# Solving Inverse Computational Imaging Problems using Deep Pixel-level Prior

Akshat Dave, Anil Kumar Vadathya, Ramana Subramanyam, Rahul Baburajan, Kaushik Mitra, *Member, IEEE*

*Abstract*—Signal reconstruction is a challenging aspect of computational imaging as it often involves solving ill-posed inverse problems. Recently, deep feed-forward neural networks have led to state-of-the-art results in solving various inverse imaging problems. However, being task specific, these networks have to be learned for each inverse problem. On the other hand, a more flexible approach would be to learn a deep generative model once and then use it as a signal prior for solving various inverse problems. We show that among the various state of the art deep generative models, autoregressive models are especially suitable for our purpose for the following reasons. First, they explicitly model the pixel level dependencies and hence are capable of reconstructing low-level details such as texture patterns and edges better. Second, they provide an explicit expression for the image prior which can then be used for MAP based inference along with the forward model. Third, they can model long range dependencies in images which make them ideal for handling global multiplexing as encountered in various compressive imaging systems. We demonstrate the efficacy of our proposed approach in solving three computational imaging problems: Single Pixel Camera (SPC), LiSens and FlatCam. For both real and simulated cases, we obtain better reconstructions than the state-of-the-art methods in terms of perceptual and quantitative metrics.

*Index Terms*—Inverse problems, compressive image recovery, deep generative models, lensless image reconstruction, autoregressive models, MAP inference.

## I. INTRODUCTION

COMPUTATIONAL imaging systems enable us to extract much more information out of the visual world as compared to the traditional imaging systems. This is achieved by jointly designing optics, to encode the desired signal information, and algorithms to reconstruct the signal back from those measurements. Signal reconstruction corresponds to inverting the forward model used in acquiring the measurements. Hence, reconstruction algorithms for different computational imaging devices amount to solving different inverse problems. Solving these inverse problems becomes challenging as they are often ill-posed. For compressive imaging setups such as Single Pixel Camera (SPC) [2], [3], high speed imaging [4], [5] and compressive hyper-spectral imaging [6], the reconstruction becomes ill-posed as the number of measurements is quite less than the signal dimension.

Akshat Dave is with the Department of Electrical and Computer Engineering, Rice University, Houston, TX, USA (e-mail: akshat.dave@rice.edu).

Anil Kumar Vadathya, Rahul Baburajan and Kaushik Mitra are with the Computational Imaging Lab, Indian Institute of Technology (IIT) Madras, Chennai, India.

Ramana Subramanyam is with Insight Centre for Data Analytics, Dublin, Ireland.

Generally, for solving an ill-posed problems, we need to incorporate the prior information about the signal to be reconstructed. Traditionally these priors are either analytically derived or hand-crafted based on the observations. For example, sparsity of image gradients [7], sparsity of coefficients in wavelet and DCT domain [8] etc. have been used for solving inverse imaging problems. However, the underlying data distribution may not precisely follow these analytic priors leading to poor solutions in challenging scenarios. Dictionary learning [9] methods being data driven are an improvement over these analytic priors. However, being limited by patch size they cannot account for long range dependencies which are necessary for handling global multiplexing in case of compressive image reconstruction.

On the other hand, deep learning based reconstruction algorithms recently have led to state-of-the-art results in solving such ill-posed problems in computational imaging [10], [11] [12] [13]. These approaches typically learn an inverse mapping from measurements to the signal by minimizing reconstruction loss on a set of training examples. However, this kind of training, popularly known as discriminative learning, makes the network task specific. Furthermore, we need to retrain the network for various parameter settings of the forward model. For example, for every new setting of measurement rate and sensing matrix in SPC, we need to relearn the network parameters. Instead of having to design/retrain a different network for each task and parameter setting, it would be more efficient to have a generalized framework which can be used for solving various inverse problems.

A more flexible approach would be to learn the natural image statistics using a generative model and use it for solving various inverse problems. Recently, deep generative models especially using *autoregressive framework* [14], [15], [16] have led to state-of-the-art performance in modeling natural image manifold. Autoregressive models factorize the image distribution as a 2D directed causal graph and hence model it as a 2-D sequence where current pixel's distribution is conditioned on the causal context. By employing deep neural networks for summarizing the causal context, autoregressive models excel at capturing long range dependencies in images. Also, being a pixel level model it explicitly accounts for higher order correlations like texture patterns, sharp edges, etc. within a neighbourhood. Thus, these models are capable of generating visually convincing and crisp images [14]. Examples of deep autoregressive image models are recurrent image density estimator (RIDE) [15], pixel recurrent neural networks (PixelRNN) and its CNN equivalent (PixelCNN) [14] and PixelCNN++ [16].
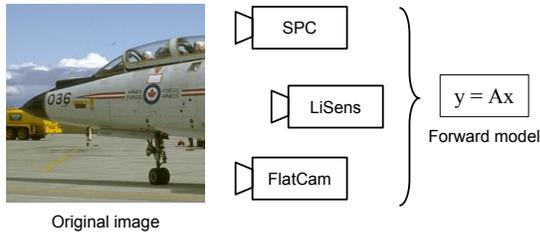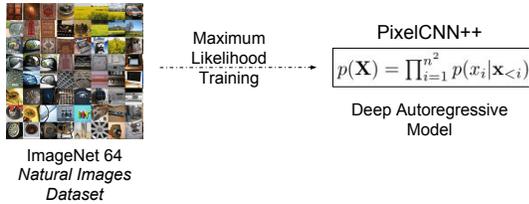
Computational imaging setups:
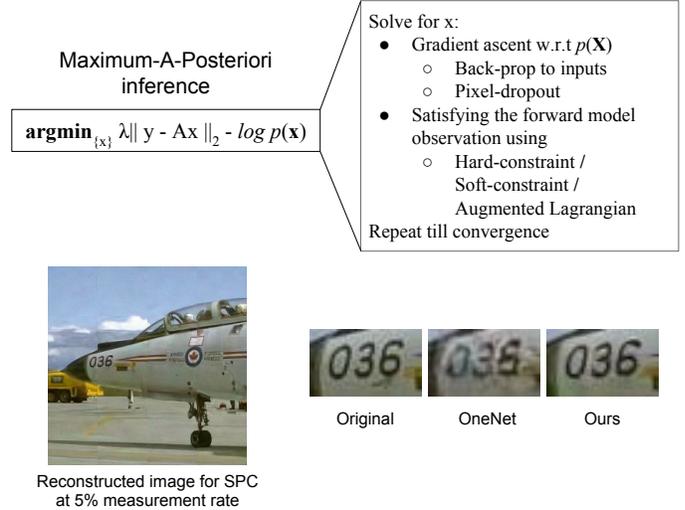
Image reconstruction:



Fig. 1. An overview of our approach. We employ a single deep autoregressive model learned on natural images for solving multiple inverse problems. From the zoomed in patch of the reconstructed image in the inset it is evident that our approach has better pixel-level consistencies as compared to existing latent representation based models like OneNet [1].

We show that deep autoregressive generative models are ideally suitable for solving various computational imaging problems for the following reasons. First, it explicitly models the distribution of each pixel in relation to its causal neighbor. Thus, when used as an image prior, this explicit pixel dependency modeling helps it to better reconstruct low level details without artifacts (see Figure 1). Second, this framework gives us an explicit expression for the image prior, which can be used for doing MAP inference. Moreover, the entire framework is differentiable, which is amenable for gradient based inference. Third, its ability to capture long range dependencies in images makes them ideal for handling global multiplexing in compressive imaging setups. Given these advantages with deep autoregressive models, we use it for solving various computational imaging problems such as - Single Pixel Camera (SPC) [2], Line Sensor (LiSens) [3] and lensless imaging - FlatCam [17]. Our results demonstrate that we perform better than the current state-of-the-art methods in both traditional and learning based approaches.

In summary we make the following contributions:

- We propose a versatile approach which employs the same learned prior model for solving various computational imaging problems.
- We propose to use a deep autoregressive model, Pixel-CNN++, as an image prior. The autoregressive nature of this prior ensures pixel-level consistencies in the reconstruction and hence provides better quality than using latent representation based models such as OneNet [1] as shown in Figure 1.
- We utilize back-propagation to the inputs for obtaining tractable estimates of the prior gradients and employ them for solving inverse problems using MAP inference.
- We observe that randomly dropping the gradient updates for a certain percentage of pixels at every iteration

helps in reconstructing the texture better. We analyze the effect of this pixel dropout ratio on the quality of reconstructions.
- We demonstrate better reconstructions than the existing state-of-the-art methods for three computational imaging problems: Single Pixel Camera, LiSens, and FlatCam.

## II. RELATED WORK

**Compressive imaging** Single Pixel Camera (SPC) [2] is a classic example of compressive imaging. It uses a programmable digital micro-mirror device (DMD) array to multiplex the scene on to a single photodetector. Using different settings on the DMD, we can sequentially acquire a set of measurements. Thus, scene at full resolution is reconstructed from much less than 100% measurements. Compressive imaging systems pose a viable solution for high resolution imaging in non-visible parts of the spectrum where full frame sensors are very expensive.

The measurement bandwidth of the SPC is limited by the operating speed of the DMDs (Tens of kHz for commercially-available units). With this speed, SPC cannot be extended for high resolution video sensing. On one end, we have exorbitant full frame sensors (Nyquist sampling) for high resolution imaging in non visible bands, and on the other, we have SPC, an inexpensive compressive sensing setup but with low measurement rates. Wang et al. [3] propose LiSens - Line Sensor based compressive camera which lies midway between these two imaging extremes. Each pixel in the line sensor is mapped to a row in DMD array. Thus, unlike SPC, where the whole scene is multiplexed, here only rows of the scene are multiplexed.

**Lensless imaging** FlatCam [17] and DiffuserCam [18] are novel imaging systems which get rid of the conventional lens optics. Instead, they use amplitude and diffuser mask

respectively to encode light coming from different parts of the scene onto the sensor. As a result, information localized at a point in the scene gets spread throughout the sensor, making priors essential for accurate recovery of the image. These works use traditional reconstruction algorithms such as Total Variation norm and Tikhonov regularization which are quick but do not provide natural looking reconstructions.

**Reconstruction with analytical priors** Many algorithms have been proposed for compressive image reconstruction. Typically, reconstruction algorithms use $l_1$ regularization, exploiting the sparsity of spatial gradients in natural images. Total Variation (TV) minimization prior [7], [19] is the most commonly used reconstruction algorithm based on this sparsity. Chengbo et al. [20] propose an efficient augmented Lagrangian based TV minimization for CS reconstruction. Recent approaches involving compressive architectures such as fpa-cs [21], LiSens [3], and video CS [22], demonstrated successful results with TV minimization prior. However, at lower measurement rates, reconstructions suffer from the piece-wise smooth modeling of TV prior and results tend to be blocky, as is noted by recent works [10], [23]. Metzler et al. [24] propose a denoiser based CS reconstruction algorithm. Specifically, use a Gaussian denoiser with approximate message passing algorithm (D-AMP). At very low measurement rates, the denoiser tends to result in overly smooth images as is recently shown by Dave et al. [23], Kulkarni et al. [10].

**Data driven CS reconstruction** Duarte et al. [25] propose an approach for simultaneous learning of the sensing matrix and dictionary atoms. Due to the small patch size of the atoms, their usage for compressive image reconstruction is limited to local multiplexing, unlike the actual SPC involving global multiplexing of the scene. Reconstruction algorithms using convolutional neural networks (CNNs) typical take input as measurements from an image patch and try to output the image back by minimizing the reconstruction loss. Kulkarni et al. [10] proposed ReconNet, Yao et al. [11] proposed $DR^2$-Net having residual connections for reconstruction. Although these approaches lead to a non-iterative and hence faster inference, being task specific, they only work for the fixed settings of the sensing matrix and measurement rates used for training. Changing the settings requires retraining the architecture which is not very appealing. Also, being patch-wise, they also fail to account for global multiplexing in SPC.

**Deep generative models** With the success of deep neural networks, there have been multiple works proposing deep generative models, which explicitly or implicitly try to model the distribution of natural images. For example, latent representation models like adversarial networks, GAN by Goodfellow et al. [26], variational auto-encoders by Kingma et al. [27] and autoregressive models like RIDE by Theis et al. [15], PixelRNN/CNN by Oord et al. [14], PixelCNN++ by Salimans et al. [16]. GANs learn to transform samples from a Gaussian distribution to a sample in the natural image manifold via a generator network, which is trained with an adversarial learning framework involving a discriminator network. VAEs are a probabilistic framework of autoencoders that learn to encode and decode the images from a distribution.

Autoregressive models factorize an image as a 2D directed graph by conditioning the current pixel $x_i$'s distribution on the pixels before it as in a raster scan $x_{<i}$. Modeling this conditional density is analogous to sequence modeling and initial methods proposed to use spatial 2D recurrent neural networks, given their efficacy in modeling sequences. RIDE by Theis et al. [15] uses 2D Long Short Term Memory (LSTM) units called Spatial-LSTMs for modeling the causal context $x_{<i}$, and GSMs for parametrizing the distribution. PixelRNN by Oord et al. [14] uses a much complex architecture using LSTMs and residual connections to better handle the causal context. Importantly, it models the conditional density as a discrete distribution with $x_i \in \{0, 1, \dots 255\}$. PixelRNN has resulted in state-of-the-art negative loglikelihood (NLL) scores. However, due to the sequential nature of distribution modeling, both training and sampling are computationally demanding with the runtime as $O(N)$, where $N$ is the total number of pixels. Oord et al. proposed PixelCNN which is a convolutional version of PixelRNN. This led to an improvement in the training time by a large factor at the cost of slight loss in the accuracy as with convolutions we can now only capture bounded context. Salimans et al. [16] proposed PixelCNN++, which builds on PixelCNN by employing a discretized mixture of logistics for modeling the distribution, and using drop-out regularization, and additional skip connections. It improves on the NLL score over PixelRNN on the CIFAR dataset leading to state-of-the-art results.

**Deep image priors** When solving linear inverse problems using the alternating direction method of multipliers (ADMM) algorithm, Venkatakrishnan et al. [28] observed that it results in two decoupled optimizations. The first one enforces the data prior while the second enforces data fidelity to the observation. The first step can be thought of as a denoising problem, thus, a denoiser can be employed to solve this step thereby avoiding the need for an explicit image prior. Venkatakrishnan et al. [28] use denoisers like BM3D [29] in ADMM setting for image restoration. Inspired by this, recent methods propose learning-based proximal operators for the denoising step of ADMM. OneNet by Chang et al. [1], CNN denoiser by Zhang et al.[30], Meinhardt et al. [31] . In this work, we compare our explicit natural image prior based MAP inference with the learned proximal operator of OneNet. Our evaluations show that our results are superior to OneNet. It is important to note that OneNet's proximal operator uses adversarial loss [26] which is known to result in sharper recovery of details.

In this paper, we extend upon our previous work, Dave et al. [23] (RIDE-CS), where we used recurrent image density estimator (RIDE) for CS reconstruction. We observed that the sequential nature of recurrent networks in RIDE makes it too slow for inference and training (computational cost is proportional to the image size). Also, in our experiments, the two layer RIDE fails to yield results comparable to recent approaches like OneNet [1]. Here, we explore sophisticated deep autoregressive models which are order faster than RIDE-CS for both training and inference. We apply the deep autoregressive model based inference to recent frameworks in computational imaging like LiSens [3] and FlatCam [17]. We enhance the inference algorithm by incorporating the augmented Lagrangian method when necessary. In addition, we

improve texture recovery using pixel-wise stochastic gradient updates.

## III. Inference with deep autoregressive models

### A. Problem Formulation

Consider $\mathbf{X}$ to be a $n \times n$ matrix corresponding to a natural image and $f$ to be a linear transformation corresponding to the forward model of a computational camera. The measurements obtained $\mathbf{Y}$ can be written as $\mathbf{Y} = f(\mathbf{X})$. Our goal is to reconstruct back the image $\mathbf{X}$ from the measurements $\mathbf{Y}$.

Discriminative networks learn the inverse mapping $\hat{\mathbf{X}} = g(\mathbf{Y})$ by modelling $g$ as a deep neural network and minimizing the reconstruction error on a set of training examples $\{\mathbf{X}_i, \mathbf{Y}_i\}$. Hence, the inverse mapping is implicitly dependant on the forward model $f$. Dealing with reconstructions for multiple forward models would require learning separate networks for each model which can be expensive.

For our generative approach, we model the distribution of natural images $p(\mathbf{X})$ using a deep autoregressive model. We formulate the inverse problem as MAP inference. Hence, the estimated image $\hat{\mathbf{X}}$ can be written as

$$\hat{\mathbf{X}} = \arg\max_{\mathbf{X}} log(p(\mathbf{X}|\mathbf{Y})) \tag{1}$$

$$= \arg\max_{\mathbf{X}} (log(p(\mathbf{Y}|\mathbf{X})) + log(p(\mathbf{X}))) \tag{2}$$

The likelihood term $p(\mathbf{Y}|\mathbf{X})$ varies for different imaging systems based on the forward model but the image prior $p(\mathbf{X})$ remains the same. Thus, we need to learn the prior only once for all the problems.

### B. Forward Models

Let the $n^2 \times 1$ column vector $\mathbf{x}$ represent the rasterized version of the $n \times n$ image matrix $\mathbf{X}$ i.e. $\mathbf{x} \triangleq vec(\mathbf{X})$ by taking pixels row by row. The forward models that we consider in this work are as follows:

*1) Randomly Missing Pixels:* Here, we randomly set certain number of pixels in an image to by missing, by setting their values to zero. Hence, $\mathbf{y}$ i.e. the vectorized version of the resultant image can be written as

$$\mathbf{y} = \mathbf{m} \circ \mathbf{x} \tag{3}$$

where $\circ$ denotes the Hadamard product and $\mathbf{m}$ is a Bernoulli random vector. The above equation can also be expressed in a matrix-vector multiplication form as :

$$\mathbf{y} = \mathbf{M}\mathbf{x} \tag{4}$$

where $\mathbf{M}$ is a sub-sampling matrix.

*2) Single Pixel Camera:* In SPC [2], the DMD array optically multiplexes the scene onto a single pixel sensor. By changing the orientation of the array, we will get different multiplexing patterns, which results in different measurements. If $\mathbf{y}$ is the vector of $m$ single pixel measurements from SPC and $\Phi$ is the $m \times n^2$ compressive sensing matrix, then we have the forward model as:

$$\mathbf{y} = \Phi\mathbf{x}. \tag{5}$$

*3) LiSens:* In Lisens [3], the 2D image of the scene formed on the DMD plane is mapped onto a 1D line-sensor which essentially captures the 1D integral of the 2D image (along rows or columns). If $Y$ is the $m \times n$ matrix formed by stacking $m$ line sensor measurements from Lisens and $\Phi$ is the $m \times n$ sensing matrix, then we have

$$\mathbf{Y} = \Phi\mathbf{X} \tag{6}$$

*4) FlatCam:* FlatCam [17] replaces the lens system by a coded amplitude mask close to the sensor. For ease of calibration, this mask is designed to be separable, i.e., it can be written as an outer product of 2 one dimensional patterns. Neglecting the diffraction effects, it was shown in [17] that using such a mask, the $m \times m$ measurements $\mathbf{Y}$ obtained on the FlatCam sensor can be written as

$$\mathbf{Y} = \Phi_L \mathbf{X} \Phi_R^T \tag{7}$$

where $\Phi_L$ and $\Phi_R$ are $m \times n$ matrices corresponding to 1-D convolution of the scene $\mathbf{X}$ along the rows and columns respectively.

### C. Deep autoregressive model

Here we model the dependencies between pixels using a directed probabilistic chain. The pixel $x_i$ depends on all the pixels before the index $i$ in $\mathbf{x}$, which we denote as $\mathbf{x}_{<i}$. Hence the joint distribution over the pixels in the image can be factorized as

$$p(\mathbf{X}) = p(x_1, x_2, \ldots, x_{n^2}) = \prod_{i=1}^{n^2} p(x_i|\mathbf{x}_{<i}) \tag{8}$$

In this work, we use state-of-the-art autoregressive generative model, PixelCNN++ [16]. Here, the context $\mathbf{x}_{<i}$ for the conditional distribution of each of the pixels is modelled using a deep convolutional neural network with residual connections. The convolution kernels are masked appropriately to ensure that the context of a pixel does not depend on the pixels after it. The conditional distribution is then modelled as a mixture of logistic distributions, where the parameters of the distribution depend on the context. This model is then learned on RGB images using maximum likelihood training.

Once the model is trained, it can be used to solve different inference tasks, as we describe below. Sampling from autoregressive models is slow because of their sequential nature which limits their utility. However, for our approach, we only require the gradients of the density $p(\mathbf{X})$ with respect to the image $X$. This can be computed efficiently using backpropagation to the inputs.

## IV. Optimization methods for deep autoregressive inference

In this section, we discuss inference methods for various forward models discussed earlier. We want the desired solution to have higher likelihood (lower NLL) under the image prior and at the same time satisfy the constraints specified by the forward model. For this, we perform projected gradient descent. We divide our approach into three categories based

on the amount of noise and the kind of forward model. Hard constraint (equality) method is used when there is less or no measurement noise (Section IV-A). For certain imaging models like FlatCam, there is no closed form for the projection operator. We instead use the Augmented Lagrangian Method (ALM), see Section IV-B. For the cases of high noise, the measurements deviate significantly from the forward model, and the soft constraint method (inequality) is used (Section IV-C). Further, in Sections IV-D and IV-E, we describe two implementation hacks which have proved useful for our approach.

### A. Hard constraint method

We first analyze the case when the measurement is directly obtained using the imaging model without any noise. $\mathbf{Y}$ is then a deterministic function of $\mathbf{X}$ and hence the likelihood term would correspond to constraints. The problem can be formulated as

$$\hat{\mathbf{X}} = \arg\max_{\mathbf{X}}(log(p(\mathbf{X})) \text{ such that } \mathbf{Y} = f(\mathbf{X}) \quad (9)$$

where $f$ is provided by the imaging model. The signal prior model is the learned autoregressive model with parameters $\theta$. Also, we constrain the intensity of the image to be between $0$ and $1$. Thus our problem is given by:

$$\hat{\mathbf{X}} = \arg\max_{\mathbf{X}}(log(p_\theta(\mathbf{X})) \text{ s.t. } \mathbf{Y} = f(\mathbf{X}), 0 \leq \mathbf{X}_{ij} \leq 1 \quad (10)$$

Let $\mathcal{C}_1$ and $\mathcal{C}_2$ denote the constraint sets $\{\mathbf{X} : \mathbf{Y} = f(\mathbf{X})\}$ and $\{\mathbf{X} : 0 \leq \mathbf{X}_{ij} \leq 1 \quad \forall i,j\}$ respectively.

We use projected gradient descent to solve this constrained optimization, which involves performing the following steps iteratively:

$$\mathbf{H}_k = \mathbf{X}_k + \alpha\nabla_{\mathbf{X}}log(p_\theta(\mathbf{X}_k)) \quad (11)$$

$$\mathbf{J}_k = \Pi_{\mathcal{C}_1}(\mathbf{H}_k) \quad (12)$$

$$\mathbf{X}_{k+1} = \Pi_{\mathcal{C}_2}(\mathbf{J}_k) \quad (13)$$

where $\Pi_{\mathcal{C}_1}$ and $\Pi_{\mathcal{C}_2}$ are projection operators to the constraint sets $\mathcal{C}_1$ and $\mathcal{C}_2$ respectively. For Eq. 11 backpropagation to the inputs is used to get the data gradients. For Eq. 13, pixels in the image are clipped between $0$ and $1$ in every iteration.

$\Pi_{\mathcal{C}_1}$ is different for different imaging problems. For the randomly missing pixels case,

$$\mathbf{j_k} = (\mathbf{1} - \mathbf{m}) \circ \mathbf{h_k} + (\mathbf{m}) \circ \mathbf{y} \quad (14)$$

where $\mathbf{1}$ is an $n^2$ vector of ones. This implies that we should only be updating the missing pixels and leave the other pixels the same, which is intuitive.

For Single Pixel Camera we have,

$$\mathbf{j_k} = \mathbf{h}_k - \Phi^T \left(\Phi\Phi^T\right)^{-1} \left(\Phi\mathbf{h}_k - \mathbf{y}\right) \quad (15)$$

where $\mathbf{j_k}$ and $\mathbf{h}_k$ are vector representations of matrices $\mathbf{J}_k$ and $\mathbf{H}_k$ respectively. We consider row-orthonormalized matrices for compressive sensing, hence $\Phi\Phi^T$ is an identity matrix.

For LiSens case, similar to SPC, we have

$$\mathbf{J_k} = \mathbf{H}_k - \Phi^T \left(\Phi\Phi^T\right)^{-1} \left(\Phi\mathbf{H}_k - \mathbf{Y}\right). \quad (16)$$

### B. Augmented Lagrangian method

For the case of FlatCam reconstruction, the matrices $\mathbf{\Phi_L\Phi_L}^T$ and $\mathbf{\Phi_R\Phi_R}^T$ are ill-conditioned and can't be inverted. A closed form solution for projection operator doesn't exist. So, we consider the augmented Lagrangian corresponding to $\mathcal{C}_1$, with a dual parameter $\lambda$.

$$\mathcal{L}(\mathbf{X}, \boldsymbol{\lambda}) = -log(p_\theta(\mathbf{X})) + \rho\|\mathbf{Y} - \Phi_L\mathbf{X}\Phi_R^T\|_F^2 \\ + \langle \boldsymbol{\lambda}, \mathbf{Y} - \Phi_L\mathbf{X}\Phi_R^T \rangle_F \quad (17)$$

However, instead of minimizing the Lagrangian with respect to the primal variable in each iteration, we just take one step of gradient descent. We further separate the gradient descent into two steps, one entirely depends on the prior while the other entirely depends on the imaging model. The update steps are as follows.

$$\mathbf{H}_k = \mathbf{X}_k + \alpha\nabla_{\mathbf{X}}log(p_\theta(\mathbf{X}_k) \quad (18)$$

$$\mathbf{J}_k = \mathbf{H}_k + \Phi_L^T(\boldsymbol{\lambda}_k - \rho(\mathbf{Y} - \Phi_L\mathbf{X}_k\Phi_R^T))\Phi_R \quad (19)$$

$$\mathbf{X}_{k+1} = \Pi_{\mathcal{C}_2}(\mathbf{J}_k) \quad (20)$$

$$\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \rho(\mathbf{Y} - \Phi_L\mathbf{X}_k\Phi_R^T) \quad (21)$$

### C. Soft constraint method

Consider the case when the sensor has measurement noise,

$$\mathbf{Y} = f(\mathbf{X}) + \boldsymbol{\eta} \quad (22)$$

Assume the measurement noise $\boldsymbol{\eta}$ to be Gaussian distributed, i.e.

$$\boldsymbol{\eta} \sim \mathcal{N}(0, \sigma) \quad (23)$$

$$\mathbf{Y} \sim \mathcal{N}(f(\mathbf{X}), \sigma) \quad (24)$$

The MAP estimation problem can hence be reduced to

$$\hat{\mathbf{X}} = \arg\max_{\mathbf{X}}(log(p_\theta(\mathbf{X})) + \lambda\|\mathbf{Y} - f(\mathbf{X})\|^2) \quad (25)$$

where $\lambda$ has to be estimated if we do not know the standard deviation of the measurement noise. Since the constraints are not exact here, we replace the step to project to the constraint space by instead taking a step towards minimizing the likelihood. Hence, we replace Eq. 12 by gradient descent over likelihood,

$$\mathbf{J}_k = \mathbf{H}_k - \alpha f'(\mathbf{H}_k)(\mathbf{Y} - f(\mathbf{H}_k)) \quad (26)$$

### D. Stochastic gradients using pixel dropout

We observe that if we update all the pixels in the gradient update (Eq. 11), then we get washed out reconstructions. The autoregressive prior directly models correlation between neighbouring pixels. Hence it tends to assign same values to neighbouring problems. We combat this problem by randomly selecting a certain amount of pixels to update in each step. Hence, not all pixels get updates at every step. We call this pixel dropout, and for incorporating that, we replace the gradient in Eq. 11 by stochastic gradients, i.e.,

$$\mathbf{H}_k = \mathbf{X}_k + \alpha\mathbf{M} \circ \nabla_{\mathbf{X}}log(p_\theta(\mathbf{X}_k)) \quad (27)$$

where $\mathbf{M}$ is a random binary mask with the percentage of zeros determined by the pixel dropout ratio. This is analogous to the case of training deep neural networks, where Stochastic Gradient Descent (SGD) helps in escaping from sharp local minima [32]. Here, the washed out reconstructions correspond to sharp local minima owing to the strong correlation between pixels. We demonstrate the effect of the amount of pixel dropout on the reconstructions in Section VI-D1.

### E. Splitting and Stitching

Our prior model is trained on $64 \times 64$ patches, hence the input for $p_\theta(\mathbf{X})$ has to be $64 \times 64$. While we perform the likelihood step on the entire image, our approach is designed such that the prior gradient update, projection, and clipping steps are separate. Before the prior gradient update, we split the image into a batch of $64 \times 64$ patches. Before performing the likelihood step, we stitch the patches back into original dimensions.

Our approach is summarized as follows:

---

**Algorithm 1:** Our image reconstruction algorithm

**Data:** Simulated or real measurements $\mathbf{Y}$, Simulated or calibrated imaging matrix $\Phi$, Learned autoregressive prior model $p_\theta(\mathbf{X})$

**Result:** Reconstructed $N \times N$ image $\mathbf{X}$

**1** Initialization: $\mathbf{X}_{ij} \sim \mathcal{U}(0,1) \; \forall$ pixels $i, j$ **while** *iterations < max_iter* **do**

**2**    Split $\mathbf{X}$ into a batch of $64 \times 64$ patches **for** *Gradient ascent w.r.t $p_\theta(\mathbf{X})$* **do**

**3**      Obtain $\nabla_{\mathbf{X}} p_\theta(\mathbf{X})$ via back-prop to inputs

**4**      Apply pixel dropout mask and update $\mathbf{X}$ (Eq. 27)

**5**    **end**

**6**    Stitch $\mathbf{X}$ back into $N \times N$ image

**7**    **for** *Satisfying constraints* **do**

**8**      Clip $\mathbf{X}_{ij}$ between 0 and 1 $\forall \; i, j$

**9**      Project the solution to the constraint space specified by the forward model and inference method appropriately (Eq. 14,15,16,19 or 25)

**10**    **end**

**11**    **if** *method is augmented Lagrangian* **then**

**12**      Update dual variable $\boldsymbol{\lambda}$ (Eq. 21)

**13**    **end**

**14 end**

---

## V. IMPLEMENTATION DETAILS

### A. Our Approach

We train PixelCNN++ on the downsampled $64 \times 64$ ImageNet data as introduced in [14] for 6 epochs. Batch size is kept as 36 and the number of filter channels as 100. The rest of the parameters are same as the ones used for training PixelCNN++ on $32 \times 32$ ImageNet in [16]. We obtain a negative log likelihood score of 3.66 on test data and 3.5 on train data which is consistent with the numbers reported in [16] for similar data.

With this learned model, we use our proposed algorithm as described in Algorithm 1, for the experiments described below. An initial image is sampled from a uniform random distribution. However, we observe that starting with different initial images doesn't have much effect on the final converged reconstruction. We use momentum in the gradient update for faster convergence, with its value set to 0.9. Step size $\alpha$, maximum iterations, likelihood weightage $\rho$ for each experiment are mentioned in the subsequent section.

For reconstructing color images, we consider multiplexing along individual color channels. Hence, we have separate $\Phi$ matrices for all the three channels and obtain three separate measurement vectors $\mathbf{Y}$ for each channel.

We have made the code of our implementation for the task Single Pixel Camera reconstruction available online[1].

### B. One Network to solve them all

We use the original implementation of [1] available online[2] with certain modification as mentioned below.

For simulating color Single Pixel Camera, the original implementation rasterizes the entire $N \times N \times 3$ image into a single vector and creates one $\Phi$ matrix to compress this into a single measurement vector. We believe that this might not be feasible to implement in a real system. Hence, we modify their implementation to instead simulate separate $\Phi$ matrices for each channel as in Section V-A.

While simulating SPC measurements on large images, the original implementation only deals with local multiplexing. It breaks them down into patches of $64 \times 64$ and compresses each of these patches separately. We modify this to deal with the more challenging case of global multiplexing, where we compress the entire image.

We extend the original implementation for LiSens and FlatCam as well, by considering the above modifications and incorporating the respective forward models.

We use model provided which was trained on $64 \times 64$ Imagenet for 2 epochs for testing the results. We found that the results were very much dependent on the alpha parameter (penalty parameter) which had to be tuned for each image to get the best solution.

### C. TVAL3

For comparisons with TVAL3 ( TV minimization by Augmented Lagrangian and ALternating direction algorithms ) [20], we use the MATLAB implementation[3] with the default parameters. The number of iterations is set to 80. For color image reconstruction, we update each channel separately using TVAL3.

## VI. EXPERIMENTS

In this section, we present the reconstructions from our approach and compare them with the existing state-of-the-art approaches. To being with, we illustrate the ability of an

---

[1]https://github.com/adaveiitm/deep-pixel-level-prior

[2]https://github.com/rick-chang/OneNet

[3]http://www.caam.rice.edu/ optimization/L1/TVAL3/

| Original image | Masked image | OneNet | Ours |
|---|---|---|---|

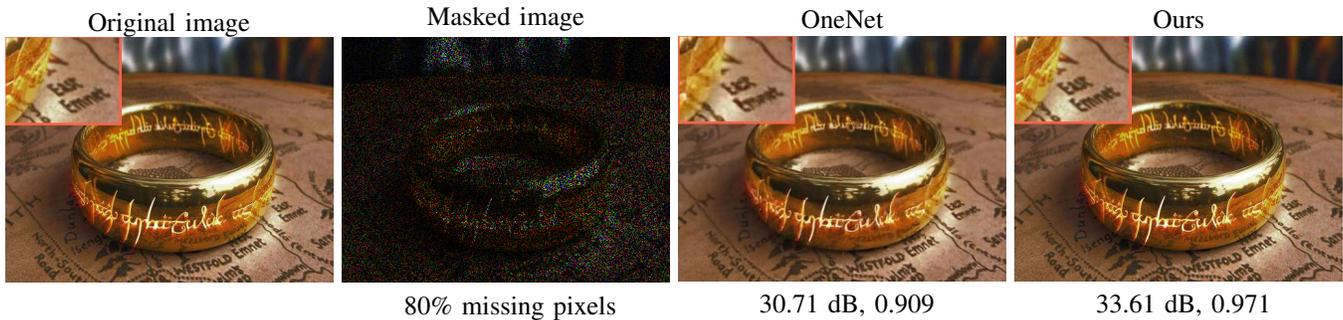| | 80% missing pixels | 30.71 dB, 0.909 | 33.61 dB, 0.971 |

Fig. 2. Random pixel inpainting with 80% missing pixels. Our approach reconstructs the finer edges better and has more consistency among neighbouring pixels, as compared to OneNet. Note the details around the text shown in zoomed patch. The difference between the two reconstructions can be perceived by further zooming into the images. The numbers reported in this and the subsequent figures are PSNR (in dB) followed by SSIM.

autoregressive prior in reconstructing pixel level details using an example of missing pixel inpainting in an image. For this, we randomly mask out pixels from the image and use our prior to reconstruct these missing pixels. We perform by keeping the observed pixel values as same and update missing pixels to maximize the the prior loglikelihood. Specifically, we take an image of size 384x512 and mask 80% of the pixels in the initial image as could be seen in Figure 2. We compare our results with that of OneNet [1], and we can observe details in our reconstruction much better like the text outlines, also quantitatively in terms of PSNR and SSIM. We use a step size of 75 and run for approximately 1000 iterations.

For all the three imaging setups of SPC, Lisens and Flatcam we perform reconstructions on both simulated data and real measurements. In case of simulation we compare our reconstructions with TVAL3 [20] and OneNet [1]. In case of reconstructions from real measurements, we compare our results with TVAL3. OneNet experiments failed to converge to a stable point in this case hence we could not provide comparison with this approach. For real Lisens at 66% measurements, although OneNet converges, results obtained were very poor compared to other approaches.

### A. Single Pixel Camera

*1) Simulation case:* We show quantitative and qualitative comparisons of simulated SPC reconstruction results on images of sizes 128×128 and 256×256 respectively as shown in Table I and Figure 10 respectively. Measurement rates considered are 10% and 25% for 128×128 and 5% and 10% for 256×256. Similar to RIDE-CS [23], we generate the $\phi$ matrix as a random Gaussian with orthonormal rows. We perform gradient descent and projection operation on the compressed image for 2000 iterations in the case of 25% measurement rate and for 2500 iterations in case of 10% measurement rate. We use a step-size of 7.5 and the hard constraint projection method. In all cases, we intialize with random image from uniform distribution. We compare our results to [1] and we are able to show significant improvement in reconstruction results in terms of PSNR and SSIM values. Our reconstructions have better edges and textures compared to the reconstructions from OneNet.

*2) Real Case:* We show our real SPC reconstruction results in Figure 4. Data for this experiment is provided to us by the
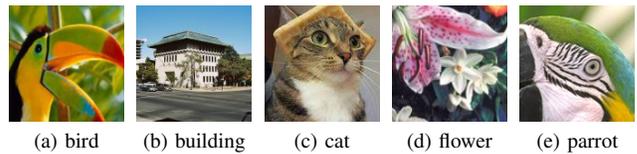
| (a) bird | (b) building | (c) cat | (d) flower | (e) parrot |
|---|---|---|---|---|

Fig. 3. Test images of $128 \times 128$ size chosen randomly for simulated SPC and LiSens reconstructions.

TABLE I
COMPARISONS OF RECONSTRUCTIONS FROM SIMULATED SPC MEASUREMENTS AT DIFFERENT MEASUREMENT RATES FOR THE IMAGES SHOWN IN FIGURE 3. OUR APPROACH OBTAINS BETTER PERFORMANCE THAN ONENET AND TVAL3 BY MODELLING PIXEL-LEVEL CONSISTENCIES. SEE FIGURE 10 FOR QUALITATIVE COMPARISONS

| Name | M.R. | TVAL3 | | OneNet | | Ours | |
|---|---|---|---|---|---|---|---|
| | | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| bird | 10 | 23.67 | 0.91 | 23.92 | 0.93 | **29.52** | **0.97** |
| | 25 | 29.67 | 0.97 | 26.89 | 0.96 | **32.96** | **0.98** |
| building | 10 | 18.81 | 0.61 | 23.85 | 0.86 | **25.93** | **0.88** |
| | 25 | 22.72 | 0.79 | 24.06 | 0.87 | **32.05** | **0.96** |
| cat | 10 | 23.27 | 0.72 | 25.15 | 0.82 | **26.68** | **0.85** |
| | 25 | 26.87 | 0.85 | 26.60 | 0.88 | **31.23** | **0.94** |
| flower | 10 | 20.07 | 0.68 | 23.39 | 0.84 | **26.22** | **0.89** |
| | 25 | 24.84 | 0.86 | 25.13 | 0.90 | **31.05** | **0.96** |
| parrot | 10 | 18.49 | 0.64 | 25.82 | 0.89 | **27.59** | **0.90** |
| | 25 | 23.67 | 0.84 | 26.79 | 0.91 | **32.18** | **0.95** |
| mean | 10 | 20.86 | 0.72 | 24.43 | 0.87 | **27.19** | **0.90** |
| | 25 | 25.55 | 0.86 | 25.74 | 0.90 | **31.89** | **0.96** |

authors of [3]. We obtain the real SPC sensor measurements at 30% and 15% measurement rate respectively. The images we reconstruct in this case are grey scale images. Here also, we use the Hard constraint projection method for inference. We compare our results with TVAL3 and RIDE-CS [23]. Our method performs better than both RIDE-CS and TVAL3 in terms of PSNR and SSIM values. Apart from these measures, we observe that our method produces a sharper reconstruction. We use the same hyperparameters and training procedures as
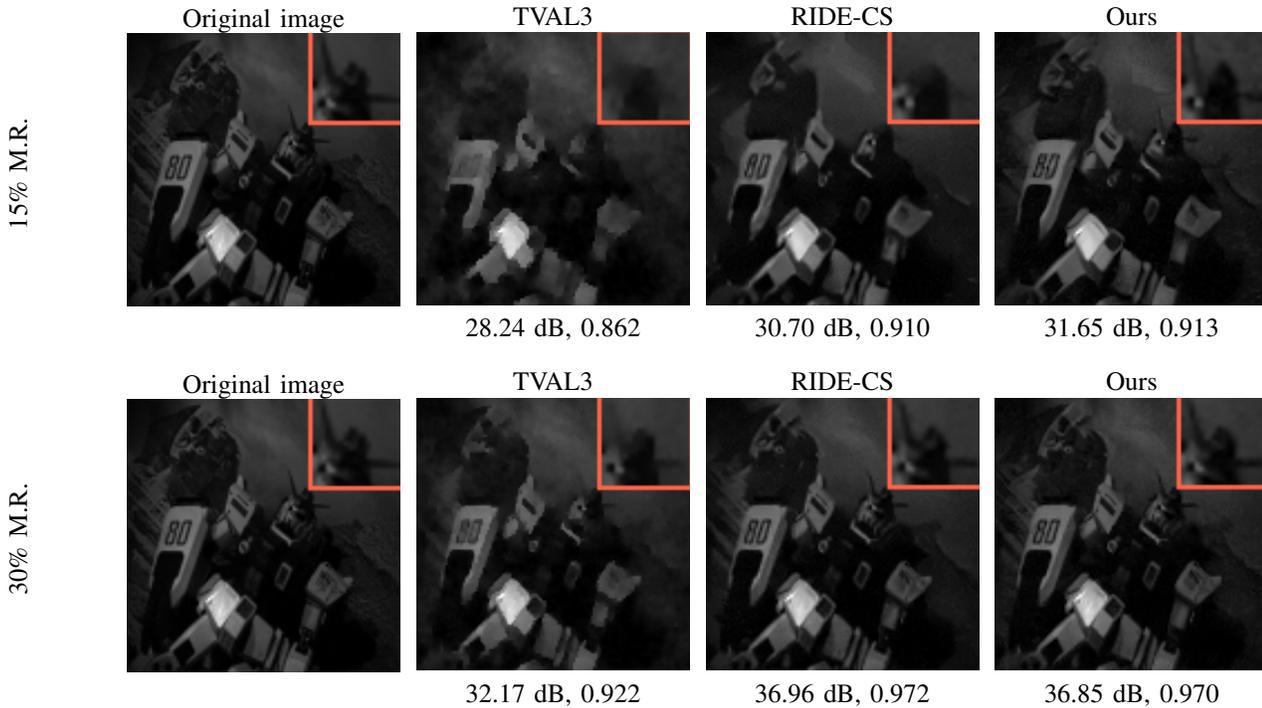
Fig. 4. Reconstructions from real Single Pixel Camera measurements at different measurement rates. Our approach recovers the low level details much better than TVAL3. Though the performance of RIDE-CS, which is also a deep autoregressive model, is similar to ours in this case, its computational complexity is much higher. Also in other simulation experiments we found RIDE-CS does not preserve fine details, see Figure 9.

in the simulated case.

### B. LiSens

*1) Simulation case:* The reconstruction in case of simulated LiSens is done at 25% and 40% measurement rates. Our LiSens experiments, similar to SPC experiments, have been done on both 128x128 and 256x256 images as shown in Table II and Figure 11 respectively. We compare our reconstructions with that obtained using OneNet. Our method provides better results in terms of visual perception as well as PSNR and SSIM values. Our reconstructions have well-defined boundaries of different objects in the image and do not produce artifacts which are observed in case of OneNet. We have used hard constraint case for the simulated LiSens reconstruction for approximately 2000 iterations with a step-size of 7.5.

*2) Real measurements:* The real LiSens experiments have been done at 16% and 33% measurement rates obtained at a resolution of $768 \times 256$ , as provided by the authors of [3]. We compare our real Lisens with TVAL3 as in Figure 7. Our method performs better reconstruction with respect to low level details in the image. Our proposed methods reconstruction has little or no blur compared to TVAL3 and the reconstruction is sharper in terms of object boundaries in the image. We use Hard constraint method for reconstruction with 25% dropout in pixel-wise update. We use an update step of 7.5 and 2000 iterations for reconstruction, similar to simulated experiment.

### C. FlatCam

*1) Simulation case:* The matrices $\Phi_L$ and $\Phi_R$ in the Flat-Cam imaging model are estimated based on the calibration

TABLE II
COMPARISONS OF RECONSTRUCTIONS FROM SIMULATED LISENS MEASUREMENTS AT DIFFERENT MEASUREMENT RATES FOR THE IMAGES SHOWN IN FIGURE 3.OUR APPROACH OBTAINS BETTER PERFORMANCE THAN ONENET AND TVAL3. SEE FIGURE 11 FOR QUALITATIVE COMPARISONS

| Name | M.R. | TVAL3 | | OneNet | | Ours | |
|---|---|---|---|---|---|---|---|
| | | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| bird | 25 | 24.59 | 0.95 | 24.98 | 0.82 | **27.13** | **0.96** |
| | 40 | 29.34 | 0.98 | 27.52 | 0.96 | **34.14** | **0.99** |
| building | 25 | 18.72 | 0.67 | 21.16 | 0.79 | **30.87** | **0.95** |
| | 40 | 23.41 | 0.82 | 22.41 | 0.84 | **35.06** | **0.98** |
| cat | 25 | 23.41 | 0.67 | 27.27 | 0.89 | **29.95** | **0.94** |
| | 40 | 25.83 | 0.87 | 29.03 | 0.92 | **34.65** | **0.97** |
| flower | 25 | 21.00 | 0.72 | **27.85** | **0.91** | 26.54 | 0.88 |
| | 40 | 23.66 | 0.83 | **30.79** | **0.95** | 30.21 | 0.93 |
| parrot | 25 | 15.27 | 0.65 | 26.02 | 0.90 | **30.17** | **0.94** |
| | 40 | 19.75 | 0.85 | 27.99 | 0.93 | **32.35** | **0.96** |
| mean | 25 | 20.60 | 0.73 | 25.45 | 0.89 | **28.93** | **0.94** |
| | 40 | 24.40 | 0.87 | 27.55 | 0.92 | **33.28** | **0.97** |

procedure mentioned in [17]. As we want to deal with RGB images, separate $\Phi_L$ and $\Phi_R$ matrices are calibrated for each of the R, G and B channels with the help of a Bayer color filter array on the sensor. We compare our results with OneNet and L2 regularisation, on two 256x256 images as shown in Figure

| Original image | L2 Reg. | OneNet | Ours |
|---|---|---|---|
| | 12.79 dB, 0.71 | 20.11 dB, 0.84 | 20.22 dB, 0.85 |
| | 10.77 dB, 0.59 | 19.52 dB, 0.81 | 25.08 dB, 0.83 |

Fig. 5. Reconstructions ($256 \times 256$) from simulated FlatCam measurements using L2 regularization, OneNet and our approach. Note the suppression of vignetting effect in our results which is clearly visible in the house image.
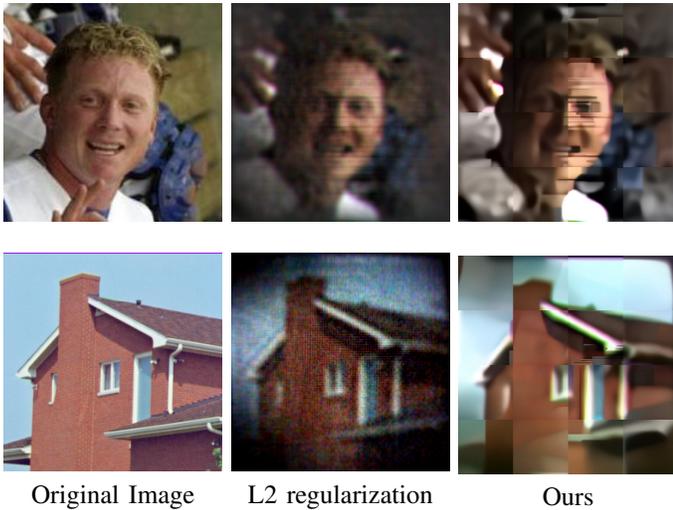
Original Image  L2 regularization  Ours

Fig. 6. Qualitative comparisons of reconstructions obtained from real FlatCam measurements using calibrated $\Phi_L$ and $\Phi_R$ using L2 regularization and our approach. Real reconstructions are not good because of calibration error and separability assumption in the forward model.

Reconstructions with 16% M.R.

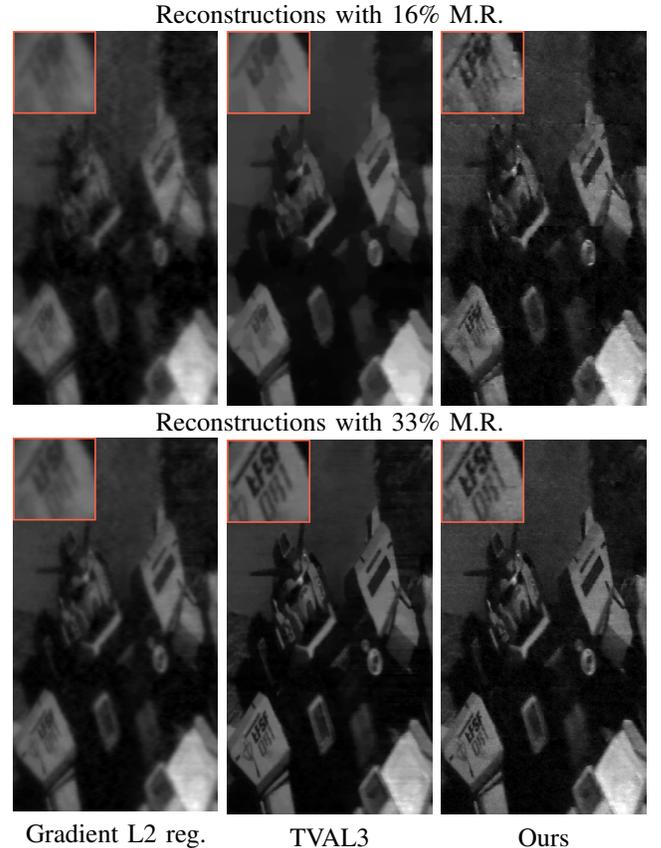Reconstructions with 33% M.R.

Gradient L2 reg.  TVAL3  Ours

Fig. 7. Qualitative comparisons of reconstructions from real LiSens measurements at different measurement rates. Reconstructions using our method are sharper and preserve the overall structure.

5. Our method shows better PSNR, SSIM, and perceptually better quality samples. Our method produces the least blurry solution and objects in the image has well defined boundaries. We use 25% pixel dropout and perform 1000 iterations of augmented Lagrangian method with the step size $\alpha$ as 60.0 and $\rho$ as 10.

*2) Real measurements:* We use the data provided by the authors of [17]. The original images were displayed on a monitor and captured using FlatCam. Using a Bayer color filter on the sensor, separate measurements for the three color channels can be obtained. We compare our reconstructions with L2 regularization as shown in Figure 6. Our reconstructions are more accurate in terms of brightness, boundaries and sharpness of the image. We use soft constraint case for reconstruction

and use the same hyperparameters as in the simulation case.

We observe that reconstructions from real FlatCam are not qualitatively as good as with real SPC and LiSens measurements. This is because the forward model assumed in this case

is erroneous. Firstly, there are calibration errors in estimating the $\Phi_L$ and $\Phi_R$ matrices. Secondly, the forward model in [17] relies on the separability assumption leading to model error.

### D. Ablation Experiments

*1) Effect of pixel-wise dropout:* In this experiment, we vary the amount of pixels not updated in each iteration and observe its effect on the reconstructed image, see Figure 8. When the dropout ratio is zero, the area in the image having texture is over smooth. With considerable dropout ratio (25%), the texture is reconstructed better amounting to a higher PSNR and SSIM. However, on increasing it further, the reconstructions appear noisy with a reduction in quality. Thus, for all our experiments, we used 25% dropout.



| Original image | 0% pixel dropout | 25% pixel dropout |
| --- | --- | --- |
| | 26.34 dB, 0.826 | 27.93 dB, 0.887 |

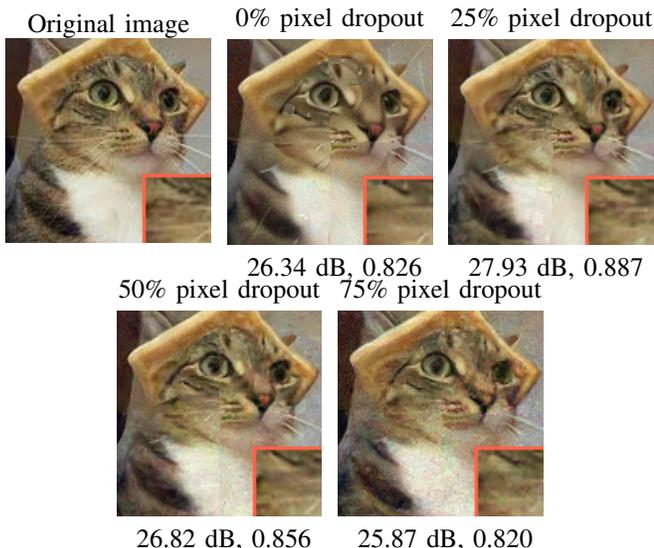| 50% pixel dropout | 75% pixel dropout |
| --- | --- |
| 26.82 dB, 0.856 | 25.87 dB, 0.820 |

Fig. 8. Effect of varying the amount of pixel dropout for SPC reconstruction at 15% measurement rate. By not updating a certain amount of pixels every iteration, the texture is reconstructed better and the image has a higher quality. However, on increasing this dropout ratio more than a certain level, the reconstructions become noisy and the quality reduces.

*2) Comparison with Ride-CS - grayscale SPC:* While we train our model on colored Imagenet data, we observe that in practice this approach works well on reconstructing grayscale images as well. We compare our reconstruction with that of RIDE-CS [23], which uses the autoregressive model RIDE [15] as image prior. In Figure 9, we compare the reconstruction of a grayscale image from Single Pixel Camera measurements using our approach and RIDE-CS for 15% measurement rate. The reconstruction obtained from our approach is better than that of RIDE-CS. This is because we use PixelCNN++ which is a deeper network than RIDE and hence has better representation power. Also, the running time of our approach ($\sim 5$ minutes) is much less than that of RIDE-CS ($\sim 30$ minutes). Our approach is CNN based and hence can be parallelized over multiple GPUs while RIDE-CS relies on a network of spatial LSTMs which are tough to parallelize.

*3) Comparison with OneNet in their original setting:* Till now we have performed all the experiments with different $\Phi$ matrix for each color channel. However in OneNet [1], the authors have considered one $\Phi$ matrix that multiplexes



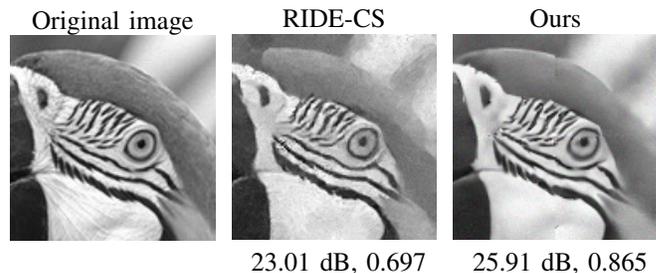| Original image | RIDE-CS | Ours |
| --- | --- | --- |
| | 23.01 dB, 0.697 | 25.91 dB, 0.865 |

Fig. 9. Comparison with Ride-CS on reconstruction of grayscale image from simulated Single Pixel Camera measurements at 15% measurement rate. Our reconstructions have a higher quality and are perceptually more closer to the true image.

across the three color channels, which might not be feasible to implement in a real system. For this ablation experiment, we consider the original setting as used in [1] and compare their reconstructions with ours for 10% SPC reconstruction on the 9 test ImageNet images mentioned in the [1]. PSNR and SSIM values for the same are mentioned in Table III. Our approach performs better than OneNet.

TABLE III
COMPARISONS OF COMPRESSIVE IMAGING RECONSTRUCTIONS FOR IMAGES PROVIDED IN [1] WITH THEIR SETTING OF MULTIPLEXING ACROSS COLOR CHANNELS. HOWEVER, THIS WAY OF MULTIPLEXING ACROSS THE COLOR CHANNELS MIGHT NOT BE FEASIBLE IN A REAL SYSTEM.

| Figure Name | OneNet | | Ours | |
| --- | --- | --- | --- | --- |
| | PSNR | SSIM | PSNR | SSIM |
| ball | 24.696 | 0.9023 | **26.656** | **0.9300** |
| dalmatian | 20.650 | 0.8314 | **21.812** | **0.8518** |
| dog | 26.873 | 0.8734 | **28.552** | **0.8952** |
| field | 26.470 | 0.9112 | **29.017** | **0.9149** |
| man | 29.152 | 0.9460 | **31.787** | **0.9540** |
| mountain | 25.484 | 0.8821 | **28.993** | **0.8912** |
| table | 19.397 | **0.8083** | **20.955** | 0.6662 |
| woman | 25.512 | 0.8518 | **27.321** | **0.8906** |
| wolf | 25.976 | 0.8839 | **28.355** | **0.9061** |

### VII. DISCUSSION AND CONCLUSION

We demonstrate the efficacy of deep pixel level image prior for ill-posed reconstruction in different computational imaging problems. Among the three proposed approaches for inference, hard and soft constraint based and ALM based, overall, soft constraint-based method works well and can handle noisy measurements by appropriately varying the tuning parameter, $\lambda$. However, when there is no noise or less noise in the measurements, the hard constraint-based method performs as good as soft constraint case with an additional advantage of being parameter free and hence is preferable. In fact, for our real experiments on SPC (Figure 4) and Lisens (Figure 7), we use hard constraint-based inference, which produces reasonable results. For cases such as Flatcam, non-invertibility of $\Phi\Phi^T$ prevents the use of hard-constraint based inference.

Our approach enjoys the versatility of image priors and rich feature representation of deep neural networks. Being pixel level, it explicitly accounts for pixel level correlations resulting in consistent texture and edges. We show our evaluations on both the simulation of forward models and data from real setups. In all cases, both quantitative and qualitative metrics suggest that our approach performs better than traditional methods and current state-of-the-art learning based methods. An interesting line of work would be to incorporate deviations from the forward model, due to calibration and model errors, in our approach to further improve the quality of reconstruction for FlatCam.

## References

[1] JH Chang, Chun-Liang Li, Barnabás Póczos, BVK Kumar, and Aswin C Sankaranarayanan, "One network to solve them all—solving linear inverse problems using deep projection models," *arXiv preprint arXiv:1703.09912*, 2017. 2, 3, 6, 7, 10

[2] Marco F Duarte, Mark A Davenport, Dharmpal Takhar, Jason N Laska, Ting Sun, Kevin E Kelly, Richard G Baraniuk, et al., "Single-pixel imaging via compressive sampling," *IEEE Signal Processing Magazine*, vol. 25, no. 2, pp. 83, 2008. 1, 2, 4

[3] Jian Wang, Mohit Gupta, and Aswin C Sankaranarayanan, "Lisens-a scalable architecture for video compressive sensing," in *Computational Photography (ICCP), 2015 IEEE International Conference on*. IEEE, 2015, pp. 1–9. 1, 2, 3, 4, 7, 8

[4] Dikpal Reddy, Ashok Veeraraghavan, and Rama Chellappa, "P2c2: Programmable pixel compressive camera for high speed imaging," in *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*. IEEE, 2011, pp. 329–336. 1

[5] Yasunobu Hitomi, Jinwei Gu, Mohit Gupta, Tomoo Mitsunaga, and Shree K Nayar, "Video from a single coded exposure photograph using a learned over-complete dictionary," in *Computer Vision (ICCV), 2011 IEEE International Conference on*. IEEE, 2011, pp. 287–294. 1

[6] Ashwin Wagadarikar, Renu John, Rebecca Willett, and David Brady, "Single disperser design for coded aperture snapshot spectral imaging," *Applied optics*, vol. 47, no. 10, pp. B44–B51, 2008. 1

[7] Leonid I Rudin, Stanley Osher, and Emad Fatemi, "Nonlinear total variation based noise removal algorithms," *Physica D: nonlinear phenomena*, vol. 60, no. 1-4, pp. 259–268, 1992. 1, 3

[8] Javier Portilla, Vasily Strela, Martin J Wainwright, and Eero P Simoncelli, "Image denoising using scale mixtures of gaussians in the wavelet domain," *IEEE Transactions on Image processing*, vol. 12, no. 11, pp. 1338–1351, 2003. 1

[9] Michal Aharon, Michael Elad, and Alfred Bruckstein, "*k*-svd: An algorithm for designing overcomplete dictionaries for sparse representation," *IEEE Transactions on signal processing*, vol. 54, no. 11, pp. 4311–4322, 2006. 1

[10] Kuldeep Kulkarni, Suhas Lohit, Pavan Turaga, Ronan Kerviche, and Amit Ashok, "Reconnet: Non-iterative reconstruction of images from compressively sensed measurements," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 449–458. 1, 3

[11] Hantao Yao, Feng Dai, Dongming Zhang, Yike Ma, Shiliang Zhang, and Yongdong Zhang, "Dr ^{2}-net: Deep residual reconstruction network for image compressive sensing," *arXiv preprint arXiv:1702.05743*, 2017. 1, 3

[12] Chris Metzler, Ali Mousavi, and Richard Baraniuk, "Learned d-amp: Principled neural network based compressive image recovery," in *Advances in Neural Information Processing Systems*, 2017, pp. 1770–1781. 1

[13] Ayan Sinha, Justin Lee, Shuai Li, and George Barbastathis, "Lensless computational imaging through deep learning," *Optica*, vol. 4, no. 9, pp. 1117–1125, 2017. 1

[14] Aaron van den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu, "Pixel recurrent neural networks," *arXiv preprint arXiv:1601.06759*, 2016. 1, 3, 6

[15] Lucas Theis and Matthias Bethge, "Generative image modeling using spatial lstms," in *Advances in Neural Information Processing Systems*, 2015, pp. 1927–1935. 1, 3, 10

[16] Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P Kingma, "Pixelcnn++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications," *arXiv preprint arXiv:1701.05517*, 2017. 1, 3, 4, 6

[17] M Salman Asif, Ali Ayremlou, Aswin Sankaranarayanan, Ashok Veeraraghavan, and Richard G Baraniuk, "Flatcam: Thin, lensless cameras using coded aperture and computation," *IEEE Transactions on Computational Imaging*, vol. 3, no. 3, pp. 384–397, 2017. 2, 3, 4, 8, 9, 10

[18] Nick Antipa, Grace Kuo, Ren Ng, and Laura Waller, "3d diffusercam: Single-shot compressive lensless imaging," in *Computational Optical Sensing and Imaging*. Optical Society of America, 2017, pp. CM2B–2. 2

[19] Antonin Chambolle, "An algorithm for total variation minimization and applications," *Journal of Mathematical imaging and vision*, vol. 20, no. 1-2, pp. 89–97, 2004. 3

[20] Chengbo Li, Wotao Yin, Hong Jiang, and Yin Zhang, "An efficient augmented lagrangian method with applications to total variation minimization," *Computational Optimization and Applications*, vol. 56, no. 3, pp. 507–530, 2013. 3, 6, 7

[21] Huaijin Chen, M Salman Asif, Aswin C Sankaranarayanan, and Ashok Veeraraghavan, "Fpa-cs: Focal plane array-based compressive imaging in short-wave infrared," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 2358–2366. 3

[22] Aswin C Sankaranarayanan, Christoph Studer, and Richard G Baraniuk, "Cs-muvi: Video compressive sensing for spatial-multiplexing cameras," in *Computational Photography (ICCP), 2012 IEEE International Conference on*. IEEE, 2012, pp. 1–10. 3

[23] Akshat Dave, Anil Kumar, Kaushik Mitra, et al., "Compressive image recovery using recurrent generative model," in *Image Processing (ICIP), 2017 IEEE International Conference on*. IEEE, 2017, pp. 1702–1706. 3, 7, 10

[24] Christopher A Metzler, Arian Maleki, and Richard G Baraniuk, "From denoising to compressed sensing," *IEEE Transactions on Information Theory*, vol. 62, no. 9, pp. 5117–5144, 2016. 3

[25] Julio Martin Duarte-Carvajalino and Guillermo Sapiro, "Learning to sense sparse signals: Simultaneous sensing matrix and sparsifying dictionary optimization," *IEEE Transactions on Image Processing*, vol. 18, no. 7, pp. 1395–1408, 2009. 3

[26] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio, "Generative adversarial nets," in *Advances in Neural Information Processing Systems*, 2014, pp. 2672–2680. 3

[27] Diederik P Kingma and Max Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013. 3

[28] Singanallur V Venkatakrishnan, Charles A Bouman, and Brendt Wohlberg, "Plug-and-play priors for model based reconstruction," in *Global Conference on Signal and Information Processing (GlobalSIP), 2013 IEEE*. IEEE, 2013, pp. 945–948. 3

[29] Kostadin Dabov, Alessandro Foi, Vladimir Katkovnik, and Karen Egiazarian, "Bm3d image denoising with shape-adaptive principal component analysis," in *SPARS'09-Signal Processing with Adaptive Sparse Structured Representations*, 2009. 3

[30] Kai Zhang, Wangmeng Zuo, Shuhang Gu, and Lei Zhang, "Learning deep cnn denoiser prior for image restoration," *arXiv preprint arXiv:1704.03264*, 2017. 3

[31] Tim Meinhardt, Michael Möller, Caner Hazirbas, and Daniel Cremers, "Learning proximal operators: Using denoising networks for regularizing inverse imaging problems," *ArXiv e-prints, Apr*, 2017. 3

[32] Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang, "On large-batch training for deep learning: Generalization gap and sharp minima," *arXiv preprint arXiv:1609.04836*, 2016. 6

Fig. 10. Qualitative comparisons of $256 \times 256$ images reconstructed from simulated Single Pixel Camera measurements using TVAL3, OneNet and our approach. Even when the measurement rate is low, our method reconstructs the sharp and prominent structures in the image better. Moreover, there are no visible artifacts in our reconstructions as the autoregressive prior ensures the nearby pixels to be consistent. This is not the case with TVAL3 and OneNet leading to poor performance.

Fig. 11. Qualitative comparisons of images reconstructed from simulated LiSens measurements using TVAL3, OneNet and our approach. Reconstructions from our approach have minimal artifacts and are closer to the original image.