

Deep End-to-End Posterior ENergy (DEEPEN) for image recovery

Jyothi Rikhab Chand, Member, IEEE, and Mathews Jacob, Fellow, IEEE

Abstract—Current end-to-end (E2E) and plug-and-play (PnP) image reconstruction algorithms approximate the maximum a posteriori (MAP) estimate but cannot offer sampling from the posterior distribution, like diffusion models. By contrast, it is challenging for diffusion models to be trained in an E2E fashion. This paper introduces a Deep End-to-End Posterior ENergy (DEEPEN) framework, which enables MAP estimation as well as sampling. We learn the parameters of the posterior, which is the sum of the data consistency error and the negative log-prior distribution, using maximum likelihood optimization in an E2E fashion. The proposed approach does not require algorithm unrolling, and hence has a smaller computational and memory footprint than current E2E methods, while it does not require contraction constraints typically needed by current PnP methods. Our results demonstrate that DEEPEN offers improved performance than current E2E and PnP models in the MAP setting, while it also offers faster sampling compared to diffusion models. In addition, the learned energy-based model is observed to be more robust to changes in image acquisition settings.

Index Terms—Energy model, MAP estimate, Memory-efficient, Parallel MRI, Uncertainty estimate

I. INTRODUCTION

Computational algorithms that can recover images from sparse and noisy Fourier measurements have revolutionized magnetic resonance (MR) imaging. Traditional compressed sensing (CS) algorithms rely on hand-crafted image priors as regularizers [1]–[3]; wherein at each iteration they alternate between the update of data consistency (DC) and the proximal map of the regularizer to recover the unknown image.

Data-driven deep learning methods have significantly improved reconstruction performance over classical methods. Plug-and-play (PnP) methods use an iterative algorithm that alternates between denoising and the DC update step, where the CNN-based pre-trained denoiser

replaces the proximal map in the CS algorithm [4]–[6]. A contraction constraint is often needed to ensure fixed-point convergence [7], which translates to reduced performance [8]. Recently, PnP methods that use explicit CNN-based energy functions to model the negative log-prior [9]–[11] and are trained using denoising score matching (DSM) [12] were introduced. These methods do not require a contraction constraint to guarantee stationary point convergence, translating into improved performance [9]–[11]. Energy models also enable the generation of samples from the target distribution. The implicit multi-scale energy model (i-MuSE), which pre-learns a single energy function using DSM at multiple noise scales, was found to offer improved convergence and hence superior performance compared to a single-scale DSM [11]. Unlike PnP models that pre-learn CNN, end-to-end (E2E) methods [13]–[17] unroll the iterative algorithm and optimize the parameters of the CNN denoiser so that the reconstructed image matches the reference image. These approaches offer better performance over PnP methods that are agnostic to the specific forward model. However, the E2E models are associated with increased memory demand, which can be reduced using the deep equilibrium (DEQ) formulation, but requires a contraction constraint similar to PnP methods that restricts its performance [8,18,19]. Since energy-based DEQ methods do not require this constraint [20], they offer improved performance.

The above approaches approximate the maximum a posteriori (MAP) estimate. In contrast, we learn the posterior distribution in an E2E fashion. Once trained, the learned posterior can be used to derive the MAP estimate and also sample from it. We derive the parameters of our Deep E2E Posterior ENergy (DEEPEN) network by maximizing the likelihood of data samples as in [21]. In particular, we minimize the negative log-posterior of the training samples, which is the sum of the negative log-likelihood and the negative log-prior modeled by the CNN. The above training strategy simplifies to the minimization of the energy of the *true* reference samples, while maximizing the energy of the *fake* samples that are obtained using the Langevin sampling algorithm using the gradient of the learned posterior. The training strategy

The authors are with the Department of Electrical and Computer Engineering, University of Virginia, Charlottesville, VA 22904, USA (e-mail: jyothi-rikhabchand@virginia.edu; mja-cob@virginia.edu). This work is supported by NIH R01AG067078, R01 EB019961, and R01 EB031169.

thus resembles a generative adversarial network (GAN) [22,23]. The key difference is that the generator involves a Langevin-based iterative algorithm specified by the energy model, and the classifier is also specified by the energy model. We show that when the Langevin noise level is small, the above training strategy does not require algorithm unrolling or DEQ strategies associated with high memory demand and computational complexity. This E2E training strategy is computationally and memory-efficient because it does not require unrolling or fixed-point iterations. Moreover, unlike current E2E methods it does not use Mean Squared Error (MSE) loss to learn a MAP estimate.

We note that the DSM approaches used in PnP models are trained to remove Gaussian noise perturbations, where the pixels of the perturbations are independent and identically distributed. However, the MR imaging inverse problem aims to recover images from corrupted undersampled Fourier measurements, where the corruption/perturbation often has highly correlated pixel values. Unfortunately, DSM-trained models are often not efficient in estimating and removing correlated noise. By contrast, our experiments show that DEEPEN gradients are more effective in estimating Gaussian as well as correlated perturbations, thus offering improved performance in MAP image recovery. Our experiments also show that the learned E2E approach also generalizes to unseen acquisition settings, unlike traditional E2E approaches that are not robust to mismatch in forward models from the training settings.

The DEEPEN model can also enable sampling from the posterior, where it can offer 10x faster sampling than diffusion models [24], with comparable image quality, and reduced uncertainty, even when CNN complexity is 5x smaller. We note that deterministic ODE flows [25,26] that learn straighter paths between probability distributions have been introduced to speed up prior sampling drastically. However, these models still require >100 steps in the posterior sampling setting. In particular, since they are trained only along the path between the distribution, noise often needs to be added after each gradient descent step involving data-consistency to ensure that the integration path remains in the regions trained originally [26].

A preliminary conference version of this approach was previously published with limited empirical results [27].

II. END-TO-END LEARNED ENERGY MODEL

Let $\mathbf{b} \in \mathbb{C}^n$ denote the noisy undersampled measurements from which we wish to recover the unknown image $\mathbf{x} \in \mathbb{C}^m$. The image and the measurements

are related through a known linear forward operator $\mathbf{A} \in \mathbb{C}^{n \times m}$:

$$\mathbf{b} = \mathbf{A}\mathbf{x} + \mathbf{n} \quad (1)$$

where $\mathbf{n} \in \mathcal{N}(0, \eta^2 \mathbf{I})$ is additive complex white Gaussian noise. The recovery of the image from the measurements is often formulated as a MAP estimation problem, where the solution is the maximum of the log-posterior:

$$\log q(\mathbf{x}|\mathbf{b}) = \log p(\mathbf{b}|\mathbf{x}) + \log q(\mathbf{x}) \quad (2)$$

Here, $q(\mathbf{x})$ is the prior distribution, and the log-likelihood term $\log p(\mathbf{b}|\mathbf{x})$ is specified by:

$$\log p(\mathbf{b}|\mathbf{x}) = -\frac{\|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2}{2\eta^2} + P \quad (3)$$

where P is the normalization constant. The log-posterior can also be used to derive samples from the posterior distribution, which can inform about uncertainties in the estimation and offer likely solutions.

A. Deep End-to-End Posterior ENergy (DEEPEN)

We approximate the prior distribution $q(\mathbf{x})$ by a parametric energy model $p_\theta(\mathbf{x})$ as in [21]:

$$p_\theta(\mathbf{x}) = \frac{1}{Z_\theta} \exp(-\mathcal{E}_\theta(\mathbf{x})), \quad (4)$$

where $\mathcal{E}_\theta(\mathbf{x}) : \mathbb{C}^m \rightarrow \mathbb{R}^+$ is modeled by a neural network and Z_θ is the normalization constant. Using (3), we obtain the parametric model for the log-posterior distribution as:

$$\log p_\theta(\mathbf{x}|\mathbf{b}) = \log p(\mathbf{b}|\mathbf{x}) + \log p_\theta(\mathbf{x}), \quad (5)$$

where $p_\theta(\mathbf{x})$ is the prior in (4). Combining (3) and (4), we obtain the negative log-posterior $\mathcal{L}_\theta(\mathbf{x}) = -\log p_\theta(\mathbf{x}|\mathbf{b})$ as:

$$\mathcal{L}_\theta(\mathbf{x}) = \underbrace{\frac{1}{2}\|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 + \mathcal{E}_\theta(\mathbf{x})}_{C_\theta(\mathbf{x}; \mathbf{A}, \mathbf{b})} + \log \tilde{Z}_\theta, \quad (6)$$

where

$$\tilde{Z}_\theta = \int e^{-C_\theta(\mathbf{x}; \mathbf{A}, \mathbf{b})} d\mathbf{x} \quad (7)$$

is a normalizing constant. For simplicity we absorbed the parameter η^2 into the definition of energy.

The common approach in energy-based models (EBMs) is to pre-learn p_θ from fully sampled images, independent of the forward operator \mathbf{A} [11,21,28]. Once the learning is complete, image recovery involves sampling from (5) using the learned p_θ or maximizing (5) as in [11]. We note that learning in such PnP methods is agnostic to the forward model \mathbf{A} . However, methods based on E2E training (for example, [8,15,16,19]) often offer better performance than PnP methods. These E2E

methods use pairs of measurements and images (\mathbf{x}, \mathbf{b}) to learn p_θ so that the maximum of (5) matches the reference image, i.e., they try to learn the MAP estimate. Motivated by the improved performance of E2E methods, we propose to train the EBM in an E2E fashion for a specific forward operator \mathbf{A} in the next section.

B. Maximum likelihood training of DEEPEN

We determine the optimal weights of $\mathcal{E}_\theta(\mathbf{x})$ by minimizing the negative log-likelihood of the training data set:

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\underbrace{-\log p_\theta(\mathbf{x}|\mathbf{b})}_{\mathcal{L}_\theta(\mathbf{x})} \right] \quad (8)$$

where the negative log-posterior $\mathcal{L}_\theta(\mathbf{x})$ is specified by (6). Evaluating the gradient of the cost function in (8), we obtain:

$$\begin{aligned} \nabla_\theta \mathcal{L}_\theta(\mathbf{x}) &= \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} [\nabla_\theta \mathcal{C}_\theta(\mathbf{x}; \mathbf{A}, \mathbf{b})] + \nabla_\theta \log \tilde{Z}_\theta \\ &= \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} [\nabla_\theta \mathcal{C}_\theta(\mathbf{x}; \mathbf{A}, \mathbf{b})] - \\ &\quad \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x}|\mathbf{b})} [\nabla_\theta \mathcal{C}_\theta(\mathbf{x}; \mathbf{A}, \mathbf{b})]. \end{aligned} \quad (9)$$

In the second step, chain rule is used to simplify $\nabla_\theta \log \tilde{Z}_\theta$ using (7) as in [21]. Here, $\mathbf{x} \sim p_\theta(\mathbf{x}|\mathbf{b})$ are termed as *fake samples*, drawn from the parametric posterior distribution $p_\theta(\mathbf{x}|\mathbf{b})$. For simplicity, we denote the fake samples as $\mathbf{x}^- \sim p_\theta(\mathbf{x}|\mathbf{b})$, while the reference samples are termed as *true samples*, denoted by $\mathbf{x}^+ \sim q(\mathbf{x})$.

Thus, the ML estimation of θ is equivalent to minimization of the loss:

$$\begin{aligned} \mathcal{L}'(\theta) &= \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \mathcal{C}_\theta(\mathbf{x}; \mathbf{A}, \mathbf{b}) - \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x}|\mathbf{b})} \mathcal{C}_\theta(\mathbf{x}; \mathbf{A}, \mathbf{b}) \\ &\approx \left(\sum_{i=1}^r \mathcal{C}_\theta(\mathbf{x}_i^+; \mathbf{A}, \mathbf{b}) - \sum_{j=1}^r \mathcal{C}_\theta(\mathbf{x}_j^-; \mathbf{A}, \mathbf{b}) \right) \end{aligned} \quad (10)$$

where we consider r samples. Intuitively, the training strategy (10) seeks to decrease the energy of the true samples $\mathcal{C}_\theta(\mathbf{x}^+; \mathbf{A}, \mathbf{b})$ and increase the energy of the fake samples $\mathcal{C}_\theta(\mathbf{x}^-; \mathbf{A}, \mathbf{b})$. Thus, this approach may be seen as an adversarial training scheme similar to [22], where the \mathcal{C}_θ serves as the classifier as shown in Fig. 1. As in GAN models, the algorithm converges when the fake samples are identical in distribution to the training samples; i.e., $\mathcal{C}_\theta(\mathbf{x}^+; \mathbf{A}, \mathbf{b}) \approx \mathcal{C}_\theta(\mathbf{x}^-; \mathbf{A}, \mathbf{b})$. Unlike GAN models that use a separate generator to create the fake samples, the production of the fake samples in the proposed approach also relies on \mathcal{C}_θ , as described in the next subsection.

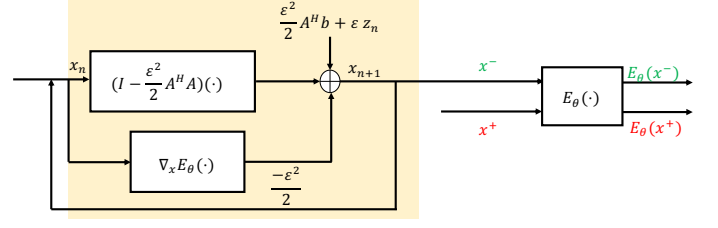


Figure 1: Demonstration of training procedure of DEEPEN. The training procedure determines the optimal weights of the energy $E_\theta(\cdot)$ by minimizing the energy difference between true and fake samples. The true samples \mathbf{x}^+ are obtained from the training data, while the fake samples \mathbf{x}^- are generated using the Langevin sampling algorithm, highlighted by the yellow box. We note that the intermediate results are not stored to evaluate the loss's gradient; therefore, a single physical layer is used for forward propagation. This keeps the training memory demand low.

C. Generation of samples from posterior

We generate the fake samples $\mathbf{x}^- \sim p_\theta(\mathbf{x}|\mathbf{b})$ in (5) using Langevin Markov Chain Monte Carlo (MCMC) method, which only requires the gradient of $\log p_\theta(\mathbf{x}|\mathbf{b})$ w.r.t. \mathbf{x} :

$$\begin{aligned} \mathbf{x}_{n+1} &= \mathbf{x}_n - \frac{\epsilon^2}{2} \nabla_{\mathbf{x}} \mathcal{C}_\theta(\mathbf{x}_n; \mathbf{A}, \mathbf{b}) + \epsilon \mathbf{z}_n \\ &= \mathbf{x}_n - \frac{\epsilon^2}{2} (\mathbf{A}^H (\mathbf{A} \mathbf{x}_n - \mathbf{b}) + \nabla_{\mathbf{x}} \mathcal{E}_\theta(\mathbf{x}_n)) + \epsilon \mathbf{z}_n \end{aligned} \quad (11)$$

where $\epsilon > 0$ is the step-size, $\mathbf{z}_n \sim \mathcal{N}(0, \mathbf{I})$, and \mathbf{x}_0 is drawn randomly from a zero mean Gaussian distribution. Note that $\nabla_{\mathbf{x}} \mathcal{E}_\theta(\mathbf{x})$ is the gradient of the energy model and does not depend on the normalization constant \tilde{Z}_θ .

We now consider the gradient of (10) w.r.t. θ :

$$\begin{aligned} \nabla_\theta \mathcal{L}'(\theta) &\approx \sum_{i=1}^r \nabla_\theta \mathcal{C}_\theta(\mathbf{x}_i^+; \mathbf{A}, \mathbf{b}) - \sum_{j=1}^r \nabla_\theta \mathcal{C}_\theta(\mathbf{x}_j^-; \mathbf{A}, \mathbf{b}) \\ &= \sum_{i=1}^r \nabla_\theta \mathcal{E}_\theta(\mathbf{x}_i^+) - \sum_{j=1}^r \nabla_\theta \mathcal{E}_\theta(\mathbf{x}_j^-) \\ &\quad - \sum_{j=1}^r \left(\nabla_{\mathbf{x}} \mathcal{C}_\theta(\mathbf{x}; \mathbf{A}, \mathbf{b})|_{\mathbf{x}_j^-} \right) \cdot \nabla_\theta (\mathbf{x}_j^-) \end{aligned} \quad (12)$$

In the second step, we used the chain rule to expand $\nabla_\theta \mathcal{E}_\theta(\mathbf{x}^-)$. In addition, we use the fact that the first term in (6) is independent of θ . We note that the last term in the second equation requires unrolling of the Langevin iterations in (11). However, we note that when $\epsilon \rightarrow 0$, the Langevin sampling in (11) simplifies to a gradient descent, which converges to the minimum of $\mathcal{C}_\theta(\mathbf{x}; \mathbf{A}, \mathbf{b})$. Thus, we have $\left(\nabla_{\mathbf{x}} \mathcal{C}_\theta(\mathbf{x}; \mathbf{A}, \mathbf{b})|_{\mathbf{x}_j^-} \right) \approx 0$

as $\epsilon \rightarrow 0$. This implies that the last term in (12) can be ignored, eliminating the need for unrolling.

Therefore we obtain

$$\nabla_{\theta} \mathcal{L}(\theta) \approx \nabla_{\theta} \mathcal{E}_{\theta}(x^+) - \nabla_{\theta} \mathcal{E}_{\theta}(x_{\theta}^-) \quad (13)$$

which does not require backpropagation through the Langevin iterations, thus being efficient from a computational and memory perspective. The Langevin generation path gradient is usually not used in maximum likelihood training of EBMs [21]. The training procedure is summarized in Fig. 1.

D. Image recovery

Once training is complete, images can be derived from the posterior using Langevin sampling, specified by (11).

The proposed formulation also facilitates the estimation of MAP, where we minimize (6) w.r.t. \mathbf{x} . In this work, we use the majorization minimization (MM) framework [29] to derive the MAP estimate, which is guaranteed to converge monotonically [29,30] to a minimum of (6). The MM framework consists of two steps [29]. First, a surrogate function $g(\mathbf{x}|\mathbf{x}_n)$ is constructed such that $L_{\theta}(\mathbf{x}) \leq g(\mathbf{x}|\mathbf{x}_n)$. Next, the surrogate function is minimized to get the next iterate. Similar to [11], we used the following quadratic surrogate function to majorize $L_{\theta}(\mathbf{x})$ in (6):

$$g(\mathbf{x}|\mathbf{x}_n) = \frac{\|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2}{2} + \mathcal{E}_{\theta}(\mathbf{x}_n) + \frac{L}{2}\|\mathbf{x} - \mathbf{x}_n\|^2 + \text{Re}(\mathbf{H}_{\theta}(\mathbf{x}_n)^H(\mathbf{x} - \mathbf{x}_n)) \quad (14)$$

where L is the Lipschitz constant of \mathcal{E}_{θ} and is approximately estimated using CLIP [31]. The above surrogate function has the following closed-form solution:

$$\mathbf{x}_{n+1} = (\mathbf{A}^H \mathbf{A} + L\mathbf{I})^{-1} (\mathbf{A}^H \mathbf{b} + L\mathbf{x}_n - \mathbf{H}_{\theta}(\mathbf{x}_n)) \quad (15)$$

For a large-dimensional \mathbf{x} , inverting the \mathbf{A} operator is computationally expensive. Therefore, we use the conjugate gradient algorithm to obtain the next iterate \mathbf{x}_{n+1} . The following result [30] shows that the proposed algorithm converges monotonically to a stationary point of the cost function (6).

Lemma 2.1: Consider the cost function $\mathcal{L}_{\theta}(\mathbf{x})$ in (6), which is bounded below by zero¹. Then the sequence of iterates $\{\mathbf{x}_n\}$ generated by MM algorithm in (15) will converge to a stationary point of (6).

Proof: Note that $L_{\theta}(\mathbf{x}) \geq 0; \forall \mathbf{x} \in \mathbb{C}^m$ and hence it is lower bounded by a finite value. Moreover, the surrogate function satisfies $g(\mathbf{x}_n|\mathbf{x}_n) = f(\mathbf{x}_n)$ and

¹The CNN implementation $\mathcal{E}_{\theta}(\mathbf{x})$ has an absolute function in the output layer, which makes the lower bound zero.

$g(\mathbf{x}|\mathbf{x}_n) \geq f(\mathbf{x})$. Then, using Theorem 1 from [30], the MM algorithm in (15) will converge to a stationary point of (6). ■

III. MAP EXPERIMENTS & RESULTS

A. Dataset

We compare the proposed DEEPEN approach with the SOTA methods described in the context of multichannel MR image reconstruction. The forward operator is defined as $\mathbf{A} = \mathbf{S}\mathbf{F}\mathbf{C}$, where \mathbf{S} is the sampling matrix, \mathbf{F} is the Fourier matrix, and \mathbf{C} is the Coil Sensitivity Map (CSM) that is estimated using [32]. MR images were obtained from the publicly available multichannel fastMRI brain dataset [33]. It is a 12-channel brain dataset and consists of complex images of size 320×320 . The data set was divided into 45 training, 5 validation and 50 test subjects, each consisting of approximately 450, 50, and 500 images, respectively. We evaluated the models on T2-weighted images using 2D and 1D undersampling masks for different acceleration factors. We now describe the details of the proposed method as well as the SOTA methods.

B. Implementation details of the algorithms

1) **DEEPEN:** We implement the energy model $\mathcal{E}_{\theta}(\mathbf{x}) : \mathbb{C}^m \rightarrow \mathbb{R}^+$ as a CNN consisting of five 3×3 convolutional layers followed by a linear layer. Each convolutional layer consists of 64 channels, and a Rectified Linear Unit (ReLU) was used between each layer, except the last linear layer. An absolute activation function was used in the linear layer to ensure that $\mathcal{E}_{\theta}(\mathbf{x})$ is lower bounded by a finite value. The score function $\nabla_{\mathbf{x}} \mathcal{E}_{\theta}(\mathbf{x})$ was evaluated using the chain rule.

To ensure stable ML training, similar to [34], we found it useful to add zero-mean Gaussian noise with standard deviation $2\epsilon^2$ to the training data. We also found it beneficial to scale the posterior $p_{\theta}(\mathbf{x}|\mathbf{b})$ by $\frac{\epsilon^2}{2}$, which gives rise to the equivalent Langevin MCMC update [34]:

$$\mathbf{x}_{n+1} = \mathbf{x}_n - (\mathbf{A}^H(\mathbf{A}\mathbf{x}_n - \mathbf{b}) + \nabla_{\mathbf{x}} \mathcal{E}_{\theta}(\mathbf{x}_n)) + \epsilon \mathbf{z}_n \quad (16)$$

The MCMC Langevin algorithm was initialized with zero-mean Gaussian Noise and 100 MCMC sampling iterations were performed to generate the fake samples.

2) **PnP MuSE:** We compare the proposed method with MuSE energy model that is trained in a PnP fashion. In particular, a single EBM is trained using DSM technique to predict Gaussian noise corresponding

to a range of noise standard deviations. We used the following network architecture for MuSE:

$$E_{\theta}^{\text{MuSE}}(\mathbf{x}) = \frac{1}{2} \|\mathbf{x} - \psi_{\theta}(\mathbf{x})\|^2 \quad (17)$$

where $\psi_{\theta}(\mathbf{x})$ is a five layer convolutional network. Each layer was followed by ReLU activation, except in the final layer. The convolution layer consists of a 3×3 filter with 64 channels. MuSE was trained to predict Gaussian noise with standard deviations ranging from 0 to 0.1.

3) *PnP-ISTA*: PnP-ISTA is motivated by the iterative soft thresholding approach [35] used in CS algorithm. The proximal of the regularizer in the ISTA algorithm is replaced by a CNN denoiser, which is pre-trained as a Gaussian denoiser [36]. The optimization algorithm alternates between the following steps:

$$\mathbf{q}_t = \mathbf{x}_{t-1} - \alpha \frac{\mathbf{A}^H(\mathbf{A}\mathbf{x}_{t-1} - \mathbf{b})}{\eta^2} \quad (18)$$

$$\mathbf{x}_t = D_{\sigma}(\mathbf{q}_t) \quad (19)$$

where $D_{\sigma}(\cdot) : \mathbb{C}^m \rightarrow \mathbb{C}^m$ is a five-layer CNN-based denoiser with each layer consisting of 64 channels with a 3×3 filter. Except for the final layer, after each layer the ReLU activation function was used. The denoiser was trained to remove Gaussian noise with standard deviation $\sigma = 0.01$.

4) *E2E MAP learning using CNN denoisers*: We also compare EBM with an E2E approach MoL [8], which approximate the MAP estimate and are trained to minimize the MSE between the recovered images and the reference images.

$$L_{\text{MSE}}(\theta) = \sum_{i=1}^{N_{\text{samples}}} \|\mathbf{x}_{\theta}^*(i) - \mathbf{x}_{\text{ref}}(i)\|^2 \quad (20)$$

Here, N_{samples} denotes the number of training samples, $\mathbf{x}_{\text{ref}}(i)$ is the reference image, and $\mathbf{x}_{\theta}^*(i)$ is the solution to a regularized optimization problem that uses the CNN as a regularizer [8]. The gradient of (20) is specified by

$$\nabla_{\theta} L_{\text{MSE}}(\theta) = \sum_{i=1}^{N_{\text{samples}}} \underbrace{(\mathbf{x}_{\theta}^*(i) - \mathbf{x}_{\text{ref}}(i))}_{l_i} \cdot \nabla_{\theta} \mathbf{x}_{\theta}^* \quad (21)$$

Unlike the last term in (12), l_i is not zero in the above case. The DEQ approach [8,18,19] is used to derive the above gradient, which assumes the iterative algorithm to derive \mathbf{x}_{θ}^* to converge to a fixed point. A monotone constraint is placed on the CNN to ensure convergence to the fixed point. We used a five-layer CNN to implement MoL and the monotone constraint was imposed using a log-barrier approach as described in [8].

5) *E2E MAP using CNN energy models*: The DEEPEN framework learns the posterior distribution using (8) in an E2E fashion, which can be used for posterior sampling or to derive MAP estimates. In contrast, ELDER in [20] learns the regularizer parameters in an E2E fashion by minimizing the loss of MSE in (20), similar to the MoL approach described in the above section. In this work, we consider the following iterative algorithm for K iterations:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha (\mathbf{A}^H(\mathbf{A}\mathbf{x}_k - \mathbf{b}) + \nabla_{\mathbf{x}_k} \mathcal{E}_{\theta}(\mathbf{x}_k)) \quad (22)$$

to obtain the MAP estimate. Here, α is the step-size and is a learnable parameter.

C. Visualization of MuSE and DEEPEN energies

In this section, we compare pre-trained MuSE with DEEPEN, whose energy model is trained in E2E fashion. For both energy models, we fed images of the form $\tilde{\mathbf{x}} = \mathbf{x} + \alpha_z \mathbf{z} + \alpha_s \mathbf{s}$, where \mathbf{x} is the test image, $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$ is the Gaussian noise perturbation, $\mathbf{s} = (\mathbf{x}_0 - \mathbf{x})$ is the structural artifact perturbation and \mathbf{x}_0 is the SENSE solution given as $\mathbf{x}_0 = (\mathbf{A}^H \mathbf{A} + \tilde{\lambda} \mathbf{I})^{-1} \mathbf{A}^H \mathbf{b}$. An example of the test image, Gaussian and structural perturbation is shown in the top rows of Fig.2.a and Fig.2.b for four-fold and six-fold acceleration, respectively. The second and third rows of each figure show the MuSE and DEEPEN energy plot as a function of α_s and α_z for four-fold and six-fold acceleration, respectively. The red cross on the plot indicates (α_s^*, α_z^*) for which the energy function achieves a minimum value. Note that, a well-trained energy function will have a lower energy value for good images.

We observe from the figure that unlike MuSE, the energy model trained via DEEPEN has (α_s^*, α_z^*) closer to zero, for both acceleration settings. This implies that the DEEPEN energy model is a better discriminator between the reference and the image with correlated perturbations. We also show the image $\hat{\mathbf{x}} = \mathbf{x} + \alpha_s^* \mathbf{s} + \alpha_z^* \mathbf{z}$ at minimum in the second and third rows of Fig.2 and the corresponding error images. We observe that DEEPEN offers an image that is closer to the reference image. The difference between the minimizer and reference can be appreciated from the error images, especially for six-fold acceleration, where the error image of MuSE has a significant structural perturbation compared to DEEPEN. This shows that DEEPEN is a better regularizer and, consequently, when employed in inverse problems, will be efficient in removing structural perturbations.

D. Reconstruction

In this section, we compare the reconstruction performance of DEEPEN with MuSE, ELDER, PnP-ISTA,

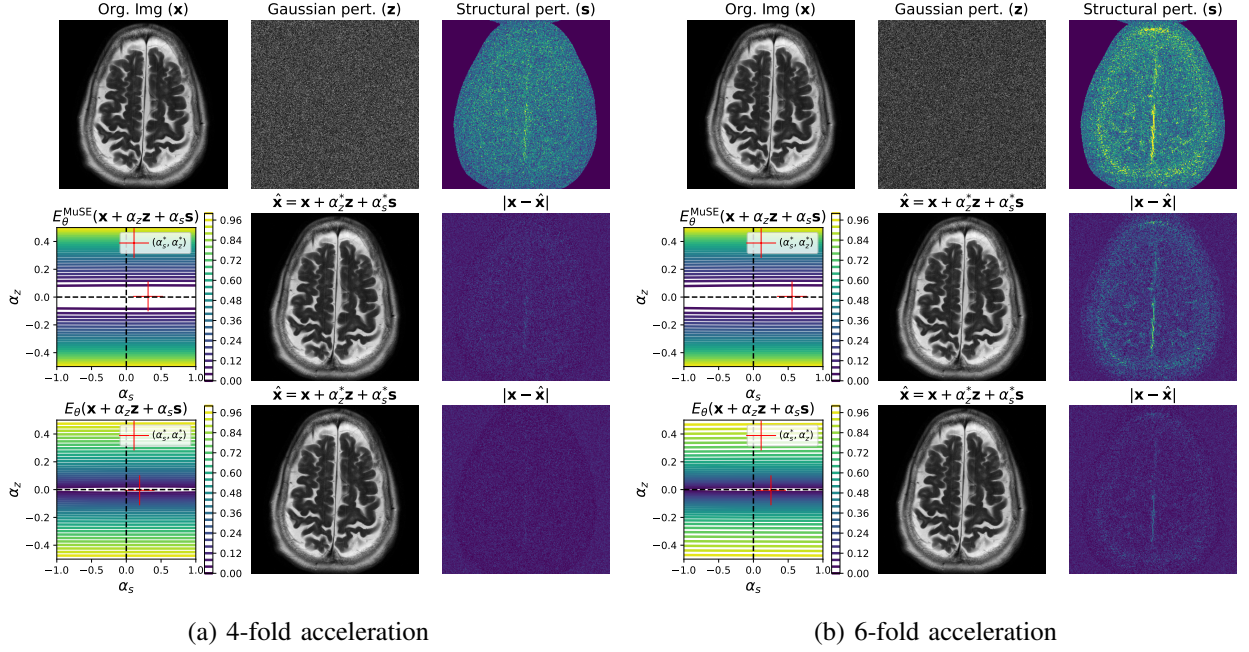


Figure 2: Comparison of pre-trained (MuSE) with E2E-trained (DEEPEN) energy models for (a) four-fold and (b) six-fold acceleration. The top row in each figure shows the original image, Gaussian noise, and the structural perturbation specified by $x_0 - x$, where x_0 denotes the sense solution. The second and third row in each figure shows the plot of MuSE and DEEPEN energy as a function of α_z and α_s , their corresponding reconstructed and the error images, respectively. The images are reconstructed by taking the combination of the form $\hat{x} = x + \alpha_z^* z + \alpha_s^* s$, where (α_s^*, α_z^*) are the minimizer (indicated by cross mark in the contour plot) of the energy function. We note that the minimum of the DEEPEN energy is closer to $\alpha_s^* \approx 0$; $\alpha_z^* = 0$, with the differences $x^* - x$ smaller than that of MuSE. This shows that the DEEPEN energy is effective in suppressing both correlated structural perturbations as well as Gaussian noise.

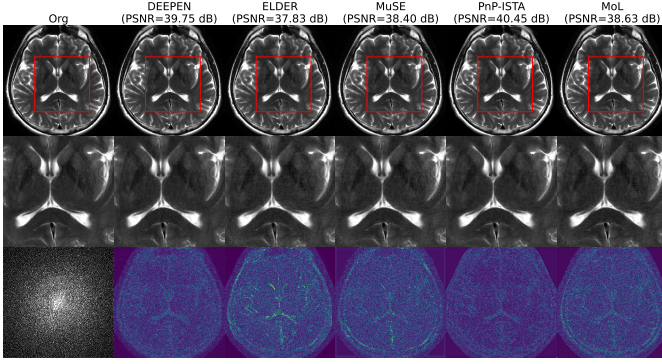
Table I: Reconstruction performance of different models for four different acquisition settings.

Algorithm	2D 4x Mask		2D 6x Mask		1D 2x Mask		1D 4x Mask	
	Avg. PSNR	SSIM	Avg. PSNR	SSIM	Avg. PSNR	SSIM	Avg. PSNR	SSIM
DEEPEN	39.15+/-1.43	0.98	37.18+/-1.32	0.975	38.91+/-2.21	0.98	31.24+/-1.64	0.93
MuSE	38.27+/-1.30	0.97	36.64+/-1.20	0.96	38.37+/-1.88	0.97	31.66+/-1.75	0.93
ISTA	39.88+/-1.38	0.97	36.76+/-5.15	0.96	39.91+/-2.93	0.98	29.04+/-10.80	0.93
MoL	38.12+/-1.27	0.97	36.91+/-1.19	0.97	37.72 +/-1.67	0.97	31.35+/-1.38	0.926
ELDER	38.66+/-1.23	0.98	37.39+/-1.17	0.97	38.02+/-1.52	0.97	31.19+/-1.21	0.93

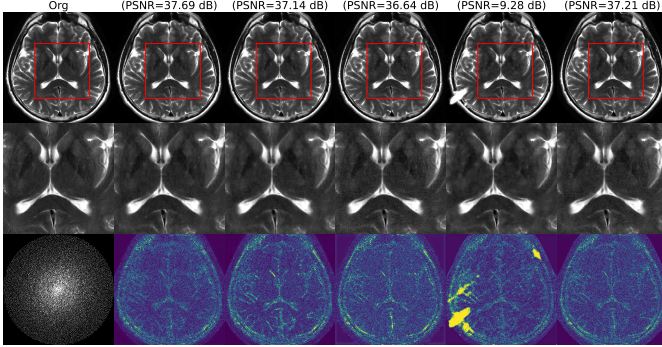
and MoL. The energy-based models DEEPEN and MuSE were run until the cost function satisfied $|L_\theta(x_{n+1}) - L_\theta(x_n)|/|L_\theta(x_n)| \leq 1e^{-6}$ or until 500 iterations were reached. The optimal hyperparameters needed to run MuSE were obtained from [11]. We use the update equation in (22) as the inference algorithm for ELDER trained with MSE loss with $K = 30$. PnP-ISTA was run until $|x_{n+1} - x_n|/|x_n| \leq 1e^{-6}$ or until 500 iterations were reached.

All reconstruction algorithms were initialized with SENSE. Table. I compares the reconstruction performance for four different settings: 4x and 6x acceleration with 2D undersampling mask; 2x and 4x acceleration with 1D undersampling mask. The reconstructions are

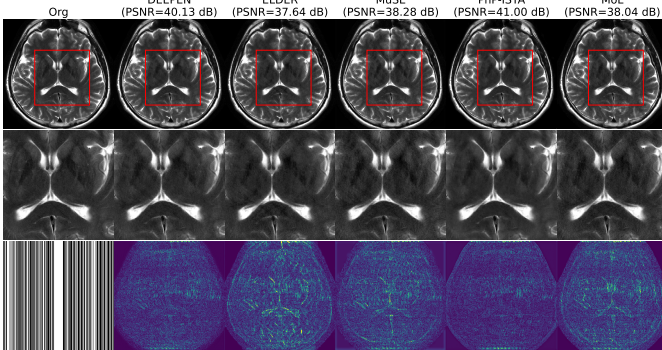
compared using two metrics: Peak Signal-to-Noise Ratio (PSNR) and Structural Similarity Index Measure (SSIM). From the table, we observe that DEEPEN offers better results compared to MuSE. The improved performance of DEEPEN can be attributed to the training strategy, which ensures that the true samples correspond to the minimum of energy, with it increases in other directions. In particular, the energy is higher for both Gaussian and structural perturbations (as demonstrated in Fig. 2), compared to MuSE. However, we note that both energy models learn a negative log-prior distribution that ensures convergence (without a contraction constraint on the score function), as discussed in Lemma 2.1. We note that the performance of DEEPEN is comparable



(a) Four-fold acceleration using 2D undersampling mask



(b) Six-fold acceleration using 2D undersampling mask



(c) Two-fold acceleration using 1D undersampling mask

Figure 3: Comparison of DEEPEN, ELDER, MuSE, PnP-ISTA, and MoL for three different acquisition settings on the fastMRI brain data set. The first, second, and the third row in each figure shows the reconstructed image, enlarged image, and the error image, respectively. The error image is scaled by a factor of 10 to highlight the differences. The first image in the third row shows the undersampling mask.

to that of PnP-ISTA at lower undersampling rates (e.g. 2D mask 4-fold setting). PnP-ISTA was implemented without contraction constraints, which offers higher performance than the implementation with constraints [11]. We also observe that the performance of PnP-ISTA drops significantly at higher accelerations, which is consistent with the observation in [11]. In particular, PnP-ISTA exhibits localized artifacts (bright regions in Fig. 3.(b))

due to convergence problems. We observe that DEEPEN offers improved performance than ELDER for the 2D 4-fold and 1D 2-fold acquisition settings, while the performance is comparable for other settings.

E. Generalization performance of E2E energy models

DEEPEN and ELDER both use energy models, which are trained in E2E fashion. The main difference is that DEEPEN learns the posterior using the maximum likelihood approach, while ELDER relies on MSE loss. We compare the generalization performance by considering a different acquisition scheme (sampling pattern) during the test setting from what was assumed during training. Table II compares the performance of DEEPEN and ELDER for different acquisitions. In the table, we report the Avg. PSNR, where we have boldfaced the performance of DEEPEN. From the table, we observe that when there is a change in the acquisition setting, the performance of ELDER drops drastically compared to DEEPEN. For example, when the energy models are trained on a 4-fold 2D undersampling mask and tested in the same setting, the performance of DEEPEN and ELDER is about 39.15 dB and 38.66 dB, respectively. However, when tested on a 2-fold 1D mask, the performance of ELDER drops by about 2.38 dB while the performance of DEEPEN drops only by 0.92 dB. Fig. 4 illustrates the generalization performance of DEEPEN for two different mismatch scenarios.

We attribute the improved generalization performance of DEEPEN to the adversarial training strategy. We note that the fake samples are generated by Langevin dynamics, where Gaussian noise is added at every iteration to generate the fake samples. This training strategy ensures that the energy function is well-trained along different random paths from the initialization. This approach also encourages the energy function to have well-defined minima in the training samples, with the energy values increasing in all directions, as seen in Fig. 2. By contrast, training with the MSE loss learns a single path from the initialization and the reference image. We note that the iterations during inference depend on the energy function as well as the forward model; a change in the forward model can result in a different path to the samples, which may make E2E methods trained using MSE loss less robust to changes in acquisition settings.

IV. POSTERIOR SAMPLING EXPERIMENTS

Most of the current E2E deep learning methods focus on learning the MAP estimate [8,15,16,20]. In contrast, DEEPEN learns the posterior distribution in an E2E fashion, which enables us to sample the distribution using (16). We note that the long sampling chain in

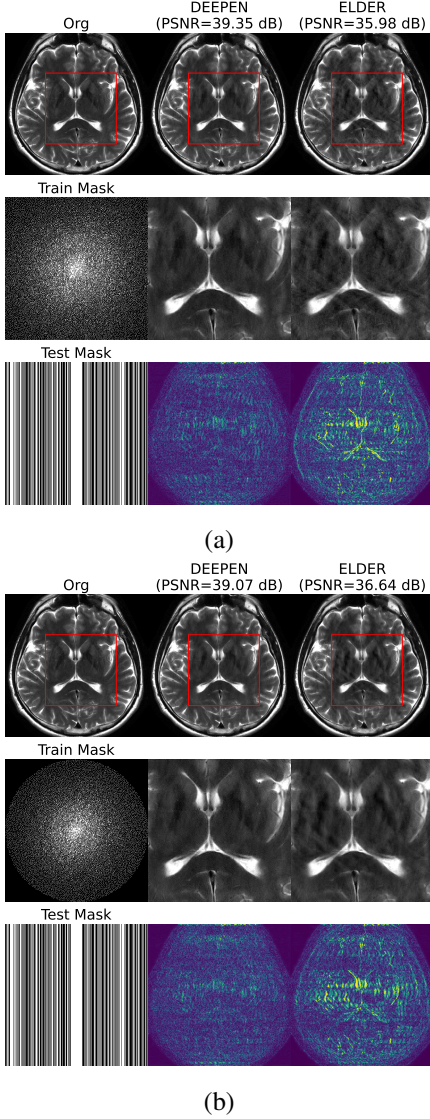


Figure 4: Generalization comparison of E2E-trained DEEPEN and ELDER models for two different settings: (a) models are trained on 4-fold 2D undersampling mask and tested on a 2-fold 1D undersampling mask. When compared to Fig. 3.a, which demonstrates the reconstruction performance of models trained and tested on 4-fold 2D mask, the performance of DEEPEN and ELDER drops by 0.4 dB and 1.85 dB, respectively (b) models are trained on 6-fold 2D mask and tested on 2-fold 1D mask. When compared to Fig. 3.b, which demonstrates the performance of models trained and tested on 6-fold 2D mask, the performance of DEEPEN improved by 1.38 dB while the performance of ELDER dropped by 0.5 dB.

diffusion models and the time dependence of the score model make it challenging to train diffusion models in an E2E fashion. We compare the sampling performance of DEEPEN with Diffusion Posterior Sampling (DPS) [24], where the diffusion model is learned in a PnP fashion

in the T2-weighted fastMRI brain dataset. We also note that it is challenging for diffusion models to realize MAP estimate. In particular, score models requires the computation of the following complex integral to estimate the log-prior [37]:

$$\log p_{\theta}(\mathbf{x}) = \int_0^T \nabla \cdot s_{\theta}(\mathbf{x}(t), t) dt + \log p_{\pi}(\mathbf{x}_T) \quad (23)$$

where $\log p_{\pi}(\mathbf{x}_T)$ represents the final Gaussian noise distribution, $s_{\theta}(\mathbf{x}(t), t)$ represents the time conditional score model at scale t , and $\nabla \cdot$ is the divergence operator. This makes it computationally expensive to realize a MAP estimate using diffusion models.

A. Architecture and implementation

1) *DEEPEN*: For sampling, we used the same E2E trained network discussed in the previous section. Langevin algorithm in (16) was used to generate 100 different samples. The algorithm was initialized with zero-mean Gaussian noise and ran for 100 iterations.

2) *Diffusion models*: A time-conditional score model was trained using the loss in [37], where the noise was added to the data according to the following perturbation kernel:

$$p_{0t}(\mathbf{x}(t)|\mathbf{x}(0)) = \mathcal{N}(\mathbf{x}(t); \mathbf{x}(0), \frac{1}{2 \log \Lambda} (\Lambda^{2t} - 1) \mathbf{I}) \quad (24)$$

where t is uniformly sampled over $[0, T]$, $\{\mathbf{x}(t)\}_{t=0}^T$ represents the diffusion process, $p_{0t}(\mathbf{x}(t)|\mathbf{x}(0))$ represents the transition probability from $\mathbf{x}(0)$ to $\mathbf{x}(t)$, and $\mathbf{x}(0)$ are the training data samples. The architecture of the score model was chosen as a time-conditional DRUnet with 64, 128, 256 and 512 channels. We used the sampling algorithm proposed in [24] with 1000 iterative steps to sample the posterior distribution. We note that DEEPEN requires five times fewer number of parameters than the time-conditional score model.

B. Sampling illustration

Fig. 5 illustrates the sampling performance of DEEPEN and DPS for different acquisition schemes. The upper row of each figure shows the MAP, the minimum MSE (MMSE), and the uncertainty estimates provided by the DEEPEN algorithm. The MMSE and the uncertainty estimate are obtained by taking the mean and variance of the generated samples, respectively. The second row of each figure in Fig. 5 shows the original image, MMSE, and the uncertainty estimates of the DPS algorithm. We do not show the MAP estimate as it is challenging to realize a MAP estimate using diffusion models, as discussed above.

Table II: Generalization performance of E2E-trained energy model: DEEPEN(boldfaced) and ELDER.

		Train setting		
		4x 2D mask	6x 2D mask	2x 1D mask
Test setting	4x 2D mask	39.15 + / - 1.43	37.60 + / - 1.53	38.64 + / - 1.55
		38.66 + / - 1.23	38.88 + / - 1.24	37.86 + / - 1.20
	6x 2D mask	37.50 + / - 1.29	37.18 + / - 1.32	37.20 + / - 1.37
		36.82 + / - 1.18	37.39 + / - 1.17	36.26 + / - 1.14
	2x 1D mask	38.23 + / - 1.99	37.66 + / - 1.98	38.91 + / - 2.21
		36.28 + / - 1.43	36.54 + / - 1.46	38.02 + / - 1.52

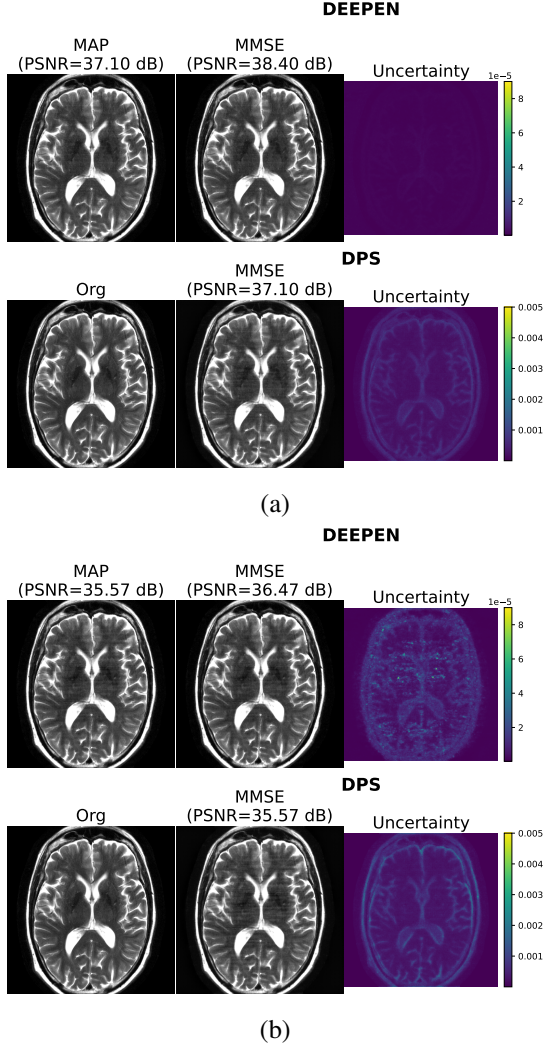


Figure 5: MAP, MMSE, and uncertainty estimate provided by DEEPEN algorithm for two different MRI acquisition settings. We compare DEEPEN’s sampling performance with DPS when 2D undersampling mask is employed for (a) 4-fold and (b) 6-fold acceleration.

We note that DEEPEN requires 10x fewer sampling steps compared to DPS and uses a network with 5x fewer parameters, while offering comparable results. The narrower unimodal posterior distribution results in faster sampling and reduced uncertainty. This can be seen in Fig. 5.a and Fig. 5.b which shows the uncertainty estimates for a 4-fold 2D and 6-fold 2D undersampling

mask, respectively. Deterministic ODE flows [25,26] are designed to accelerate prior sampling by learning straighter paths between the source and target distributions. However, these models still require long chains for posterior sampling. A key challenge is that they are trained solely along the paths between the distributions, necessitating the addition of noise after each DC gradient descent step to project back to the originally trained regions [26]. By contrast, DEEPEN gradients are efficient in removing both correlated and Gaussian perturbations, which we hypothesize to be another reason for the faster sampling.

We note that a benefit of the diffusion setting is that the same trained model can be re-used in multiple settings. In contrast, DEEPEN uses an E2E approach, which requires a customized model for each acquisition setting. However, the proposed E2E training strategy is more computationally and memory-efficient than traditional E2E models. In addition, the DEEPEN model is more generalizable than traditional E2E models to changes in the acquisition setting. Furthermore, it offers better performance and reduced uncertainty than diffusion models, even though it uses 10x fewer sampling steps and CNN models with 5x fewer parameters.

V. CONCLUSION

We proposed DEEPEN, an E2E training framework to learn the posterior probability distribution for imaging inverse problems. DEEPEN can be used to derive the MAP estimate and sample from the posterior distribution. The proposed E2E model does not need to be unrolled, which results in a more memory-efficient training procedure. Unlike PnP and DEQ methods, this approach does not require a Lipschitz constraint for convergence, which translates to improved representation power and thus image quality. Our experiments show that the maximum likelihood training strategy offers a better-defined energy landscape, translating to improved image recovery algorithms. The experiments also show that the proposed approach can provide MAP estimates with improved image quality compared to E2E methods while being superior to PnP methods. We observe that the E2E approach can offer 10x faster sampling with

models with reduced complexity compared to diffusion models, despite offering comparable reconstructions with reduced uncertainty.

REFERENCES

- [1] A. N. Tikhonov and V. I. Arsenin, *Solutions of Ill-posed Problems: Andrey N. Tikhonov and Vasilii Y. Arsenin. Translation Editor Fritz John.* Wiley, 1977.
- [2] D. L. Donoho, “Compressed sensing,” *IEEE Transactions on information theory*, vol. 52, no. 4, pp. 1289–1306, 2006.
- [3] E. J. Candès, J. Romberg, and T. Tao, “Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information,” *IEEE Transactions on information theory*, vol. 52, no. 2, pp. 489–509, 2006.
- [4] S. V. Venkatakrishnan, C. A. Bouman, and B. Wohlberg, “Plug-and-play priors for model based reconstruction,” in *2013 IEEE Global Conference on Signal and Information Processing*, 2013, pp. 945–948.
- [5] G. T. Buzzard, S. H. Chan, S. Sreehari *et al.*, “Plug-and-play unplugged: Optimization-free reconstruction using consensus equilibrium,” *SIAM Journal on Imaging Sciences*, vol. 11, no. 3, pp. 2001–2020, 2018.
- [6] R. Ahmad, C. A. Bouman, G. T. Buzzard *et al.*, “Plug-and-play methods for magnetic resonance imaging: Using denoisers for image recovery,” *IEEE signal processing magazine*, vol. 37, no. 1, pp. 105–116, 2020.
- [7] E. Ryu, J. Liu, S. Wang *et al.*, “Plug-and-play methods provably converge with properly trained denoisers,” in *International Conference on Machine Learning*. PMLR, 2019, pp. 5546–5557.
- [8] A. Pramanik, M. B. Zimmerman, and M. Jacob, “Memory-efficient model-based deep learning with convergence and robustness guarantees,” *IEEE Transactions on Computational Imaging*, vol. 9, pp. 260–275, 2023.
- [9] R. Cohen, Y. Blau, D. Freedman *et al.*, “It has potential: Gradient-driven denoisers for convergent solutions to inverse problems,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 18 152–18 164, 2021.
- [10] S. Hurault, A. Leclaire, and N. Papadakis, “Gradient step denoiser for convergent plug-and-play,” *arXiv preprint arXiv:2110.03220*, 2021.
- [11] J. R. Chand and M. Jacob, “Multi-scale energy (muse) framework for inverse problems in imaging,” *IEEE transactions on computational imaging*, 2024.
- [12] P. Vincent, “A connection between score matching and denoising autoencoders,” *Neural computation*, vol. 23, no. 7, pp. 1661–1674, 2011.
- [13] V. Monga, Y. Li, and Y. C. Eldar, “Algorithm unrolling: Interpretable, efficient deep learning for signal and image processing,” *IEEE Signal Processing Magazine*, vol. 38, no. 2, pp. 18–44, 2021.
- [14] J. Sun, H. Li, Z. Xu *et al.*, “Deep admm-net for compressive sensing mri,” *Advances in neural information processing systems*, vol. 29, 2016.
- [15] K. Hammernik, T. Klatzer, E. Kobler *et al.*, “Learning a variational network for reconstruction of accelerated mri data,” *Magnetic resonance in medicine*, vol. 79, no. 6, pp. 3055–3071, 2018.
- [16] H. K. Aggarwal, M. P. Mani, and M. Jacob, “Modl: Model-based deep learning architecture for inverse problems,” *IEEE transactions on medical imaging*, vol. 38, no. 2, pp. 394–405, 2018.
- [17] T. Küstner, N. Fuin, K. Hammernik *et al.*, “Cinenet: deep learning-based 3d cardiac cine mri reconstruction with multi-coil complex-valued 4d spatio-temporal convolutions,” *Scientific reports*, vol. 10, no. 1, p. 13710, 2020.
- [18] S. Bai, J. Z. Kolter, and V. Koltun, “Deep equilibrium models,” *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [19] D. Gilton, G. Ongie, and R. Willett, “Deep equilibrium architectures for inverse problems in imaging,” *IEEE Transactions on Computational Imaging*, vol. 7, pp. 1123–1133, 2021.
- [20] Z. Zou, J. Liu, B. Wohlberg *et al.*, “Deep equilibrium learning of explicit regularizers for imaging inverse problems,” *arXiv preprint arXiv:2303.05386*, 2023.
- [21] Y. Song and D. P. Kingma, “How to train your energy-based models,” *arXiv preprint arXiv:2101.03288*, 2021.
- [22] I. Goodfellow, J. Pouget-Abadie, M. Mirza *et al.*, “Generative adversarial nets,” *Advances in neural information processing systems*, vol. 27, 2014.
- [23] C. Han, L. Rundo, K. Murao *et al.*, “Madgan: Unsupervised medical anomaly detection gan using multiple adjacent brain mri slice reconstruction,” *BMC bioinformatics*, vol. 22, pp. 1–20, 2021.
- [24] H. Chung, J. Kim, M. T. Mccann *et al.*, “Diffusion posterior sampling for general noisy inverse problems,” *arXiv preprint arXiv:2209.14687*, 2022.
- [25] J. Song, C. Meng, and S. Ermon, “Denoising diffusion implicit models,” *arXiv preprint arXiv:2010.02502*, 2020.
- [26] S. Martin, A. Gagneux, P. Hagemann *et al.*, “Pnp-flow: Plug-and-play image restoration with flow matching,” *arXiv preprint arXiv:2410.02423*, 2024.
- [27] J. R. Chand and M. Jacob, “Memory-efficient deep end-to-end posterior network (deepen) for inverse problems,” in *2024 IEEE International Symposium on Biomedical Imaging (ISBI)*, 2024, pp. 1–5.
- [28] Z. Li, Y. Chen, and F. T. Sommer, “Learning energy-based models in high-dimensional spaces with multiscale denoising-score matching,” *Entropy*, vol. 25, no. 10, p. 1367, 2023.
- [29] Y. Sun, P. Babu, and D. P. Palomar, “Majorization-minimization algorithms in signal processing, communications, and machine learning,” *IEEE Transactions on Signal Processing*, vol. 65, no. 3, pp. 794–816, 2016.
- [30] M. Razaviyayn, M. Hong, and Z.-Q. Luo, “A unified convergence analysis of block successive minimization methods for nonsmooth optimization,” *SIAM Journal on Optimization*, vol. 23, no. 2, pp. 1126–1153, 2013.
- [31] L. Bungert, R. Raab, T. Roith *et al.*, “Clip: Cheap lipschitz training of neural networks,” in *International Conference on Scale Space and Variational Methods in Computer Vision*. Springer, 2021, pp. 307–319.
- [32] M. Uecker, P. Lai, M. J. Murphy *et al.*, “Espirit—an eigenvalue approach to autocalibrating parallel mri: where sense meets grappa,” *Magnetic resonance in medicine*, vol. 71, no. 3, pp. 990–1001, 2014.
- [33] J. Zbontar, F. Knoll, A. Sriram *et al.*, “fastmri: An open dataset and benchmarks for accelerated mri,” *arXiv preprint arXiv:1811.08839*, 2018.
- [34] E. Nijkamp, M. Hill, T. Han *et al.*, “On the anatomy of mcmc-based maximum likelihood learning of energy-based models,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 04, 2020, pp. 5272–5280.
- [35] A. Beck and M. Teboulle, “A fast iterative shrinkage-thresholding algorithm for linear inverse problems,” *SIAM journal on imaging sciences*, vol. 2, no. 1, pp. 183–202, 2009.
- [36] U. S. Kamilov, C. A. Bouman, G. T. Buzzard *et al.*, “Plug-and-play methods for integrating physical and learned models in computational imaging,” *arXiv preprint arXiv:2203.17061*, 2022.
- [37] Y. Song, J. Sohl-Dickstein, D. P. Kingma *et al.*, “Score-based generative modeling through stochastic differential equations,” *arXiv preprint arXiv:2011.13456*, 2020.