 16×16 matrix for better visualization by lexicography ordering each 4×4 mask at a given time instance into a 16×1 vector, which becomes a column of the 16×16 matrix. That is, the vertical direction denotes pixel location while the horizontal direction denotes time. From the visualization we deduce that although the masks generated by the network converge to the same nonzero percentage (as shown in Figure 10), their support is different. The fact that the optimized masks contain a very similar percentage of nonzeros while producing improved reconstruction quality with various different reconstruction algorithms implies that such percentage is the most appropriate one for the task at hand. Similar observations about the ideal percentage of nonzeros for VCS measurement matrices have been made in [20], albeit deduced through heuristic experimentation.

Finally, from the visualization we deduct two important findings:

- First, it is apparent that regardless of the initial realization (shown in RandomMasks), the trained DeepMasks produce a similar number of nonzero elements which confirms our findings discussed earlier.
- Second, an important observation from Figure 11 is that DeepMasks are *smoother* over time than the RandomMasks. In other words in many rows the binary weights seem to be sequential (or more structured) forming runs of 1s and 0s. Again, such finding was heuristically observed in [20] and some studies cited therein, further strengthening our findings which are here obtained through a machine learning approach.

To summarize, our observations above suggest that an optimized mask design Φ for temporal VCS incorporates the following two characteristics: 1) *smoothness* as explained above and 2) percentage of *nonzero elements around 40%*.

Random seed selection. Finally, we wanted to confirm that the results presented herein are not due to a specific

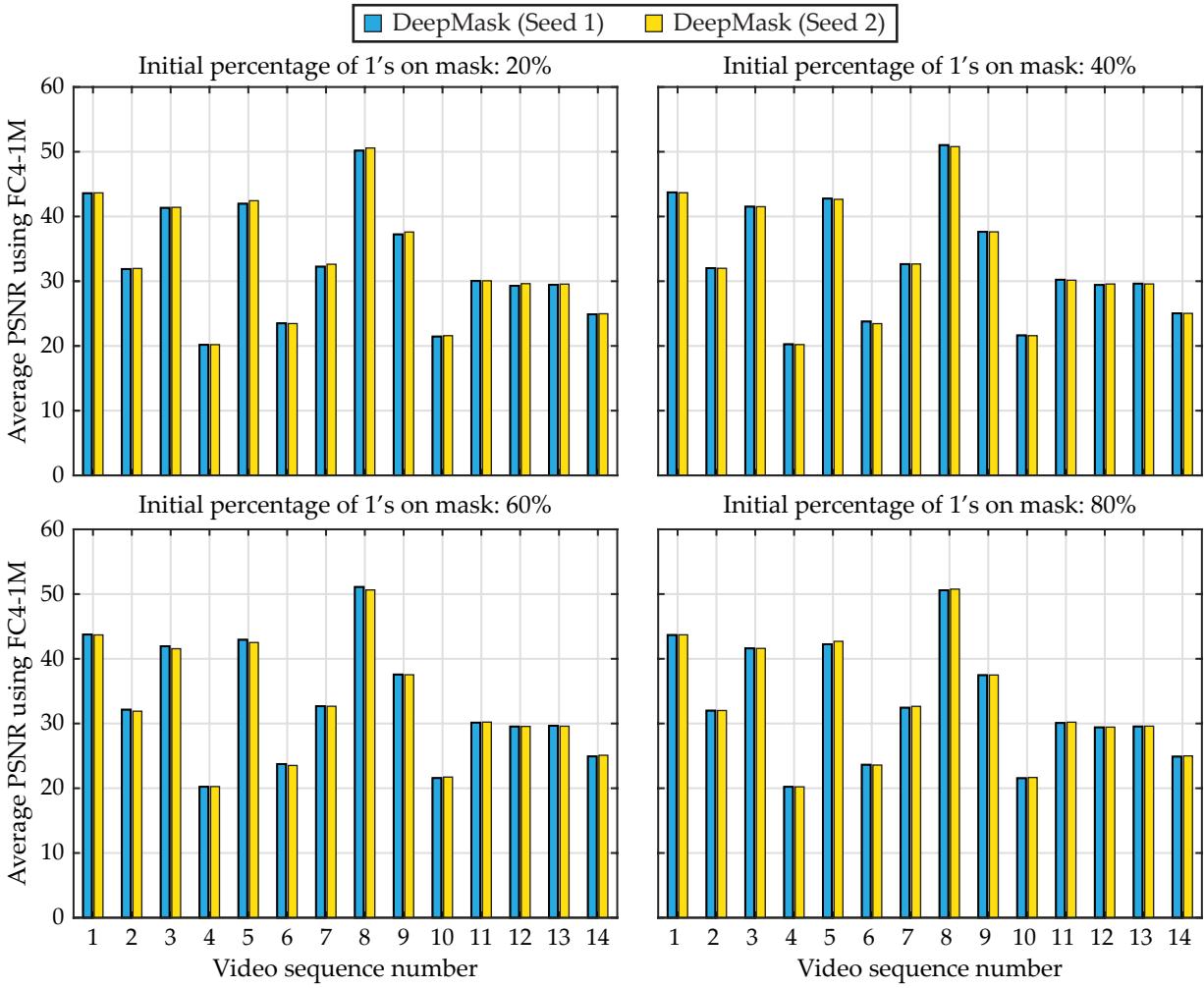


Fig. 12. Comparison of the reconstruction performance of the proposed encoder-decoder architecture when using the optimized masks trained after initialization with two different seeds. Average PSNR for the reconstruction of the first 32 frames of each one of the 14 test video sequences is presented and the values are found to be very similar regardless of the starting binary values of the measurement matrix.

selection of the random seed, used to produce the initial random masks. Therefore, we performed the whole training process described in section IV using a second seed for four new masks and compared the average reconstruction performance using the trained network for all test video sequences. The corresponding results are presented in Figure 12 where it can be observed that the final performance is very similar for both seeds. We do not include the results of the competitive algorithms with this second initialization but observed performance improvements through the new optimized masks similar to the ones presented in Figure 4.

VI. CONCLUSIONS

In this paper, we proposed a new encoder-decoder neural network architecture for video compressive sensing that is able to learn an optimized binary sensing matrix. We evaluated the proposed model on several video sequences and we documented the superiority of the trained sensing matrices over the random ones both quantitatively and qualitatively. Our qualitative analysis of the trained model shows that the optimized sensing masks converge to a similar number of nonzero elements regardless of their initial parameters and that

they exhibit a smoothness property. The proposed architecture has large potential for further analysis.

One limitation of our approach is that given a new compression ratio to be tested, a new model must be trained. Our next step is to explore ways to overcome this practical issue. Perhaps, a first step towards this goal is to explore whether fine-tuning of our networks to new compression ratios could alleviate the requirement of training from scratch. Another future direction is to examine the reconstruction performance in real video sequences acquired by a temporal compressive sensing camera.

REFERENCES

- [1] D. L. Donoho, "Compressed sensing," *IEEE Trans. Inf. Theory*, vol. 52, no. 4, pp. 1289–1306, April 2006.
- [2] E. J. Candès, J. Romberg, and T. Tao, "Robust uncertainty principles: exact signal reconstruction from highly incomplete frequency information," *IEEE Trans. Inf. Theory*, vol. 52, no. 2, pp. 489–509, Feb. 2006.
- [3] M. F. Duarte, M. A. Davenport, D. Takhar, J. N. Laska, T. Sun, K. F. Kelly, and R. G. Baraniuk, "Single-Pixel imaging via compressive sampling," *IEEE Signal Process. Mag.*, vol. 25, no. 2, pp. 83–91, Mar. 2008.
- [4] A. C. Sankaranarayanan, C. Studer, and R. G. Baraniuk, "CS-MUVI: Video compressive sensing for spatial-multiplexing cameras," in *Proc. IEEE Conf. Computational Photography*, Apr. 2012, pp. 1–10.

- [5] J. Wang, M. Gupta, and A. C. Sankaranarayanan, "LiSens - A scalable architecture for video compressive sensing," in *Proc. IEEE Conf. Computational Photography*, April 2015, pp. 1–9.
- [6] H. Chen, M. S. Asif, A. C. Sankaranarayanan, and A. Veeraraghavan, "FPA-CS: Focal plane array-based compressive imaging in short-wave infrared," in *Proc. IEEE Conf. Comp. Vision Pattern Recognition*, June 2015, pp. 2358–2366.
- [7] E. Tsiligianni, L. P. Kondi, and A. K. Katsaggelos, "Preconditioning for underdetermined linear systems with sparse solutions," *IEEE Signal Process. Lett.*, vol. 22, no. 9, pp. 1239–1243, Sept. 2015.
- [8] R. Baraniuk, M. Davenport, R. DeVore, and M. Wakin, "A simple proof of the restricted isometry property for random matrices," *Constructive Approximation*, vol. 28, no. 3, pp. 253–263, 12 2008.
- [9] M. Elad, "Optimized projections for compressed sensing," *IEEE Trans. Signal Process.*, vol. 55, no. 12, pp. 5695–5702, Dec. 2007.
- [10] E. V. Tsiligianni, L. P. Kondi, and A. K. Katsaggelos, "Construction of incoherent unit norm tight frames with application to compressed sensing," *IEEE Trans. Inf. Theory*, vol. 60, no. 4, pp. 2319–2330, April 2014.
- [11] J. Xu, Y. Pi, and Z. Cao, "Optimized projection matrix for compressive sensing," *EURASIP J. Adv. Signal Process.*, vol. 2010, pp. 43:1–43:8, Feb. 2010.
- [12] Y. LeCun, L. Bottou, G. B. Orr, and K.-R. Müller, "Efficient backprop," in *Neural Networks: Tricks of the Trade, This Book is an Outgrowth of a 1996 NIPS Workshop*. London, UK, UK: Springer-Verlag, 1998, pp. 9–50. [Online]. Available: <http://dl.acm.org/citation.cfm?id=645754.668382>
- [13] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, May 2015, insight.
- [14] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *CoRR*, vol. abs/1512.03385, 2015.
- [15] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 6, pp. 1137–1149, June 2017.
- [16] D. Reddy, A. Veeraraghavan, and R. Chellappa, "P2C2: Programmable pixel compressive camera for high speed imaging," in *Proc. IEEE Conf. Comp. Vision Pattern Recognition*, June 2011, pp. 329–336.
- [17] P. Llull, X. Liao, X. Yuan, J. Yang, D. Kittle, L. Carin, G. Sapiro, and D. J. Brady, "Coded aperture compressive temporal imaging," *Opt. Express*, vol. 21, no. 9, pp. 10 526–10 545, May 2013.
- [18] L. Gao, J. Liang, C. Li, and L. V. Wang, "Single-Shot compressed ultrafast photography at one hundred billion frames per second," *Nature*, vol. 516, pp. 74–77, 2014.
- [19] D. Liu, J. Gu, Y. Hitomi, M. Gupta, T. Mitsunaga, and S. K. Nayar, "Efficient space-time sampling with pixel-wise coded exposure for high-speed imaging," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 36, no. 2, pp. 248–260, Feb. 2014.
- [20] R. Koller, L. Schmid, N. Matsuda, T. Niederberger, L. Spinoulas, O. Cossairt, G. Schuster, and A. K. Katsaggelos, "High spatio-temporal resolution video with compressed sensing," *Opt. Express*, vol. 23, no. 12, pp. 15 992–16 007, June 2015.
- [21] A. Chakrabarti, "Learning sensor multiplexing design through back-propagation," in *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, 2016, pp. 3081–3089. [Online]. Available: <http://papers.nips.cc/paper/6251-learning-sensor-multiplexing-design-through-back-propagation>
- [22] C. J. Schuler, H. C. Burger, S. Harmeling, and B. Scholkopf, "A machine learning approach for non-blind image deconvolution," in *Proc. IEEE Conf. Comp. Vision Pattern Recognition*, June 2013, pp. 1067–1074.
- [23] J. Sun, W. Cao, Z. Xu, and J. Ponce, "Learning a convolutional neural network for non-uniform motion blur removal," in *Proc. IEEE Conf. Comp. Vision Pattern Recognition*, June 2015, pp. 769–777.
- [24] L. Xu, J. S. Ren, C. Liu, and J. Jia, "Deep convolutional neural network for image deconvolution," in *Adv. Neural Inf. Process. Syst. 27*. Curran Associates, Inc., 2014, pp. 1790–1798.
- [25] H. C. Burger, C. J. Schuler, and S. Harmeling, "Image denoising: Can plain neural networks compete with BM3D?" in *Proc. IEEE Conf. Comp. Vision Pattern Recognition*, June 2012, pp. 2392–2399.
- [26] F. Agostinelli, M. R. Anderson, and H. Lee, "Adaptive multi-column deep neural networks with application to robust image denoising," in *Adv. Neural Inf. Process. Syst. 26*. Curran Associates, Inc., 2013, pp. 1493–1501.
- [27] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, and P.-A. Manzagol, "Stacked Denoising Autoencoders: Learning useful representations in a deep network with a local denoising criterion," *J. Mach. Learn. Res.*, vol. 11, pp. 3371–3408, Dec. 2010.
- [28] J. R. Chang, C. Li, B. Póczos, B. V. K. V. Kumar, and A. C. Sankaranarayanan, "One network to solve them all - solving linear inverse problems using deep projection models," *CoRR*, vol. abs/1703.09912, 2017. [Online]. Available: <http://arxiv.org/abs/1703.09912>
- [29] C. Ledig, L. Theis, F. Huszar, J. Caballero, A. Aitken, A. Tejani, J. Totz, Z. Wang, and W. Shi, "Photo-realistic single image super-resolution using a generative adversarial network," *CVPR*, 2017.
- [30] D. Pathak, P. Krhenbhl, J. Donahue, T. Darrell, and A. A. Efros, "Context Encoders: Feature learning by inpainting," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016, pp. 2536–2544.
- [31] Z. Cui, H. Chang, S. Shan, B. Zhong, and X. Chen, *Computer Vision – ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V*. Cham: Springer International Publishing, 2014, ch. Deep Network Cascade for Image Super-resolution, pp. 49–64.
- [32] C. Dong, C. Loy, K. He, and X. Tang, "Image super-resolution using deep convolutional networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 2, pp. 295–307, Feb. 2016.
- [33] Y. Huang, W. Wang, and L. Wang, "Bidirectional recurrent convolutional networks for multi-frame super-resolution," in *Adv. Neural Inf. Process. Syst. 28*. Curran Associates, Inc., 2015, pp. 235–243.
- [34] A. Kappeler, S. Yoo, Q. Dai, and A. K. Katsaggelos, "Video super-resolution with convolutional neural networks," *IEEE Transactions on Computational Imaging*, vol. 2, no. 2, pp. 109–122, June 2016.
- [35] A. Mousavi, A. B. Patel, and R. G. Baraniuk, "A deep learning approach to structured signal recovery," *CoRR*, vol. abs/1508.04065, 2015.
- [36] K. Kulkarni, S. Lohit, P. K. Turaga, R. Kerviche, and A. Ashok, "Reconnet: Non-iterative reconstruction of images from compressively sensed random measurements," *CoRR*, vol. abs/1601.06892, 2016.
- [37] H. Yao, F. Dai, D. Zhang, Y. Ma, S. Zhang, and Y. Zhang, "Dr²-net: Deep residual reconstruction network for image compressive sensing," *CoRR*, vol. abs/1702.05743, 2017. [Online]. Available: <http://arxiv.org/abs/1702.05743>
- [38] M. Iliadis, L. Spinoulas, and A. K. Katsaggelos, "Deep fully-connected networks for video compressive sensing," *CoRR*, vol. abs/1603.04930, 2016.
- [39] M. Courbariaux and Y. Bengio, "BinaryNet: Training deep neural networks with weights and activations constrained to +1 or -1," *CoRR*, vol. abs/1602.02830, 2016.
- [40] M. Courbariaux, Y. Bengio, and J.-P. David, "BinaryConnect: Training deep neural networks with binary weights during propagations," in *Adv. Neural Inf. Process. Syst. 28*. Curran Associates, Inc., 2015, pp. 3123–3131.
- [41] Z. Lin, M. Courbariaux, R. Memisevic, and Y. Bengio, "Neural networks with few multiplications," *CoRR*, vol. abs/1510.03009, 2015.
- [42] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadji, "XNOR-Net: ImageNet classification using binary convolutional neural networks," *CoRR*, vol. abs/1603.05279, 2016.
- [43] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *Proc. Int. Conf. Machine Learning*, 2010, pp. 807–814.
- [44] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," in *Proc. Int. Conf. Artificial Intelligence and Statistics*, vol. 9, May 2010, pp. 249–256.
- [45] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Neurocomputing: Foundations of research." Cambridge, MA, USA: MIT Press, 1988, ch. Learning Representations by Back-propagating Errors, pp. 696–699.
- [46] R. Pascanu, T. Mikolov, and Y. Bengio, "On the difficulty of training recurrent neural networks," in *ICML (3)*, ser. JMLR Proceedings, vol. 28. JMLR.org, 2013, pp. 1310–1318.
- [47] J. Yang, X. Liao, X. Yuan, P. Llull, D. J. Brady, G. Sapiro, and L. Carin, "Compressive sensing by learning a gaussian mixture model from measurements," *IEEE Trans. Image Processing*, vol. 24, no. 1, pp. 106–119, Jan. 2015.
- [48] J. Yang, X. Yuan, X. Liao, P. Llull, D. J. Brady, G. Sapiro, and L. Carin, "Video compressive sensing using gaussian mixture models," *IEEE Trans. Image Processing*, vol. 23, no. 11, pp. 4863–4878, Nov. 2014.
- [49] J. M. Bioucas-Dias and M. A. T. Figueiredo, "A New TwIST: Two-step iterative shrinkage/thresholding algorithms for image restoration," *IEEE Trans. Image Process.*, vol. 16, no. 12, pp. 2992–3004, Dec. 2007.
- [50] R. Collobert, K. Kavukcuoglu, and C. Farabet, "Torch7: A matlab-like environment for machine learning," in *BigLearn, NIPS Workshop*, 2011.