

PATCH-BASED DIFFUSION MODELS BEAT WHOLE-IMAGE MODELS FOR MISMATCHED DISTRIBUTION INVERSE PROBLEMS

Anonymous authors

Paper under double-blind review

ABSTRACT

Diffusion models have achieved excellent success in solving inverse problems due to their ability to learn strong image priors, but existing approaches require a large training dataset of images that should come from the same distribution as the test dataset. When the training and test distributions are mismatched, artifacts and hallucinations can occur in reconstructed images due to the incorrect priors. In this work, we systematically study out of distribution (OOD) problems where a known training distribution is first provided. We first study the setting where only a single measurement obtained from the unknown test distribution is available. Next we study the setting where a very small sample of data belonging to the test distribution is available, and our goal is still to reconstruct an image from a measurement that came from the test distribution. In both settings, we use a patch-based diffusion prior that learns the image distribution solely from patches. Furthermore, in the first setting, we include a self-supervised loss that helps the network output maintain consistency with the measurement. Extensive experiments show that in both settings, the patch-based method can obtain high quality image reconstructions that can outperform whole-image models and can compete with methods that have access to large in-distribution training datasets. Furthermore, we show how whole-image models are prone to memorization and overfitting, leading to artifacts in the reconstructions, while a patch-based model can resolve these issues.

1 INTRODUCTION

In image processing, inverse problems are of paramount importance and consist of reconstructing a latent image \mathbf{x} from a measurement $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \varepsilon$. Here, \mathcal{A} represents a forward operator and ε represents random unknown noise. By Bayes' rule, $\log p(\mathbf{x}|\mathbf{y})$ is proportional to $\log p(\mathbf{x}) + \log p(\mathbf{y}|\mathbf{x})$, so obtaining a good prior $p(\mathbf{x})$ is crucial for recovering \mathbf{x} when \mathbf{y} contains far less information than \mathbf{x} . Diffusion models obtain state-of-the-art results for learning a strong prior and sampling from it, so similarly competitive results can be obtained when using them to solve inverse problems (Chung et al., 2022a; 2023a; Song et al., 2024; Wang et al., 2022; Kawar et al., 2021; Li et al., 2023a).

However, training diffusion models well requires vast amounts of clean training data (Song et al., 2021; Ho et al., 2020), which is infeasible to collect in many applications such as medical imaging (Chung et al., 2022b; Song et al., 2022; Jalal et al., 2021), black hole imaging (Feng et al., 2023; 2024), and phase retrieval (Li et al., 2023a; Wu et al., 2019). **In particular, no ground truth images are known for very challenging inverse problems such as black hole imaging (Feng et al., 2023) and Fresnel phase retrieval (Gureyev et al., 2004), so one must obtain a reconstruction from only a single measurement \mathbf{y} .** In other applications such as dynamic CT reconstruction (Reed et al., 2021) and single photo emission CT (Li et al., 2023b), obtaining a high quality measurement that can lead to a reconstruction that closely approximates the ground truth can be slow or potentially harmful to the patient, so only a very small dataset of clean images are available. Thus, in this paper we consider two settings: the *single measurement* setting in which we are given one measurement \mathbf{y} whose corresponding \mathbf{x} belongs to a different distribution from the training dataset, and the *small dataset* setting in which we are only given a small number of samples \mathbf{x} that belong to the same distribution as the test dataset.

In recent years, some works have aimed to address these problems by demonstrating that diffusion models have a stronger generalization ability than other deep learning methods (Jalal et al., 2021), so slight distribution mismatches between the training data and test data may not significantly degrade the reconstructed image quality. However, in cases of particularly compressed or noisy measurements, as well as when the test data is severely out of distribution (OOD) with a significant domain shift, an improper choice of training data leads to an incorrect prior that causes substantial image degradation and hallucinations (Feng et al., 2023; Barbano et al., 2023). To address these challenges in the single measurement case, recent works use each measurement \mathbf{y} to adjust the weights of a diffusion network at reconstruction time Barbano et al. (2023); Chung & Ye (2024), aiming to shift the underlying prior learned by the network toward the appropriate prior for the latent image in the test case. However, as the networks have huge numbers of weights, an intricate and parameter-sensitive refining process of the network is required during reconstruction to avoid overfitting to the measurement. Furthermore, there is still a substantial gap in performance between methods using an OOD prior and methods using an in-distribution prior. Finally, these methods have only been tested in medical imaging applications (Barbano et al., 2023; Chung & Ye, 2024). On the other hand, in the small dataset case, various methods (Moon et al., 2022; Zhang et al., 2024) have been devised to fine-tune a diffusion model on an OOD dataset, but these methods still require several hundred images and have not used the fine-tuned network to solve inverse problems.

Patch-based diffusion models have shown success both for image generation (Wang et al., 2023; Ding et al., 2023) and for inverse problem solving (Hu et al., 2024). In particular, the method of Hu et al. (2024) involves training networks that take in only patches of images at training and reconstruction time, learning priors of the entire images from only priors of patches. In cases of limited data, Hu et al. (2024) shows that patch-based diffusion models outperform whole image models for solving certain inverse problems. These works motivate our key insight that patch-based diffusion priors potentially obtain stronger generalizability than whole-image diffusion priors for both the single measurement setting and the small dataset setting due to a severe lack of data. Inspired by this, we propose to utilize patch-based diffusion models to tackle the challenges arising from mismatched distributions and lack of data in a unified way. We first develop a method to take a network trained on patches of a mismatched distribution and adjust it on the fly in the single measurement setting. We also show how in the small dataset setting, fine-tuning a patch-based network results in a much better prior than fine-tuning a whole-image network, leading to higher quality reconstructed images.

In summary, our contributions are as follows:

- We integrate the patch-based diffusion model framework with the deep image prior (DIP) framework to correct for mismatched distributions in the single measurement setting. Experimentally, we find this approach beats using whole-image models in terms of quantitative metrics and visual image quality in image reconstruction tasks, as well as achieving competitive results with methods using in-distribution diffusion models.
- We show that in the small dataset setting, fine-tuning patch-based diffusion models is much more robust than whole-image models and very little data is required to obtain a reasonable prior for solving inverse problems.
- We demonstrate experimentally that when fine-tuning on very small datasets, whole image diffusion models are prone to overfitting and memorization, which severely degrades reconstructed images, while patch-based models are much less sensitive to this problem.

2 BACKGROUND AND RELATED WORK

Diffusion models and inverse problems. In a general framework, diffusion models involve the forward stochastic differential equation (SDE)

$$d\mathbf{x}_t = -\frac{\beta(t)}{2} \mathbf{x}_t dt + \sqrt{\beta(t)} d\mathbf{w}_t, \quad (1)$$

where $t \in [0, T]$, $\mathbf{x}(t) \in \mathbb{R}^d$, and $\beta(t)$ is the noise variance schedule of the random process $d\mathbf{w}(t)$. This process adds noise to a clean image and ends with an image indistinguishable from Gaussian noise (Song et al., 2021). Thus, the distribution of $\mathbf{x}(0)$ is the data distribution and the distribution

of $\mathbf{x}(T)$ is (approximately) a standard Gaussian. Then the reverse SDE has the form (Anderson, 1982):

$$d\mathbf{x}_t = \left(-\frac{\beta(t)}{2}\mathbf{x}_t - \beta(t)\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) \right) dt + \sqrt{\beta(t)} d\bar{\mathbf{w}}_t. \quad (2)$$

Score-based diffusion models involve training a neural network to learn the score function $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)$, from which one can start with noise and run the reverse SDE to obtain samples from the learned data distribution.

When solving inverse problems, it is necessary to instead sample from $p(\mathbf{x}_t|\mathbf{y})$, so the reverse SDE becomes

$$d\mathbf{x}_t = \left(-\frac{\beta(t)}{2}\mathbf{x}_t - \beta(t)\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{y}) \right) dt + \sqrt{\beta(t)} d\bar{\mathbf{w}}_t. \quad (3)$$

Unfortunately, the term $\log p_t(\mathbf{x}_t|\mathbf{y})$ seems difficult to compute from the unconditional score $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)$ alone. Liu et al. (2023), Chung et al. (2023b), and Ozdenizci & Legenstein (2023) among others proposed directly learning this conditional score $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{y})$ instead. However, this process requires paired data (\mathbf{x}, \mathbf{y}) between the image domain and measurement domain for training, instead of just clean image data. Furthermore, the learned conditional score function is suitable only for the particular inverse problem for which it was trained, limiting its flexibility.

For greater generalizability, it is desirable to apply the unconditional score $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)$ to be able to solve a wide variety of inverse problems. Thus, many works have been proposed to approximate the conditional score in terms of the unconditional one (Wang et al., 2022; Chung et al., 2023a; 2024; Kawar et al., 2022). Notably, Peng et al. (2024) unified various diffusion inverse solvers (DIS) into two categories: the first consists of direct approximations to $p_t(\mathbf{y}|\mathbf{x}_t)$, and the second consists of first approximating $\mathbb{E}[\mathbf{x}_0|\mathbf{x}_t, \mathbf{y}]$ (typically through an optimization problem balancing the prior and measurement) and then applying Tweedie’s formula (Efron, 2011) to obtain

$$\nabla \log p_t(\mathbf{x}_t|\mathbf{y}) = \frac{\mathbb{E}[\mathbf{x}_0|\mathbf{x}_t, \mathbf{y}] - \mathbf{x}_t}{\sigma_t^2}, \quad (4)$$

where σ_t is the noise level of \mathbf{x}_t . All of these methods require a large quantity of clean training data that should come from the distribution $p(\mathbf{x})$ whose score is to be learned, which may not be available in practice.

Methods without clean training data. When no in-distribution data is available, one approach is to use traditional methods that do not require any training data, such as total variation (TV) (Li et al., 2019) or wavelet transform (Daubechies, 1992) regularizers that encourage image sparsity. More recently, plug and play (PnP) methods have risen in popularity (Sun et al., 2021; Sreehari et al., 2016; Hong et al., 2020; 2024b); these methods use a denoiser to solve general inverse problems. Although these methods often use a trained denoiser, Ryu et al. (2019) found that using an off-the-shelf denoiser such as block matching 3D (Dabov et al., 2006) can yield competitive results. Nevertheless, with the rise of deep learning in image processing applications, methods that harness the power of these tools may be desirable.

The deep image prior (DIP) is an extensively studied self-supervised method that is popular when no training data is available and reconstruction from a single measurement \mathbf{y} is desired. The method consists of training a network f_θ using the loss function

$$L(\theta) = \|\mathbf{y} - \mathcal{A}(f_\theta(\mathbf{z}))\|_2^2, \quad \mathbf{z} \sim \mathcal{N}(0, \mathbf{I}), \quad (5)$$

so that $f_\theta(\mathbf{z})$ produces the reconstruction. Although the neural network acts as an implicit regularizer whose output tends to lie in the manifold of clean images, DIP is prone to overfitting (Ulyanov et al., 2020). Various methods have been proposed involving early stopping, regularization, and network initialization (Liu et al., 2018; Jo et al., 2021; Barbano et al., 2022). Nevertheless, the method is very sensitive to parameter selection and implementation and can take a long time to train (Jo et al., 2021).

Most DIS methods learn a prior from a large collection of clean in-distribution training images, but recently Barbano et al. (2023) and Chung & Ye (2024) proposed self-supervised diffusion model methods that are based off the DIP framework. These methods involve alternating between the usual reverse diffusion update step to gradually denoise the image and a network refining step in which the score network parameters are updated via the loss function

$$L(\theta) = \|\mathbf{y} - \mathcal{A}(\text{CG}(\hat{\mathbf{x}}_{0|t}(\mathbf{x}_t; \theta)))\|_2^2 \quad (6)$$

where conjugate gradient (CG) descent is used to enforce data fidelity. This CG step consists of solving an optimization of the form

$$\arg \min_{\mathbf{x}} \frac{\gamma}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{x})\|_2^2 + \frac{1}{2} \|\mathbf{x} - \hat{\mathbf{x}}_{0|t}\|_2^2, \quad (7)$$

where γ is a tradeoff parameter controlling the strength of the prior versus the measurement. Crucially, these methods introduce an additional LoRA module (Hu et al., 2021) to the network and the original network parameters are frozen when backpropagating the loss, which helps to avoid overfitting the whole-image model. Nevertheless, many technical tricks are required (Chung & Ye, 2024) involving noisy initializations and early stopping to obtain good results and avoid artifacts. Our patch-based model avoids this overfitting issue.

Diffusion model fine-tuning. In the small dataset setting, various fine-tuning methods exist to shift the underlying prior learned by a score network away from a mismatched distribution and toward a target distribution. Given a pretrained diffusion network on a mismatched distribution, Moon et al. (2022), Zhang et al. (2024), and Zhu et al. (2024) among others have studied ways to fine-tune the network to the desired dataset. These methods generally involve freezing certain layers of the original network, appending extra modules that contain relatively few weights, or modifying the loss function to capture details that differ greatly between distributions. However, these methods usually still require thousands of images from the desired distribution and focus on image generation. When solving inverse problems, the reconstructed image should be consistent with the measurement \mathbf{y} , reducing the number of degrees of the freedom for the image compared to generation, so with proper fine-tuning the data requirement should be lower.

3 METHODS

3.1 PATCH-BASED PRIOR

We adapt the patch-based diffusion model framework of Hu et al. (2024). We first zero pad the image by an amount P on each side and model the resulting image \mathbf{x} . When choosing the i th patch offset tuple $(o_1, o_2) \in \{0, \dots, P-1\}^2$ in Figure 1, \mathbf{x} is partitioned into many square patches and one bordering region consisting of all zeros. Since $k = N/P$ patches are needed in one direction to perfectly cover the image, a model for the data distribution takes the form

$$p(\mathbf{x}) = \prod_{i=1}^{P^2} p_{i,B}(\mathbf{x}_{i,B}) \prod_{r=1}^{(k+1)^2} p_{i,r}(\mathbf{x}_{i,r}) / Z, \quad (8)$$

where $\mathbf{x}_{i,B}$ represents the bordering region of \mathbf{x} that depends on the specific value of i , $p_{i,B}$ is the probability distribution of that region, $\mathbf{x}_{i,r}$ is the r th $P \times P$ patch when using the partitioning scheme corresponding to the value of i , $p_{i,r}$ is the probability distribution of that region, and Z is a normalizing factor. This model allows for a variety of possible tilings of the image, eliminating boundary artifacts that would occur if only one tiling was used.

For training, we use a neural network $D_{\theta}(\mathbf{x}, \sigma_t)$ that accepts a noisy image \mathbf{x} and the noise level σ_t . For each patch, we define the x positional array as the 2D array consisting of the x positions of each pixel of the image scaled between -1 and 1, and the y positional array is similarly defined for the y positions. To allow the network to learn different patch distributions at different locations in the image, we extract the corresponding patches of these positional arrays and concatenate them along the channel dimension of the noisy image patch and treat the entire array as the network input.

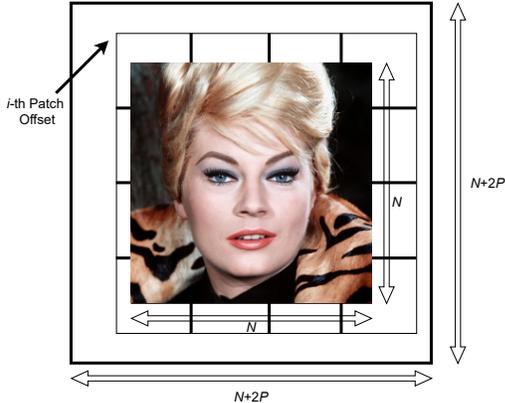


Figure 1: Schematic for zero padding and partitioning image into patches. Each index i represents one of P^2 possible ways to choose a patch offset tuple.

Since we are using a patch-based prior, we perform denoising score matching on patches of an image instead of the whole image. Hence, the training loss is given by

$$\arg \min_{\theta} \mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \sigma_t^2 \mathbf{I})} \|D_{\theta}(\mathbf{x} + \epsilon, \sigma_t) - \mathbf{x}\|_2^2, \quad (9)$$

where $\mathbf{x} \sim p(\mathbf{x})$ represents a patch drawn from a sample of the training dataset, σ_t is a predetermined noise schedule, and \mathcal{U} denotes the uniform distribution.

3.2 SINGLE MEASUREMENT SETTING

Consider the first case where only the measurement \mathbf{y} is given, and no in-sample training data is available. For each specific measurement \mathbf{y} , the DIP framework optimizes the network parameters θ via the self-supervised loss (5) from the predicted reconstructed image. Diffusion models provide a prediction of the reconstructed image at each timestep: namely, the expectation of the clean image $\mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t]$ is approximated by the denoiser $D_{\theta}(\mathbf{x}_t)$ via Tweedie’s formula. Then the expectation conditioned on the measurement $\mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t, \mathbf{y}]$ can be obtained through one of many methods of enforcing the data fidelity constraint.

We begin with the unconditional expectation by leveraging the patch-based prior. Following (8), we apply Tweedie’s formula to express the denoiser of \mathbf{x} in terms of solely the denoisers of the patches of \mathbf{x} . Because the outermost product is computationally very expensive, in practice we approximate $D_{\theta}(\mathbf{x})$ using only a single randomly selected value of i for each denoiser evaluation:

$$D_{\theta}(\mathbf{x}) \approx D_{i, B}(\mathbf{x}_{i, B}) + \sum_{r=1}^{(k+1)^2} D_{i, r}(\mathbf{x}_{i, r}). \quad (10)$$

By definition, $D_{i, B}(\mathbf{x}_{i, B}) = 0$ and we compute each $D_{i, r}(\mathbf{x}_{i, r})$ with the network. Note that (10) provides an *unconditional* estimate of the clean image; to obtain a conditional estimate $D_{\theta}(\mathbf{x}_t | \mathbf{y})$ of the clean image, we run C iterations of the conjugate gradient descent algorithm for minimizing $\|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2$, initialized with the unconditional estimate (Chung et al., 2024). The image that is being reconstructed might not come from the distribution of the training images. Hence, the estimate $D_{\theta}(\mathbf{x}_t | \mathbf{y})$ may be far from the true denoised image. Thus, we use \mathbf{y} to update the parameters of the network in a way such that $D_{\theta}(\mathbf{x}_t | \mathbf{y})$ becomes more consistent with the measurement:

$$\theta \leftarrow \arg \min_{\theta} \|\mathbf{y} - \mathbf{A} D_{\theta}(\mathbf{x}_t | \mathbf{y})\|_2^2. \quad (11)$$

Previously, additional LoRA parameters (Hu et al., 2021) have been used as an injection to the network to leave the original parameters unchanged during this process (Barbano et al., 2023; Chung & Ye, 2024). However, the effect of using different ranks for LoRA versus other methods of network fine-tuning on DIS has not been studied extensively, so we opt to update all the weights of the network in this step. Appendix A.3 shows results from using the LoRA module.

Crucially, iterative usage of CG for computing the conditional denoiser allows for simple and efficient backpropagation through this loss function, a task that would be much more computationally challenging if another DIS such as Chung et al. (2023a) or Wang et al. (2022) were used. Furthermore, because the number of diffusion steps is large and the change in \mathbf{x}_t is small between consecutive timesteps, we apply this network refining step only for certain iterations of the diffusion process, reducing the computational burden.

Algorithm 1 Single Measurement Inverse Solver

Require: $\sigma_1 < \sigma_2 < \dots < \sigma_T$, $\epsilon > 0$, P, C, \mathbf{y}, K
Initialize $\mathbf{x} \sim \mathcal{N}(0, \sigma_T^2 \mathbf{I})$
for $t = T : 1$ **do**
 if $t \bmod K = 0$ **then**
 Compute $D_{\theta}(\mathbf{x}_t)$ using (10) with a random index i
 Run C iterations of CG initialized with $D_{\theta}(\mathbf{x}_t)$ to obtain $D_{\theta}(\mathbf{x}_t | \mathbf{y})$
 Define $L(\theta) = \|\mathbf{y} - \mathbf{A} D_{\theta}(\mathbf{x}_t | \mathbf{y})\|_2^2$
 Update θ by backpropagating $L(\theta)$
 end if
 Sample $\mathbf{z} \sim \mathcal{N}(0, \sigma_t^2 \mathbf{I})$
 Set $\alpha_t = \epsilon \cdot \sigma_t^2$
 Compute $D(\mathbf{x}_t)$ using (10) with a random index i
 Run C iterations of CG for (7) initialized with $D(\mathbf{x}_t)$
 Set $\mathbf{s}_t = (D - \mathbf{x}_t) / \sigma_t^2$
 Set \mathbf{x}_{t-1} to $\mathbf{x}_t + \frac{\alpha_t}{2} \mathbf{s}_t + \sqrt{\alpha_t} \mathbf{z}$
end for

270 After this step, we apply the refined network to
 271 compute a new estimate of the score of \mathbf{x}_t and
 272 then use it to update \mathbf{x}_t . Similar to the network
 273 refining step, we use the stochastic version of the denoiser given by (10) rather than the full version.
 274 Hu et al. (2024) showed that for patch-based priors, Langevin dynamics Song & Ermon (2019)
 275 works particularly well as a sampling algorithm, so we use it here in conjunction with CG steps
 276 to enforce data fidelity. Algorithm 1 summarizes the entire method for cases where only a single
 277 measurement \mathbf{y} is available.

278 3.3 SMALL DATASET SETTING

279 Now turn to the case where we have trained a diffusion model on OOD data, but we also have a very
 280 small dataset of in-distribution test data that we can use to fine-tune the model. When fine-tuning,
 281 we initialize the network with the checkpoint trained on OOD data and then use the denoising score
 282 matching loss function to fine-tune the network on in-distribution data. Wang et al. (2023) found
 283 that improved image generation performance can be obtained by training with varying patch sizes,
 284 as opposed to fixing the patch size to the one used during inference. Here, we apply a varying patch
 285 size scheme during fine-tuning also as a method of data augmentation. We use the UNet architecture
 286 in Ho et al. (2020) that can accept images of different sizes. Hence, the loss becomes
 287

$$288 \arg \min_{\theta} \mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x} \sim p_d(\mathbf{x})} \mathbb{E}_{\varepsilon \sim \mathcal{N}(0, \sigma_t^2 I)} \|D_{\theta}(\mathbf{x} + \varepsilon, \sigma_t) - \mathbf{x}\|_2^2, \quad (12)$$

289 where $\mathbf{x} \sim p_d(\mathbf{x})$ represents the drawing a randomly sized patch from an image belonging to the
 290 fine-tuning dataset. Appendix A.5 provides full details of the training process.

291 At reconstruction time, we assume that our network has been fine-tuned reasonably to our dataset.
 292 Thus, we remove the network refining step in Algorithm 1 and keep the weights fixed throughout
 293 the entire process. We still use the same CG descent algorithm to enforce data fidelity with the
 294 measurement.

295 4 EXPERIMENTS

296 **Experimental setup.** For the CT experiments, we used the AAPM 2016 CT challenge data from
 297 McCollough et al. (2017). We applied the same data processing methods as in Hu et al. (2024) with
 298 the exception that we used all the XY slices from the 9 training volumes to train the in distribution
 299 networks, yielding a total of approximately 5000 slices. For the deblurring and superresolution
 300 experiments, we used the CelebA-HQ dataset (Liu et al., 2015) with each image having size $256 \times$
 301 256 . The test data was a randomly selected subset of 10 of the images not used for training. In all
 302 cases, we report the average metrics across the test images: peak SNR (PSNR) in dB, and structural
 303 similarity metric (SSIM) (Wang et al., 2004). For the training data, we trained networks on generated
 304 phantom images consisting of randomly placed ellipses of different shapes and sizes. See Fig. 20
 305 for examples. These phantoms can be generated on the fly in large quantities. We used networks
 306 trained on grayscale phantoms for the CT experiments and networks trained on RGB phantoms for
 307 the deblurring and superresolution experiments. Appendix A.4 contains precise specifications of the
 308 phantoms.

309 We trained the patch-based networks with 64×64 patches and used a zero padding value of 64,
 310 so that 5 patches in both directions were used to cover the target image. We used the network
 311 architecture in Karras et al. (2022) for both the patch-based networks and whole-image networks.
 312 All networks were trained on PyTorch using the Adam optimizer with 2 A40 GPUs.

313 **Single measurement setting.** In cases where no training data is available and we only have the
 314 measurement \mathbf{y} , we applied Algorithm 1 to solve a variety of inverse problems: CT reconstruction,
 315 deblurring, and superresolution. For the forward and backward projectors in CT reconstruction, we
 316 used the implementation provided by the ODL Team (2022). We performed two sparse-view CT
 317 (SVCT) experiments: one using 20 projection views, and one using 60 projection views. Both of
 318 these were done using a parallel beam forward projector where the detector size was 512 pixels. For
 319 the deblurring experiments, we used a uniform blur kernel of size 9×9 and added white Gaussian
 320 noise with $\sigma = 0.01$ where the clean image was scaled between 0 and 1. For the superresolution ex-

periments, we used a scaling factor of 4 with downsampling by averaging and added white Gaussian noise with $\sigma = 0.01$.

For the comparison methods, we ran experiments that naively used the OOD diffusion model without the self-supervised network refining process. For reference, we also ran experiments using a diffusion model trained on the entire in-distribution training set (the “correct” model). In practice, it would not be possible to obtain such a large training dataset of in-distribution images. Additionally, for these diffusion model methods, we implemented both the patch-based version as well as the whole-image version. The whole-image networks were trained with the loss function in (9) and used the same network architecture as the patch-based models, but the input of the network was the entire image and did not contain positional encoding information.

We also compared with more traditional methods: applying a simple baseline, reconstructing via the total variation regularizer (ADMM-TV), and two plug and play (PnP) methods: PnP-ADMM (Xu et al., 2020) and PnP-RED (Hu et al., 2022). For CT, the baseline was obtained by applying the filtered back-projection method to the measurement \mathbf{y} . For deblurring, the baseline was simply equal to the blurred image. For superresolution, the baseline was obtained by upsampling the low resolution image and using nearest neighbor interpolation. The implementation of ADMM-TV can be found in Hong et al. (2024a). Finally, since we assume we do not have access to any clean training data, we used the off the shelf denoiser BM3D (Dabov et al., 2006). Appendix A.5 contains the values of all the parameters of the algorithms.

Table 1 shows the main results for single-measurement inverse problem solving. The bottom two rows show the hypothetical performance if it were possible to train a diffusion model on a large dataset of in distribution images, which is not available in practice. Our self-supervised patch-based diffusion approach achieved significantly higher quantitative results when averaged across the test dataset than the self-supervised whole-image approach in all the inverse problems. Furthermore, although the diffusion model that was initially used in this algorithm was trained on completely different images, by applying the self-supervised loss, the patch-based approach is able to achieve results that are close to (and for the deblurring case, even surpassing) those using the in-distribution networks. The table also shows that by including the self-supervised step, a dramatic improvement over naively using the OOD model is achieved. Lastly, Fig. 2 shows that some artifacts appear in the whole-image SS method that are not present in our patch SS method.

Table 1: Comparison of quantitative results on three different inverse problems in the single measurement setting. Results are averages across all images in the test dataset. Best results for practical use are in bold.

Method	CT, 20 Views		CT, 60 Views		Deblurring		Superresolution	
	PSNR \uparrow	SSIM \uparrow						
Baseline	24.93	0.613	30.15	0.784	23.93	0.666	25.42	0.724
ADMM-TV	26.81	0.750	31.14	0.862	27.58	0.773	25.22	0.729
PnP-ADMM (Xu et al., 2020)	30.20	0.838	36.75	0.932	28.98	0.815	27.29	0.796
PnP-RED (Hu et al., 2022)	27.12	0.682	32.68	0.876	28.37	0.793	27.73	0.809
Whole image, naive	28.11	0.800	33.10	0.911	25.85	0.742	25.65	0.742
Patches, naive (Hu et al., 2024)	27.44	0.719	33.97	0.934	26.77	0.782	26.12	0.759
Self-supervised, whole (Barbano et al., 2023)	33.19	0.861	40.47	0.957	29.50	0.831	27.07	0.701
Self-supervised, patch (Ours)	33.77	0.874	41.45	0.969	30.34	0.860	28.10	0.827
Whole image, correct*	33.99	0.886	41.67	0.969	29.87	0.851	28.33	0.801
Patches, correct*	34.02	0.889	41.70	0.967	30.12	0.865	28.49	0.835

*not available in practice for mismatched distribution inverse problems

To demonstrate that our method also works well even when the mismatched distribution is closer to the true distribution, we also ran an experiment where the networks were initially trained on the LIDC-IDRI dataset (Armato et al., 2011). We extracted 10000 2D slices from the 3D volumes and rescaled all the images so that the pixel values were between 0 and 1. We then ran Algorithm 1 to perform CT reconstruction where the test dataset was the same as the one used in Table 1. Table 4 shows the results of this experiment. Our method achieved better quantitative results than the whole image method and even outperformed the reconstructions using the in distribution network but without any self-supervision. Appendix A.1 shows the visual results of these experiments. Appendix A.2 further discusses using self-supervision in cases where the initial network was trained on in-distribution data and shows improved image quality.

We also ran ablation studies to examine the effect of various parameters on the proposed method. Barbano et al. (2023) and Chung & Ye (2024) used the LoRA module for solving single-measurement inverse problems with diffusion models. We tested this method for CT reconstruction and deblurring with different rank adjustments and found this method to be inferior to modifying the weights of the entire network. We also ran experiments using networks with different numbers of weights. Appendix A.3 shows the results of these experiments.

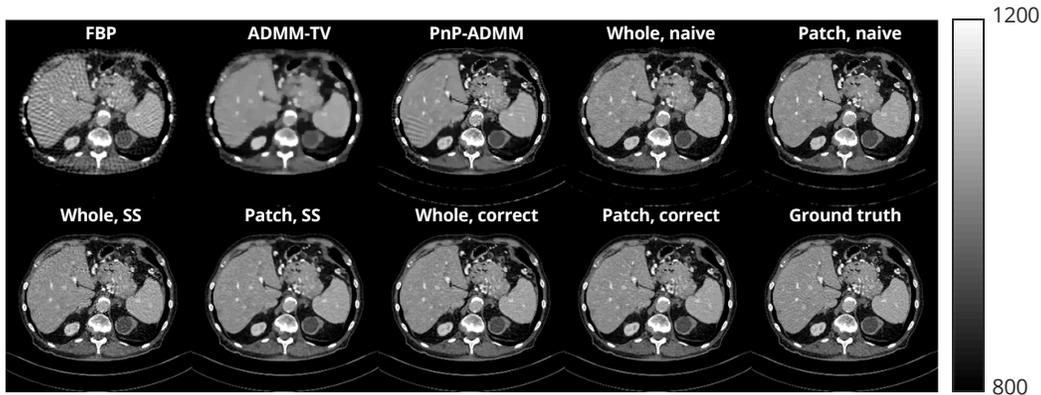


Figure 2: **Single measurement setting:** Results of 60 view CT reconstruction using self supervised (SS) approach. The display uses modified HU units to show more contrast between organs.



Figure 3: **Single measurement setting:** Results of deblurring using self supervised (SS) approach and comparison methods.

Small dataset setting. We ran experiments on the same inverse problems as the single measurement case. The OOD networks were fine-tuned with 10 images randomly selected from the in-distribution training set; we also ran ablation studies using different quantities of in-distribution data in Appendix A.3. Figures 4 and 5 show that the patch-based model is much less prone to overfitting than the whole-image model. Hence, to evaluate the best possible performance of the whole-image model compared to the patch-based model, for both models we chose the checkpoint yielding the best results for solving inverse problems.

Table 2 shows the main results for solving inverse problems using the fine-tuned diffusion model. We compared the results of fine-tuning the whole-image model with fine-tuning the patch-based model as well as the best baseline out of the four baselines shown in Table 1, namely "baseline", ADMM-TV, PnP-ADMM, and PnP-RED. The results show that the proposed patch-based method achieved the best performance in terms of quantitative metrics for all of the inverse problems. Figure 6 shows the visual results of this experiment. The patch-based model is able to learn an acceptable prior using the very small in-distribution dataset and the reconstructed images contain fewer artifacts than the comparison methods.

Table 2: Comparison of results for using diffusion models fine-tuned on 10 in-distribution images to solve inverse problems in small dataset setting. Best results are in bold.

Method	CT, 20 Views		CT, 60 Views		Deblurring		Superresolution	
	PSNR \uparrow	SSIM \uparrow						
Best baseline	30.20	0.838	36.75	0.932	28.98	0.815	27.73	0.809
Whole image	33.09	0.875	40.54	0.964	28.41	0.812	27.29	0.775
Patches (Ours)	33.44	0.875	41.21	0.965	29.25	0.840	28.10	0.827
Patches, correct*	34.02	0.889	41.70	0.967	30.12	0.865	28.49	0.835

*not available in practice for mismatched distribution inverse problems

Figures 4 and 5 further investigate the effect of overfitting. For different amounts of training time using the small in-distribution dataset, we ran the reconstruction algorithm for 60-view CT. While the whole-image model exhibited substantial image degradation when the network was fine-tuned for too long, the patch-based model retained relatively stable performance throughout the entire training process. This illustrates that whole-image diffusion models exhibit severe overfitting problems when only a small amount of training data is available. Furthermore, patch-based diffusion models assist greatly with this problem and the results are evident for solving inverse problems. Appendix A.1 shows the visual results of these experiments.

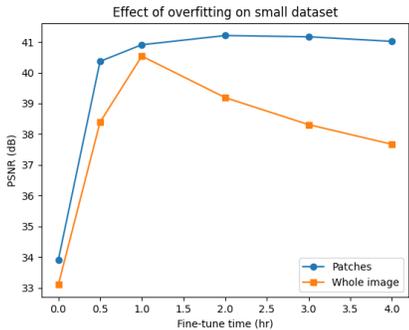


Figure 4: Comparison of PSNR between patch-based model and whole-image model for overfitting in small dataset setting.

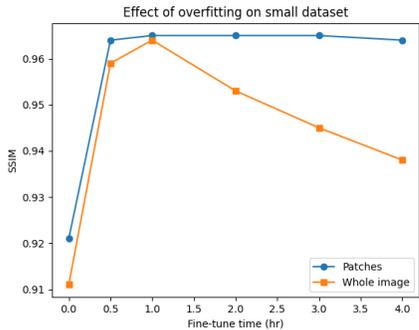
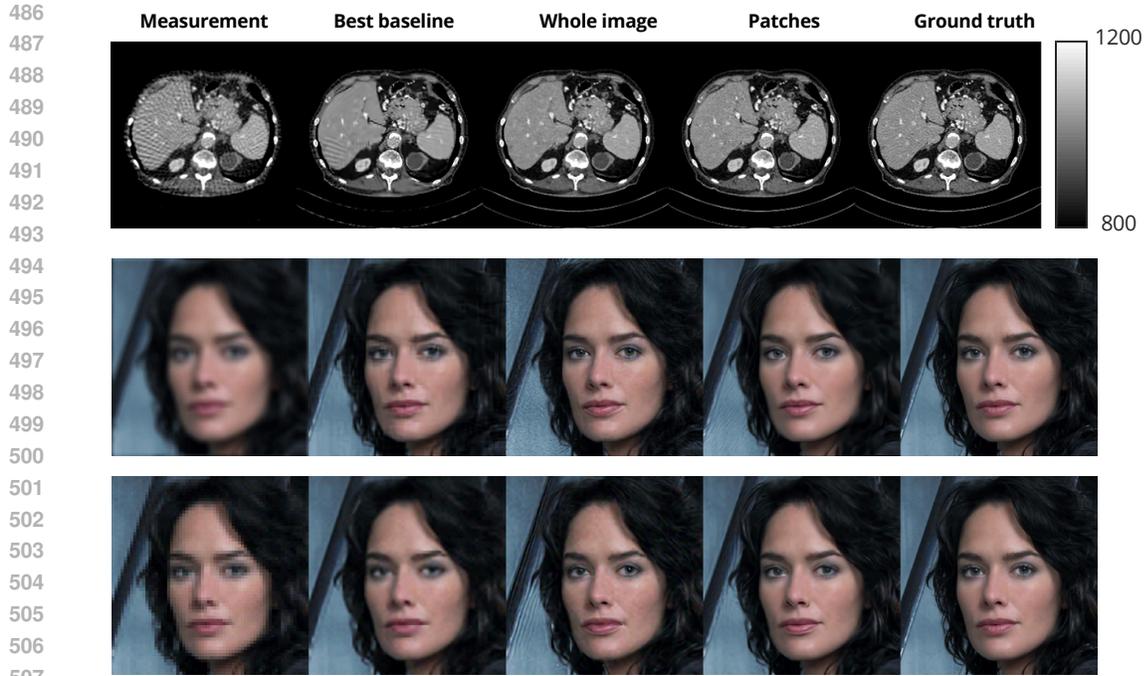


Figure 5: Comparison of SSIM between patch-based model and whole-image model for overfitting in small dataset setting.

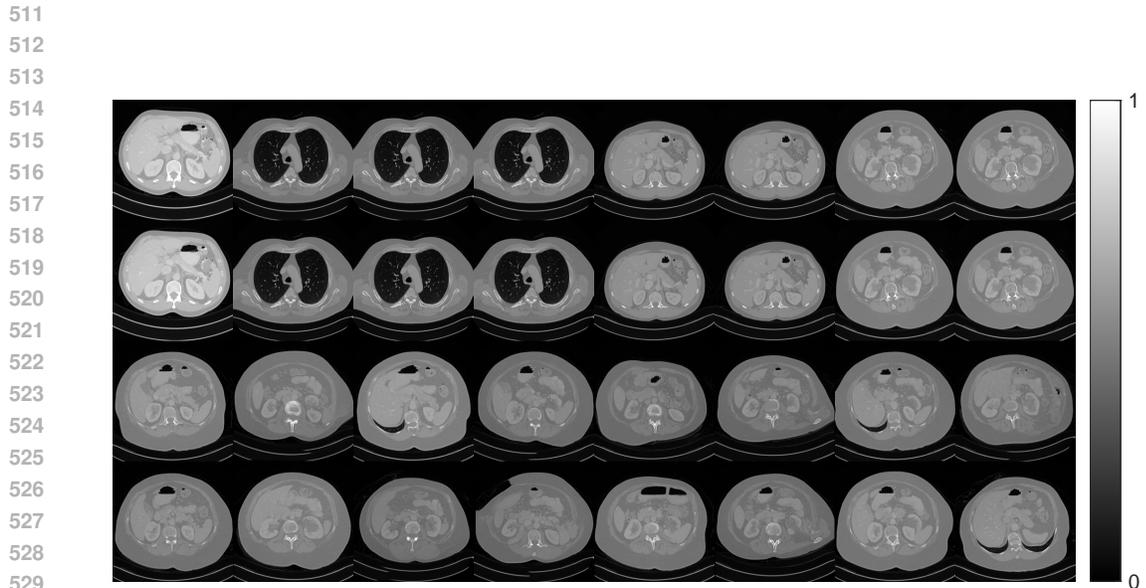
To look at the priors learned by the different models from fine-tuning, we unconditionally generated images from the checkpoints obtained by fine-tuning on the 10 image CT dataset. Figure 7 shows a subset of the generated images where we used the checkpoints obtained after 4 hours of training. The top two rows consist of images generated by the whole-image model and the bottom two rows consist of images generated by the patch diffusion model. To emphasize the memorization point, we grouped together similar looking images in the top two rows: it can be seen that the images in each group look virtually identical, despite the fact that the pure white noise initializations for each sample was different. On the other hand, while the samples generated by the patch diffusion model also show some unrealistic features, they all show some distinct features, which implies that this model has much better generalization ability.

5 CONCLUSION

This paper presented a method of using patch-based diffusion models to solve inverse problems when the data distribution is mismatched from the trained network. In particular, we conducted experiments in the setting when only a single measurement is available as well as the setting when a very small subset of in-distribution data is available. In both settings, the proposed patch-based method outperformed whole-image methods in a variety of inverse problems. In the future, more work could be done on using acceleration methods for faster reconstruction, exploring other less computationally expensive methods of fine-tuning the network geared toward inverse problem solving, and methods of refining the prior when a set of measurements are available (Yaman et al., 2020).



508 **Figure 6: Small dataset setting:** Results of inverse problem solving. Top row is 60 view CT recon,
509 middle row is deblurring, and bottom row is superresolution. For CT, measurement refers to FBP.
510



530 **Figure 7: Unconditional generation of CT images** from networks fine-tuned in the small dataset
531 setting. Top two rows were generated with the whole image model; bottom two rows were generated
532 with the patch-based model.
533

534
535
536
537
538 Limitations of the work include a slow runtime for the self-supervised algorithm and a lack of theo-
539 retical guarantees for dataset size requirements. [Providing uncertainty quantification is also an open problem for such self-supervised methods.](#)

ACKNOWLEDGMENTS

REFERENCES

- Brian D.O. Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326, 1982. ISSN 0304-4149. doi: [https://doi.org/10.1016/0304-4149\(82\)90051-5](https://doi.org/10.1016/0304-4149(82)90051-5).
- Samuel G Armato, Geoffrey McLennan, Luc Bidaut, Michael F McNitt-Gray, and Charles R Meyer. The lung image database consortium (lidc) and image database resource initiative (idri): A completed reference database of lung nodules on ct scans. *Medical Physics*, 38:915–931, 2011.
- Daniel Otero Bagger, Johannes Leuschner, and Maximilian Schmidt. Computed tomography reconstruction using deep image prior and learned reconstruction methods. *Inverse Problems*, 36(9):094004, September 2020. ISSN 1361-6420. doi: 10.1088/1361-6420/aba415.
- Riccardo Barbano, Johannes Leuschner, Maximilian Schmidt, Alexander Denker, Andreas Hauptmann, Peter Maass, and Bangti Jin. An educated warm start for deep image prior-based micro ct reconstruction. *IEEE Transactions on Computational Imaging*, 8:1210–1222, 2022. ISSN 2573-0436. doi: 10.1109/tci.2022.3233188.
- Riccardo Barbano, Alexander Denker, Hyungjin Chung, Tae Hoon Roh, Simon Arridge, Peter Maass, Bangti Jin, and Jong Chul Ye. Steerable conditional diffusion for out-of-distribution adaptation in imaging inverse problems, 2023. URL <https://arxiv.org/abs/2308.14409>.
- Hyungjin Chung and Jong Chul Ye. Deep diffusion image prior for efficient ood adaptation in 3d inverse problems, 2024. URL <https://arxiv.org/abs/2407.10641>.
- Hyungjin Chung, Byeongsu Sim, Dohoon Ryu, and Jong Chul Ye. Improving diffusion models for inverse problems using manifold constraints, 2022a. URL <https://arxiv.org/abs/2206.00941>.
- Hyungjin Chung, Byeongsu Sim, Dohoon Ryu, and Jong Chul Ye. Improving diffusion models for inverse problems using manifold constraints. In *Advances in Neural Information Processing Systems*, volume 35, pp. 25683–25696, 2022b. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/a48e5877c7bf86a513950ab23b360498-Paper-Conference.pdf.
- Hyungjin Chung, Jeongsol Kim, Michael Thompson Mccann, Marc Louis Klasky, and Jong Chul Ye. Diffusion posterior sampling for general noisy inverse problems. In *The Eleventh International Conference on Learning Representations*, 2023a. URL <https://openreview.net/forum?id=OnD9zGAGT0k>.
- Hyungjin Chung, Jeongsol Kim, and Jong Chul Ye. Direct diffusion bridge using data consistency for inverse problems, 2023b. URL <https://arxiv.org/abs/2305.19809>.
- Hyungjin Chung, Suhyeon Lee, and Jong Chul Ye. Decomposed diffusion sampler for accelerating large-scale inverse problems, 2024. URL <https://arxiv.org/abs/2303.05754>.
- Kostadin Dabov, Alessandro Foi, Vladimir Katkovnik, and Karen Egiazarian. Image denoising with block-matching and 3d filtering. In *Proc. SPIE 6064, Image Processing: Algorithms and Systems, Neural Networks, and Machine Learning*, volume 6064, pp. 606414, February 2006. doi: 10.1117/12.643267.
- Ingrid Daubechies. *Ten Lectures on Wavelets*. Society for Industrial and Applied Mathematics, 1992. doi: 10.1137/1.9781611970104.
- Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. In *Advances in Neural Information Processing Systems*, volume 34, pp. 8780–8794, 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/file/49ad23d1ec9fa4bd8d77d02681df5cfa-Paper.pdf.
- Zheng Ding, Mengqi Zhang, Jiajun Wu, and Zhuowen Tu. Patched denoising diffusion models for high-resolution image synthesis, 2023. URL <https://arxiv.org/abs/2308.01316>.

- 594 Bradley Efron. Tweedie’s formula and selection bias. *Journal of the American Statistical Association*, 106(496):1602–1614, 2011.
- 595
- 596
- 597 Berthy T. Feng, Jamie Smith, Michael Rubinstein, Huiwen Chang, Katherine L. Bouman, and
598 William T. Freeman. Score-based diffusion models as principled priors for inverse imaging, 2023.
599 URL <https://arxiv.org/abs/2304.11751>.
- 600
- 601 Berthy T. Feng, Ricardo Baptista, and Katherine L. Bouman. Neural approximate mirror maps for
602 constrained diffusion models, 2024. URL <https://arxiv.org/abs/2406.12816>.
- 603 T.E Gureyev, A Pogany, D.M Paganin, and S.W Wilkins. Linear algorithms for phase retrieval
604 in the fresnel region. *Optics Communications*, 231(1):53–70, 2004. ISSN 0030-4018. doi:
605 <https://doi.org/10.1016/j.optcom.2003.12.020>.
- 606
- 607 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In
608 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neu-*
609 *ral Information Processing Systems*, volume 33, pp. 6840–6851. Curran Associates, Inc.,
610 2020. URL [https://proceedings.neurips.cc/paper_files/paper/2020/
611 file/4c5bcfec8584af0d967f1ab10179ca4b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/4c5bcfec8584af0d967f1ab10179ca4b-Paper.pdf).
- 612
- 613 Tao Hong, Irad Yavneh, and Michael Zibulevsky. Solving red with weighted proximal methods.
IEEE Signal Processing Letters, 27:501–505, 2020. doi: 10.1109/LSP.2020.2979062.
- 614
- 615 Tao Hong, Luis Hernandez-Garcia, and Jeffrey A. Fessler. A complex quasi-newton proximal
616 method for image reconstruction in compressed sensing mri. *IEEE Transactions on Computa-*
617 *tional Imaging*, 10:372–384, 2024a. ISSN 2573-0436. doi: 10.1109/tci.2024.3369404.
- 618
- 619 Tao Hong, Xiaojian Xu, Jason Hu, and Jeffrey A Fessler. Provable preconditioned plug-and-play
620 approach for compressed sensing mri reconstruction. *arXiv preprint arXiv:2405.03854*, 2024b.
621 URL <https://arxiv.org/abs/2405.03854>.
- 622
- 623 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,
624 and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. URL <https://arxiv.org/abs/2106.09685>.
- 625
- 626 Jason Hu, Bowen Song, Xiaojian Xu, Liyue Shen, and Jeffrey A. Fessler. Learning image priors
627 through patch-based diffusion models for solving inverse problems, 2024. URL <https://arxiv.org/abs/2406.02462>.
- 628
- 629 Yuyang Hu, Jiaming Liu, Xiaojian Xu, and Ulugbek S. Kamilov. Monotonically convergent regu-
630 larization by denoising, 2022. URL <https://arxiv.org/abs/2202.04961>.
- 631
- 632 Ajil Jalal, Marius Arvinte, Giannis Daras, Eric Price, Alexandros G Dimakis, and Jon
633 Tamir. Robust compressed sensing mri with deep generative priors. In *Advances in Neu-*
634 *ral Information Processing Systems*, volume 34, pp. 14938–14954. Curran Associates, Inc.,
635 2021. URL [https://proceedings.neurips.cc/paper_files/paper/2021/
636 file/7d6044e95a16761171b130dcb476a43e-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/7d6044e95a16761171b130dcb476a43e-Paper.pdf).
- 637
- 638 Yeonsik Jo, Se Young Chun, and Jonghyun Choi. Rethinking deep image prior for denoising, 2021.
639 URL <https://arxiv.org/abs/2108.12841>.
- 640
- 641 Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative
642 adversarial networks, 2019. URL <https://arxiv.org/abs/1812.04948>.
- 643
- 644 Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-
645 based generative models, 2022. URL <https://arxiv.org/abs/2206.00364>.
- 646
- 647 Bahjat Kawar, Gregory Vaksman, and Michael Elad. Snips: Solving noisy inverse problems stochas-
648 tically, 2021. URL <https://arxiv.org/abs/2105.14951>.
- 649
- 650 Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising diffusion restoration
651 models, 2022. URL <https://arxiv.org/abs/2201.11793>.

- 648 Sui Li, Dong Zeng, Jiangjun Peng, Zhaoying Bian, Hao Zhang, Qi Xie, Yongbo Wang, Yuting Liao,
649 Shanli Zhang, Jing Huang, Deyu Meng, Zongben Xu, and Jianhua Ma. An efficient iterative
650 cerebral perfusion ct reconstruction via low-rank tensor decomposition with spatial–temporal total
651 variation regularization. *IEEE Transactions on Medical Imaging*, 38(2):360–370, 2019. doi:
652 10.1109/TMI.2018.2865198.
- 653 Zongyu Li, Jason Hu, Xiaojian Xu, Liyue Shen, and Jeffrey A. Fessler. Poisson-gaussian holo-
654 graphic phase retrieval with score-based image prior, 2023a. URL [https://arxiv.org/
655 abs/2305.07712](https://arxiv.org/abs/2305.07712).
- 656 Zongyu Li, Xiaojian Xu, Jason Hu, Jeffrey Fessler, and Yuni Dewaraja. Reducing spect ac-
657 quisition time by predicting missing projections with single-scan self-supervised coordinate-
658 based learning. *Journal of Nuclear Medicine*, 64(supplement 1):P1014–P1014, 2023b. URL
659 https://jnm.snmjournals.org/content/64/supplement_1/P1014.
- 660 Guan-Horng Liu, Arash Vahdat, De-An Huang, Evangelos A. Theodorou, Weili Nie, and Anima
661 Anandkumar. I²sb: Image-to-image schrödinger bridge, 2023. URL [https://arxiv.org/
662 abs/2302.05872](https://arxiv.org/abs/2302.05872).
- 663 Jiaming Liu, Yu Sun, Xiaojian Xu, and Ulugbek S. Kamilov. Image restoration using total variation
664 regularized deep image prior, 2018. URL <https://arxiv.org/abs/1810.12864>.
- 665 Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild.
666 In *Proceedings of International Conference on Computer Vision (ICCV)*, 12 2015.
- 667 Cynthia H. McCollough, Adam C. Bartley, Rickey E. Carter, Baiyu Chen, Tammy A. Drees, Phillip
668 Edwards, David R. Holmes, Alice E. Huang, Farhana Khan, Shuai Leng, Kyle L. McMillan,
669 Gregory J. Michalak, Kristina M. Nunez, Lifeng Yu, and Joel G. Fletcher. Results of the 2016
670 low dose ct grand challenge. *Medical physics*, 44(10):e339–e352, October 2017. ISSN 0094-
671 2405. doi: 10.1002/mp.12345.
- 672 Taehong Moon, Moonseok Choi, Gayoung Lee, Jung-Woo Ha, and Juho Lee. Fine-tuning diffusion
673 models with limited data. In *NeurIPS 2022 Workshop on Score-Based Methods*, 2022. URL
674 <https://openreview.net/forum?id=0J6afk9DqrR>.
- 675 O. Ozdenizci and R. Legenstein. Restoring vision in adverse weather conditions with patch-based
676 denoising diffusion models. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 45
677 (08):10346–10357, 8 2023. ISSN 1939-3539. doi: 10.1109/TPAMI.2023.3238179.
- 678 Xinyu Peng, Ziyang Zheng, Wenrui Dai, Nuoqian Xiao, Chenglin Li, Junni Zou, and Hongkai
679 Xiong. Improving diffusion models for inverse problems using optimal posterior covariance,
680 2024. URL <https://arxiv.org/abs/2402.02149>.
- 681 Albert W. Reed, Hyojin Kim, Rushil Anirudh, K. Aditya Mohan, Kyle Champley, Jingu Kang, and
682 Suren Jayasuriya. Dynamic ct reconstruction from limited views with implicit neural representa-
683 tions and parametric motion fields, 2021. URL <https://arxiv.org/abs/2104.11745>.
- 684 Ernest K. Ryu, Jialin Liu, Sicheng Wang, Xiaohan Chen, Zhangyang Wang, and Wotao Yin. Plug-
685 and-play methods provably converge with properly trained denoisers, 2019. URL [https://
686 arxiv.org/abs/1905.05406](https://arxiv.org/abs/1905.05406).
- 687 Bowen Song, Jason Hu, Zhaoxu Luo, Jeffrey A. Fessler, and Liyue Shen. Diffusionblend: Learning
688 3d image prior through position-aware diffusion score blending for 3d computed tomography
689 reconstruction, 2024. URL <https://arxiv.org/abs/2406.10211>.
- 690 Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the
691 data distribution. In *Advances in Neural Information Processing Systems*, volume 32,
692 2019. URL [https://proceedings.neurips.cc/paper_files/paper/2019/
693 file/3001ef257407d5a371a96dcd947c7d93-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/3001ef257407d5a371a96dcd947c7d93-Paper.pdf).
- 694 Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and
695 Ben Poole. Score-based generative modeling through stochastic differential equations. In *9th
696 International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May
697 3-7, 2021*, 2021. URL <https://openreview.net/forum?id=PXTIG12RRHS>.
- 700
701

- 702 Yang Song, Liyue Shen, Lei Xing, and Stefano Ermon. Solving inverse problems in medical imaging
703 with score-based generative models. In *International Conference on Learning Representations*,
704 2022. URL <https://openreview.net/forum?id=vaRCHVj0uGI>.
705
- 706 S. Sreehari, S. V. Venkatakrisnan, B. Wohlberg, G. T. Buzzard, L. F. Drummy, J. P. Simmons, and
707 C. A. Bouman. Plug-and-play priors for bright field electron tomography and sparse interpolation.
708 *IEEE Transactions on Computational Imaging*, 2(4):408–423, December 2016.
- 709 Yu Sun, Zihui Wu, Xiaojian Xu, Brendt Egon Wohlberg, and Ulugbek Kamilov. Scalable plug-and-
710 play admm with convergence guarantees. *IEEE Transactions on Computational Imaging*, 7, 7
711 2021. ISSN 2573-0436. doi: 10.1109/TCI.2021.3094062.
712
- 713 ODL Development Team. Odl: Operator discretization library. https://odlgroup.github.io/odl/guide/geometry_guide.html, 2022. Accessed: April 2024.
714
- 715 Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Deep image prior. *International Jour-
716 nal of Computer Vision*, 128(7):1867–1888, March 2020. ISSN 1573-1405. doi: 10.1007/
717 s11263-020-01303-4.
- 718 Yinhuai Wang, Jiwen Yu, and Jian Zhang. Zero-shot image restoration using denoising diffusion
719 null-space model, 2022. URL <https://arxiv.org/abs/2212.00490>.
720
- 721 Zhendong Wang, Yifan Jiang, Huangjie Zheng, Peihao Wang, Pengcheng He, Zhangyang Wang,
722 Weizhu Chen, and Mingyuan Zhou. Patch diffusion: Faster and more data-efficient training of
723 diffusion models, 2023. URL <https://arxiv.org/abs/2304.12526>.
724
- 725 Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment:
726 From error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):
727 600–612, 4 2004.
- 728 Zihui Wu, Yu Sun, Jiaming Liu, and Ulugbek Kamilov. Online regularization by denoising with
729 applications to phase retrieval. In *2019 IEEE/CVF International Conference on Computer Vision
730 Workshop (ICCVW)*, pp. 3887–3895, 2019. doi: 10.1109/ICCVW.2019.00482.
- 731 Xiaojian Xu, Jiaming Liu, Yu Sun, Brendt Wohlberg, and Ulugbek S. Kamilov. Boosting the perfor-
732 mance of plug-and-play priors via denoiser scaling. In *54th Asilomar Conf. on Signals, Systems,
733 and Computers*, pp. 1305–1312, 2020. doi: 10.1109/IEEECONF51394.2020.9443410.
734
- 735 B. Yaman, S. A. H. Hosseini, S. Moeller, J. Ellermann, Kamil Ugurbil, and M. Akcakaya. Self-
736 supervised learning of physics-based reconstruction neural networks without fully-sampled refer-
737 ence data. *Mag. Res. Med.*, 84(6):3172–91, December 2020. doi: 10.1002/mrm.28378.
- 738 Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable
739 effectiveness of deep features as a perceptual metric, 2018. URL [https://arxiv.org/abs/
740 1801.03924](https://arxiv.org/abs/1801.03924).
- 741 Xinxi Zhang, Song Wen, Ligong Han, Felix Juefei-Xu, Akash Srivastava, Junzhou Huang, Hao
742 Wang, Molei Tao, and Dimitris N. Metaxas. Spectrum-aware parameter efficient fine-tuning for
743 diffusion models, 2024. URL <https://arxiv.org/abs/2405.21050>.
744
- 745 Jingyuan Zhu, Huimin Ma, Jiansheng Chen, and Jian Yuan. Domainstudio: Fine-tuning diffusion
746 models for domain-driven image generation using limited data, 2024. URL [https://arxiv.
747 org/abs/2306.14153](https://arxiv.org/abs/2306.14153).
748
749
750
751
752
753
754
755

A APPENDIX

A.1 ADDITIONAL INVERSE PROBLEM SOLVING FIGURES

Figure 8 shows the results of various methods applied to superresolution in the single measurement setting.

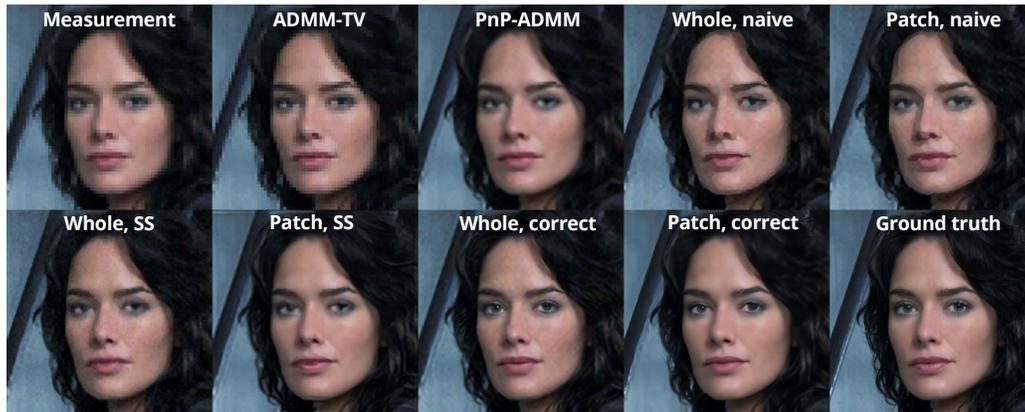


Figure 8: Results of superresolution using self supervised (SS) approach and comparison methods.

Figure 9 shows the results of 20 view CT reconstruction using Algorithm 1. This very sparse view CT recon problem is made more challenging by the lack of any training data. Artifacts can clearly be seen in all the comparison methods. Despite this challenge, reconstructions such as this one can still be useful for medical applications such as patient positioning.

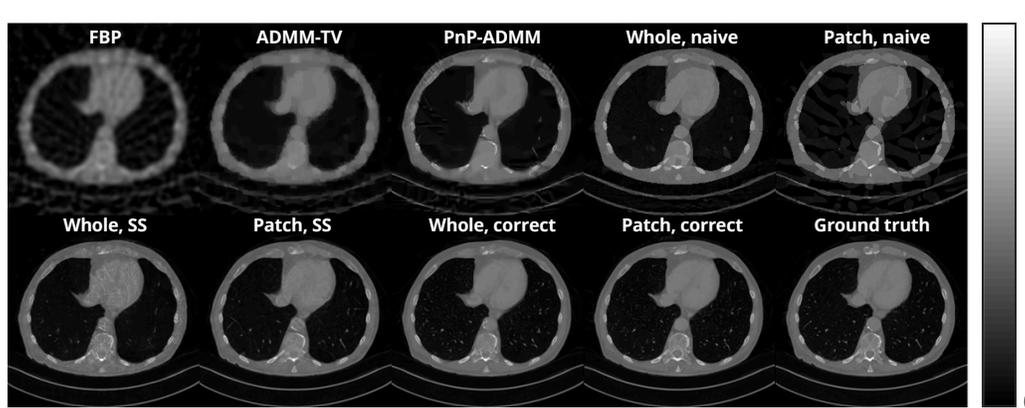
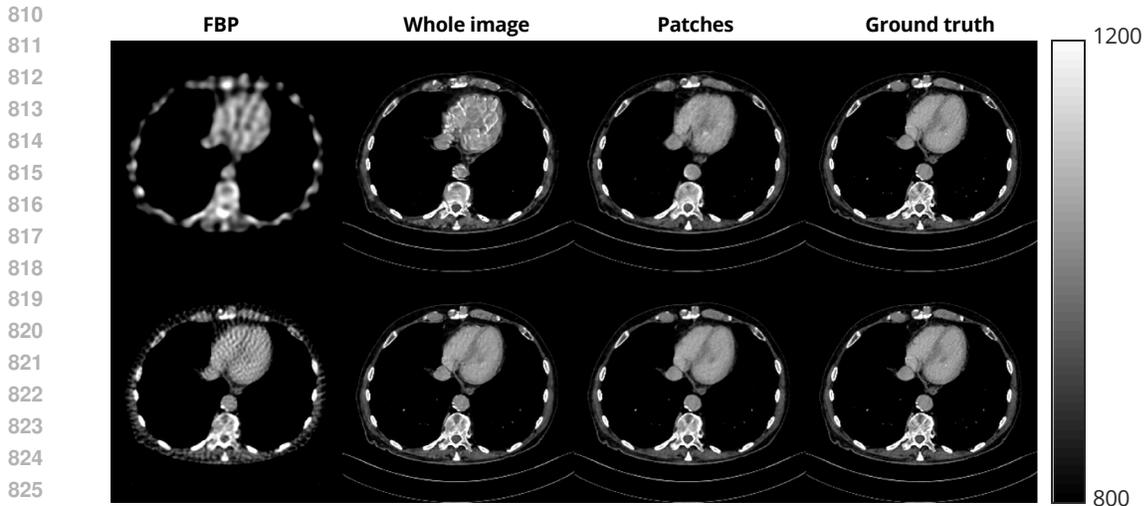


Figure 9: Results of 20 view CT reconstruction in a self-supervised setting. For clarity, the images are plotted on the same scale as the diffusion models were trained.

Figure 10 shows the results of running self-supervised CT reconstruction with 20 views and 60 views where the starting checkpoint was obtained through training on a large (but out of distribution) CT dataset: 10000 LIDC-IDRI slices. Particularly for 20 views, the artifacts from using the whole image model are apparent, while the patch-based model obtains a much higher quality reconstruction. Thus, regardless of whether the starting network has a severely mismatched distribution (ellipses) or a slightly mismatched distribution (different CT dataset), our proposed method outperforms the whole image model.

Figure 11 shows the results of performing 60 view CT reconstruction in an unsupervised manner from checkpoints fine-tuned using the small in distribution CT dataset. The images on the bottom row shows the progressively worsening degradation and increasing number of artifacts resulting



827 Figure 10: Results of CT reconstruction in a self-supervised setting when the starting network was
828 trained on the LIDC dataset. Top row used 20 views and bottom row used 60 views.

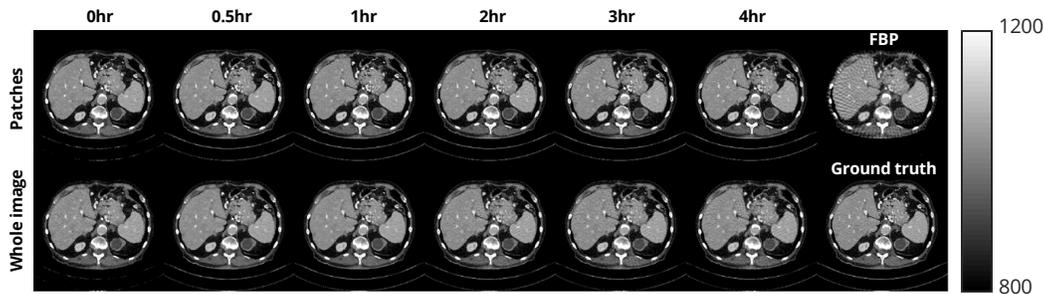
829
830 from overfitting exhibited by whole image model. On the other hand, the top row shows relatively
831 stable performance exhibited by the patch-based model as it is able to avoid overfitting much better.
832

833 Table 3: Performance of fine-tuning on 60 view CT using checkpoints trained for different lengths
834 of time. Best results are in bold.

835
836

Train time (hr)	Patches		Whole image	
	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow
0	33.91	0.921	33.10	0.911
0.5	40.37	0.964	38.39	0.959
1	40.91	0.965	40.54	0.964
2	41.21	0.965	39.19	0.953
3	41.17	0.965	38.31	0.945
4	41.02	0.964	37.67	0.938

837
838
839
840
841
842
843
844

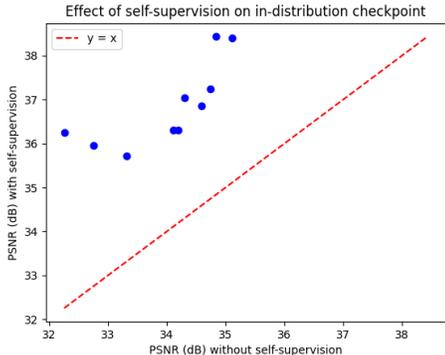


856 Figure 11: Results of 60 view CT recon with networks fine-tuned on 10 in distribution CT images
857 for varying amounts of training time.

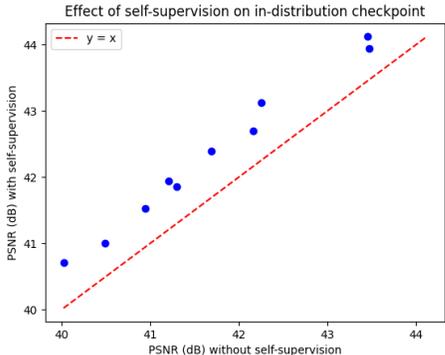
858
859
860 A.2 EFFECT OF SELF-SUPERVISION FOR DIFFERENT DISTRIBUTIONS

861
862 Recall that in the single measurement setting, Algorithm 1 is used to adjust the underlying dis-
863 tribution of the network away from the originally trained OOD data and toward the ground truth
image. We investigated the effect of applying this method even when the network was trained on the

864 in-distribution data. Figures 12 and 13 show the results of this experiment for CT reconstruction,
 865 where each point represents the specific PSNR for one of the images in the test dataset. If the addi-
 866 tional self-supervision step had no effect on the image quality, the points would lie on the red line.
 867 However, in both cases, all of the points are above the red line, indicating that the self-supervision
 868 step of the algorithm improves the image quality even when the network was already trained on
 869 in-distribution data. Furthermore, the improvement is more substantial for for the 20 view case than
 870 the 60 view case, as the predicted clean images $D_{\theta}(x_t|\mathbf{y})$ at each step for the 60 view case are likely
 871 to be more closely aligned with the measurement, so the network refining step becomes less signif-
 872 icant. Importantly, this shows that in practice, one may directly apply Algorithm 1 to solve inverse
 873 problems without knowledge of the severity of the mismatch in distribution between training and
 874 testing data: even when there is no mismatch, the additional self-supervision step does not degrade
 875 the image quality.



887
888
889 Figure 12: PSNR of 20 view CT reconstruction in single-measurement setting using a patch-
 890 based in-distribution network.



891
892 Figure 13: PSNR of 60 view CT reconstruction in single-measurement setting using a patch-
 893 based in-distribution network.

893 Table 5 summarizes the results of using different training datasets while keeping the same test dataset
 894 (AAPM CT images). The distribution shift is greatest when the network is trained on ellipse phan-
 895 tomograms and used to reconstruct the AAPM CT images, so the reconstruction quality is the lowest in
 896 this case. The LIDC dataset consists of CT images which belong to a distribution that is reason-
 897 ably similar to the distribution of AAPM CT images, so when using the network trained on LIDC
 898 images, the quality drop over using an in-distribution network is not substantial. Finally, the im-
 899 provements obtained by using more in-distribution networks is more apparent for the 20 view case
 900 as the measurements are sparser for this case, so the prior plays a larger role in obtaining an accurate
 901 reconstruction.

902 Table 4: Single measurement CT reconstruction results where the initial checkpoint was trained on
 903 LIDC dataset and refined on the fly with the AAPM measurement.

904

Dataset size	CT, 20 views		CT, 60 views	
	PSNR↑	SSIM↑	PSNR↑	SSIM↑
Whole image	35.01	0.894	41.95	0.967
Patches (Ours)	36.34	0.918	42.32	0.972

905
906
907
908
909

Table 5: Performance of patch-based model in single measurement setting for CT reconstruction for different OOD training datasets.

Train time (hr)	20 views		60 views	
	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow
Ellipses	33.77	0.874	41.45	0.966
LIDC	36.34	0.918	42.32	0.970
AAPM	36.82	0.923	42.33	0.970

A.3 ABLATION STUDIES

We performed four ablation studies to evaluate the impact of various parameters on the proposed methods. Similar to the main text, all quantitative results are averaged across the test dataset.

Low rank adaptation. To avoid overfitting to the measurement in self-supervised settings, Barbano et al. (2023) proposed using a low rank adaptation to the weights of the neural network, reducing the number of weights that are adjusted during reconstruction by a factor of around 100. Here we investigate the effect of using different ranks of adaptations on two inverse problems: 60 view CT reconstruction and deblurring. Consistent with Barbano et al. (2023) and Chung & Ye (2024), we only used the LoRA module for attention and convolution layers. We also allowed the biases of the network to be changed.

Tables 6 and 7 show the quantitative results of these experiments, where a rank of “full” represents fine-tuning all the weights of the network. In all cases, using LoRA for this fine-tuning process results in worse reconstructions than simply fine-tuning the entire network. The visual results are especially apparent in Figure 15: the reconstructed image becomes oversmoothed when using LoRA and artifacts become present when using the whole image model. This is likely due to the large distribution shift between the initial distribution of images and target distribution of faces: the low rank adaptation to the mismatched network is not sufficient to represent the new distribution and thus the self-supervised loss function results in smoothed images.

Table 6: Performance of 60 view CT recon using self-supervised network refining with LoRA module. Best results are in bold.

Rank	Parameters (%)	Patches		Whole image	
		PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow
2	1.1	40.37	0.963	39.25	0.952
4	2.0	40.32	0.963	39.10	0.951
8	3.8	40.33	0.963	39.18	0.951
16	7.2	40.32	0.963	39.33	0.953
Full	100	41.45	0.966	40.47	0.957

Table 7: Performance of deblurring using self-supervised network refining with LoRA module. Best results are in bold.

Rank	Parameters (%)	Patches		Whole image	
		PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow
2	1.1	29.31	0.830	29.19	0.811
4	2.0	29.31	0.829	29.35	0.817
8	3.8	29.38	0.831	29.19	0.810
16	7.2	29.31	0.830	29.33	0.815
Full	100	30.34	0.860	29.50	0.831

Effect of network size. In the self-supervised case, another potential method to avoid overfitting is to use a smaller network. We trained networks with differing numbers of base channels (but no other

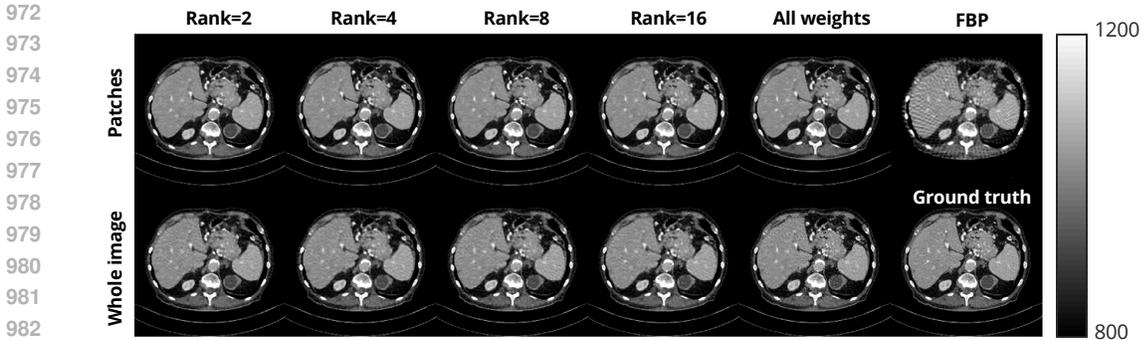


Figure 14: Results of using LoRA module for 60 view CT reconstruction in a single measurement setting. All weights refers to adjusting all the weights of the network at reconstruction time.

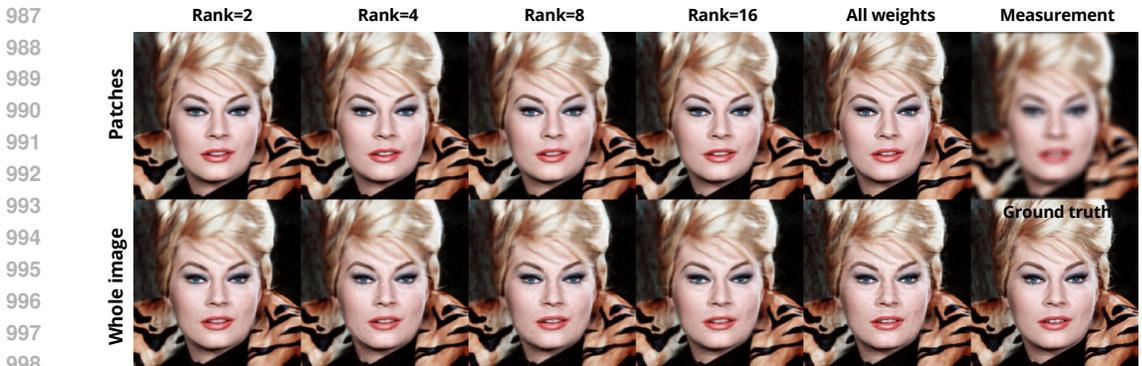


Figure 15: Results of using LoRA module for deblurring in a single measurement setting. All weights refers to adjusting all the weights of the network at reconstruction time.

modifications) on the ellipse phantom dataset and then used Algorithm 1 to perform self-supervised 60 view CT reconstruction. Table 8 shows the quantitative results of this experiment. For both the patch-based model and the whole image model, the network with 128 base channels obtained the best result, so we used this network architecture for all the main experiments. Figure 16 again shows evidence of overfitting in the form of artifacts in the otherwise smooth regions of the organs when using the network with 256 base channels. These artifacts are less obvious in the patch-based model.

Table 8: Performance of 60 view CT recon in a self-supervised manner with networks of different sizes. Best results are in bold.

Base Channels	Parameters (Millions)	Patches		Whole image	
		PSNR↑	SSIM ↑	PSNR↑	SSIM ↑
32	3.4	39.73	0.958	39.69	0.957
64	14	40.37	0.961	40.07	0.958
128	60	41.45	0.966	40.47	0.957
256	217	40.29	0.959	39.28	0.954

Fine-tuning with a larger dataset. To examine the effect of fine-tuning the networks on differing sizes of in-distribution datasets, we started with the same checkpoint trained on ellipses and fine-tuned them using various sizes of datasets consisting of CT images. Each small dataset consisted of randomly selected images from the entire 5000 image AAPM dataset. Next we used these checkpoints to perform 60 view CT reconstruction (without any self supervision). Table 9 shows the results of these experiments, where we also included the results of using the in-distribution network

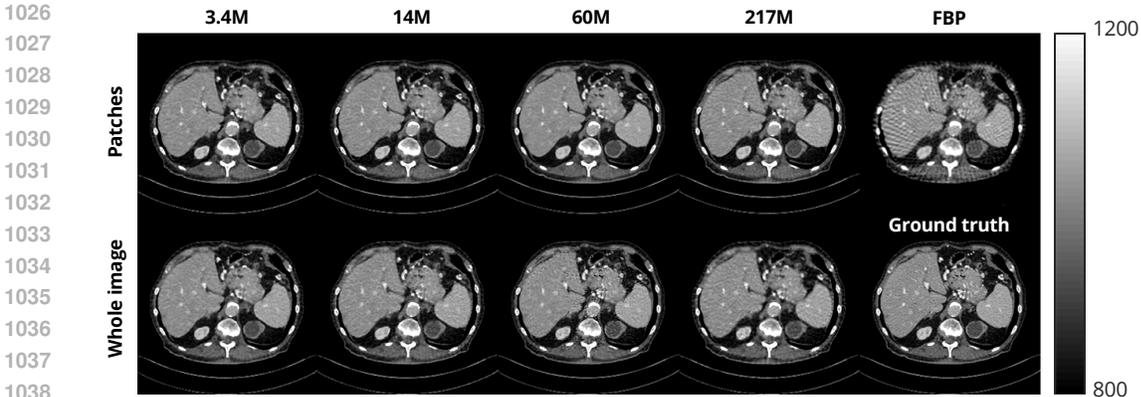


Figure 16: Results of 60 view CT recon using networks with different numbers of parameters in the single-measurement setting. The top numbers show the number of total parameters in the network.

trained on the entire 5000 image dataset. This shows that for a wide range of different fine-tuning dataset sizes our proposed method obtained better metrics than the whole-image model.

We emphasize the difference between the results of Hu et al. (2024), which showed that patch-based models outperform whole image models in cases of limited data, and the results here. Since the networks in Hu et al. (2024) were trained from scratch, more data was required: the smallest datasets used in Hu et al. (2024) contained 144 images. In contrast, we are able to fine-tune networks in our work using only 10 images. Consequently, the training time is also much lower for our work: Figure 4 shows that we fine-tuned a patch-based model in only about 2 hours, whereas Hu et al. (2024) required 12-24 hours to train the patch-based models from scratch. Thus, our results complement the work of Hu et al. (2024) by showing that, compared to whole-image models, patch-based diffusion models easier to train from scratch in settings of limited data, and they are also easier to fine-tune when data is very limited.

Table 9: Performance of fine-tuning on 60 view CT using checkpoints fine-tuned from different dataset sizes. Best results are in bold.

Dataset size	Patches		Whole image	
	PSNR↑	SSIM↑	PSNR↑	SSIM↑
3	40.93	0.964	40.45	0.964
10	41.21	0.965	40.54	0.964
30	41.31	0.966	40.66	0.967
100	41.46	0.967	40.96	0.968
5000*	41.70	0.967	41.67	0.969

Backpropagation iterations during self-supervision. In the single measurement setting, the self-supervised loss is crucial to ensuring that the OOD network output is consistent with the measurement. Backpropagation through the network is necessary to minimize this loss, but too much network refining during this step could lead to overfitting to the measurement and image degradation. We ran experiments examining the effect of the number of backpropagation iterations during each step for the patch-based model and the whole image model. Figures 18 and 19 show that in both cases, performance generally improved when increasing the number of backpropagation iterations and overfitting is avoided. Additionally, the patch-based model always outperformed the whole image model and exhibited more improvement as the number of backpropagation iterations increased. For our main experiments, we used 5 iterations as the improved performance became marginal compared to the extra runtime.

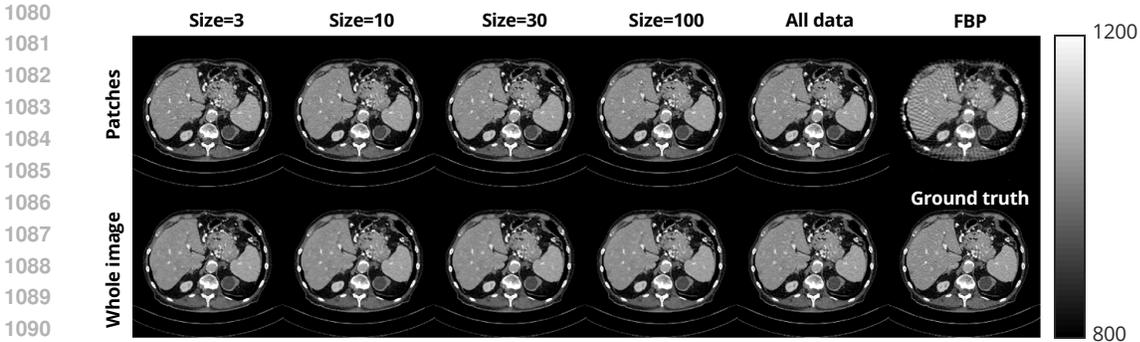


Figure 17: Results of 60 view CT recon in the small dataset setting where the size of the small dataset is varied.

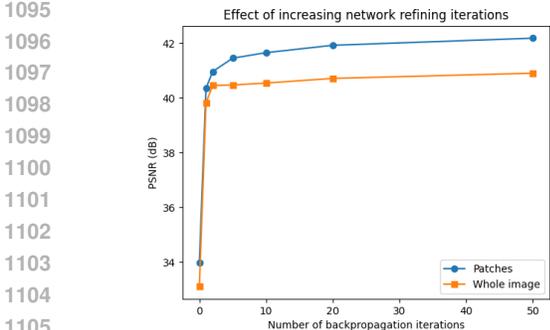


Figure 18: Comparison of PSNR between patch-based model and whole-image model for number of network refining iterations in single measurement setting.

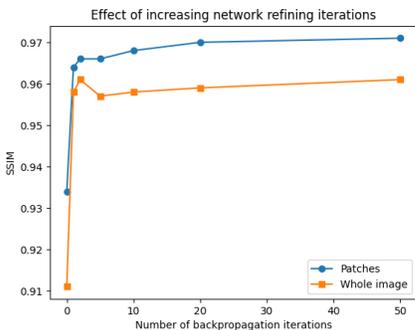


Figure 19: Comparison of SSIM between patch-based model and whole-image model for number of network refining iterations in single measurement setting.

A.4 PHANTOM DATASET DETAILS

We used two phantom datasets of 10000 images each: one consisting of grayscale phantoms and the other consisting of colored phantoms. The grayscale phantoms consisted of 20 ellipses with a random center within the image, each with minor and major axis having length equal to a random number chosen between 2 and 20 percent of the width of the image. The grayscale value of each ellipse was randomly chosen between 0.1 and 0.5; if two or more ellipses overlapped, the grayscale values were summed for the overlapped area with all values exceeding 1 set to 1. Finally, all ellipses were set to a random angle of rotation. The colored phantoms were generated in the same way, except the RGB values for each ellipse were set independently and then multiplied by 255 at the end. Figure 20 shows some of the sample phantoms.

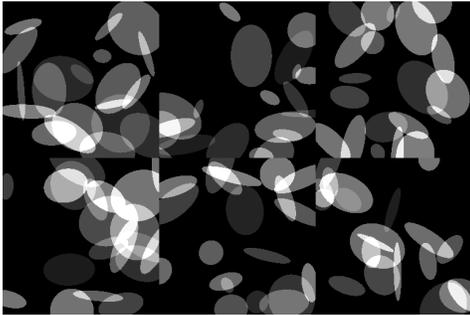
A.5 EXPERIMENT PARAMETERS

We applied the framework of Karras et al. (2022) to train the patch-based networks and whole image networks. Since images were scaled between 0 and 1 for both grayscale images and RGB channels, we chose a maximum noise level of $\sigma = 40$ and a minimum noise level of $\sigma = 0.002$ for training. We used the same UNet architecture for all the networks consisting of a base channel multiplier size of 128 and 2, 2, and 2 channels per resolution for the three layers. We also used dropout connections with a probability of 0.05 and exponential moving average for weight decay with a half life of 500K images to avoid overfitting.

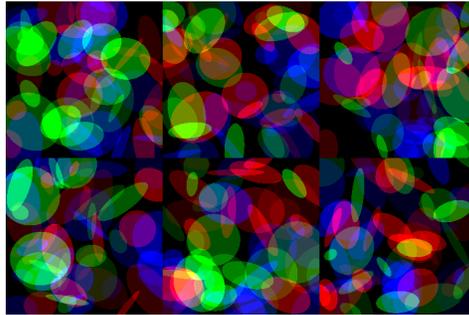
The learning rate was chosen to be $2 \cdot 10^{-4}$ when training networks from scratch and was $1 \cdot 10^{-4}$ for the fine-tuning experiments. For the patch-based networks, the batch size for the main patch size (64×64) was 128, although batch sizes of 256 and 512 were used for the two smaller patch sizes

Table 10: Performance of Algorithm 1 for 60 view CT reconstruction in single measurement setting with different numbers of backpropagation iterations. Best results are in bold.

Backprop iterations	Patches		Whole image	
	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow
0	33.97	0.934	33.10	0.911
1	40.35	0.964	39.81	0.958
2	40.96	0.966	40.45	0.961
5	41.45	0.966	40.47	0.957
10	41.65	0.968	40.54	0.958
20	41.92	0.970	40.71	0.959
50	42.18	0.971	40.90	0.961



(a) Six grayscale phantoms



(b) Six colored phantoms

Figure 20: Six sample grayscale phantoms and colored phantoms used to train the mismatched distribution diffusion models

of 32×32 and 16×16 . The probabilities of using these three patch sizes were 0.5, 0.3, and 0.2 respectively. For the whole image model, we kept all the parameters the same, but used a batch size of 8.

For image generation and inverse problem solving, we used a geometrically spaced descending noise level that was fine tuned to optimize the performance for each type of problem. We used the same set of parameters for the patch-based model and whole image model. The values without the self-supervised loss are as follows:

- CT with 20 and 60 views: $\sigma_{\max} = 10, \sigma_{\min} = 0.005$
- Deblurring: $\sigma_{\max} = 40, \sigma_{\min} = 0.005$
- Superresolution: $\sigma_{\max} = 40, \sigma_{\min} = 0.01$.

The values with the self-supervised loss are as follows:

- CT with 20 and 60 views: $\sigma_{\max} = 10, \sigma_{\min} = 0.01$
- Deblurring: $\sigma_{\max} = 1, \sigma_{\min} = 0.01$
- Superresolution: $\sigma_{\max} = 1, \sigma_{\min} = 0.01$.

Finally, for generating the CT images we used $\sigma_{\max} = 40, \sigma_{\min} = 0.005$.

When running Algorithm 1, we set $K = 10$ for all experiments and $M = 5$ for CT reconstruction and $M = 1$ for deblurring and superresolution. We ran 5 iterations of network backpropagation with a learning rate of 10^{-5} . When using the LoRA module as in the ablation studies (see Tables 7 and 6), we ran 10 iterations of network backpropagation with a learning rate of 10^{-3} .

The ADMM-TV method for linear inverse problems consists of solving the optimization problem

$$\operatorname{argmax}_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - A\mathbf{x}\|_2^2 + \lambda \operatorname{TV}(\mathbf{x}), \quad (13)$$

where $\text{TV}(\mathbf{x})$ represents the L1 norm total variation of $v\mathbf{x}$, and the problem is solved with the alternating direction method of multipliers. For CT reconstruction, deblurring, and superresolution, we chose λ to be 0.001, 0.002, and 0.006 respectively.

The PnP-ADMM method consists of solving the intermediate optimization problem

$$\operatorname{argmax}_{\mathbf{x}} f(\mathbf{x}) + (\rho/2)\|\mathbf{x} - (\mathbf{z} - \mathbf{u})\|_2^2, \quad (14)$$

where ρ is a constant. The values for ρ we used for CT reconstruction, deblurring, and superresolution were 0.05, 0.1, and 0.1 respectively. We used BM3D as the denoiser with a parameter representing the noise level: this parameter was set to 0.02 for 60 view CT and 0.05 for the other inverse problems. A maximum of 50 iterations of conjugate gradient descent was run per outer loop. The entire algorithm was run for 100 outer iterations at maximum and the PSNR was observed to decrease by less than 0.005dB per iteration by the end.

The PnP-RED method consists of the update step

$$\mathbf{x} \leftarrow \mathbf{x} + \mu(\nabla f - \lambda(\mathbf{x} - D(\mathbf{x}))), \quad (15)$$

where $D(\mathbf{x})$ represents a denoiser. The stepsize μ was set to 0.01 for the CT experiments and 1 for deblurring and superresolution. We set λ to 0.01 for the CT experiments and 0.2 for deblurring and superresolution. Finally, the denoiser was kept the same as the PnP-ADMM experiments with the same denoising strength.

Table 11 shows the average runtimes of each of the implemented methods when averaged across the test dataset for 60 view CT reconstruction.

Table 11: Average runtimes of different methods across images in the test dataset for 60 view CT recon.

Method	Runtime (s) ↓
Baseline	0.1
ADMM-TV	1
PnP-ADMM	73
PnP-RED	121
Whole diffusion	112
Whole SS	248
Whole LoRA	329
Patch diffusion	123
Patch SS	289
Patch LoRA	377

1242 A.6 SELF-SUPERVISED INVERSE PROBLEM FIGURES
1243

1244 The following figures show additional examples of self-supervised inverse problem solving.

1245 Figure 21 shows additional example slices of CT reconstruction from 60 views.
1246

1247 Figure 22 shows additional example slices of CT reconstruction from 20 views.

1248 Figure 23 shows additional examples of deblurring with face images.
1249

1250 Figure 24 shows additional examples of superresolution with face images.
1251

1252

1253

1254

1255

1256

1257

1258

1259

1260

1261

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

1289

1290

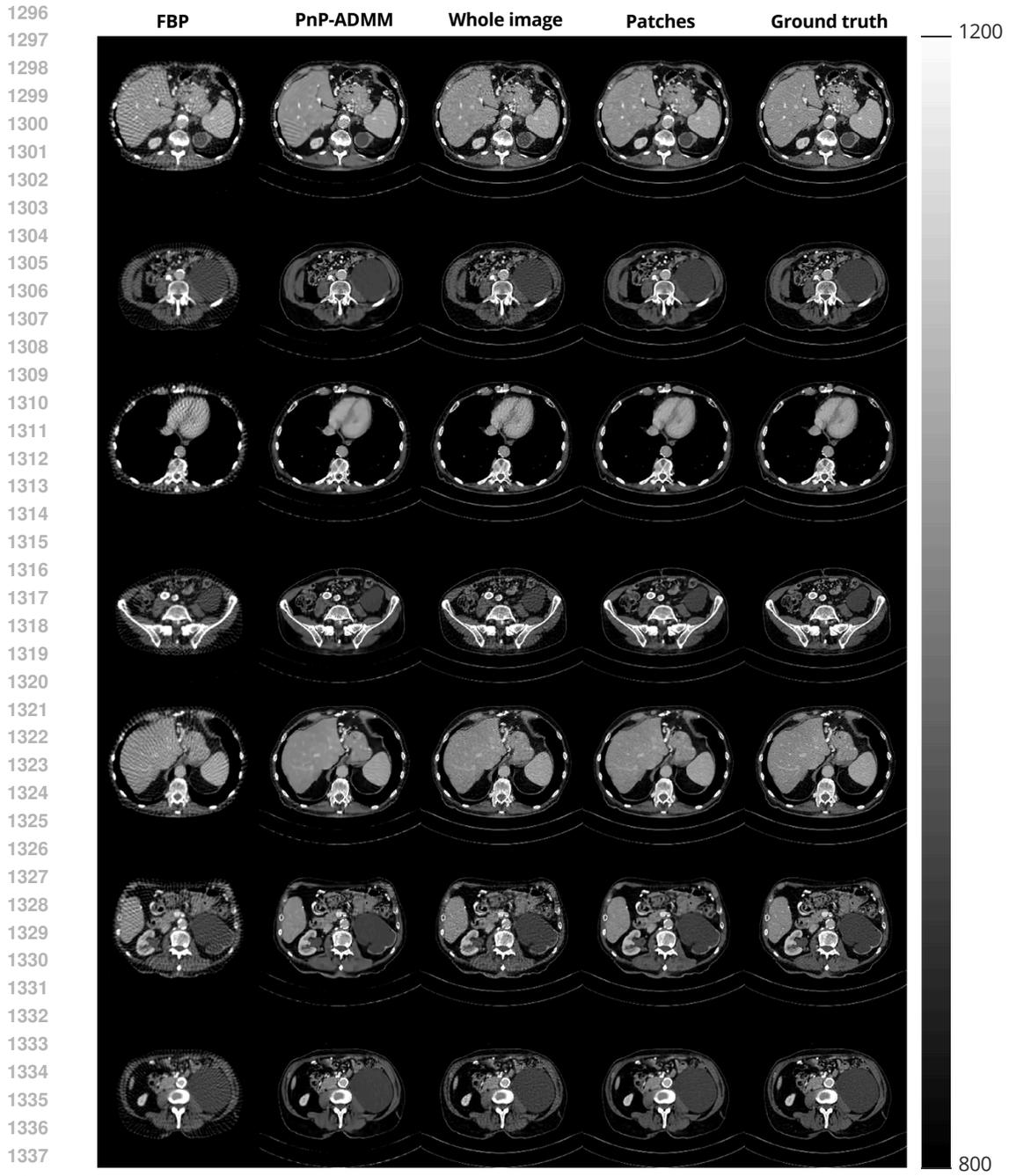
1291

1292

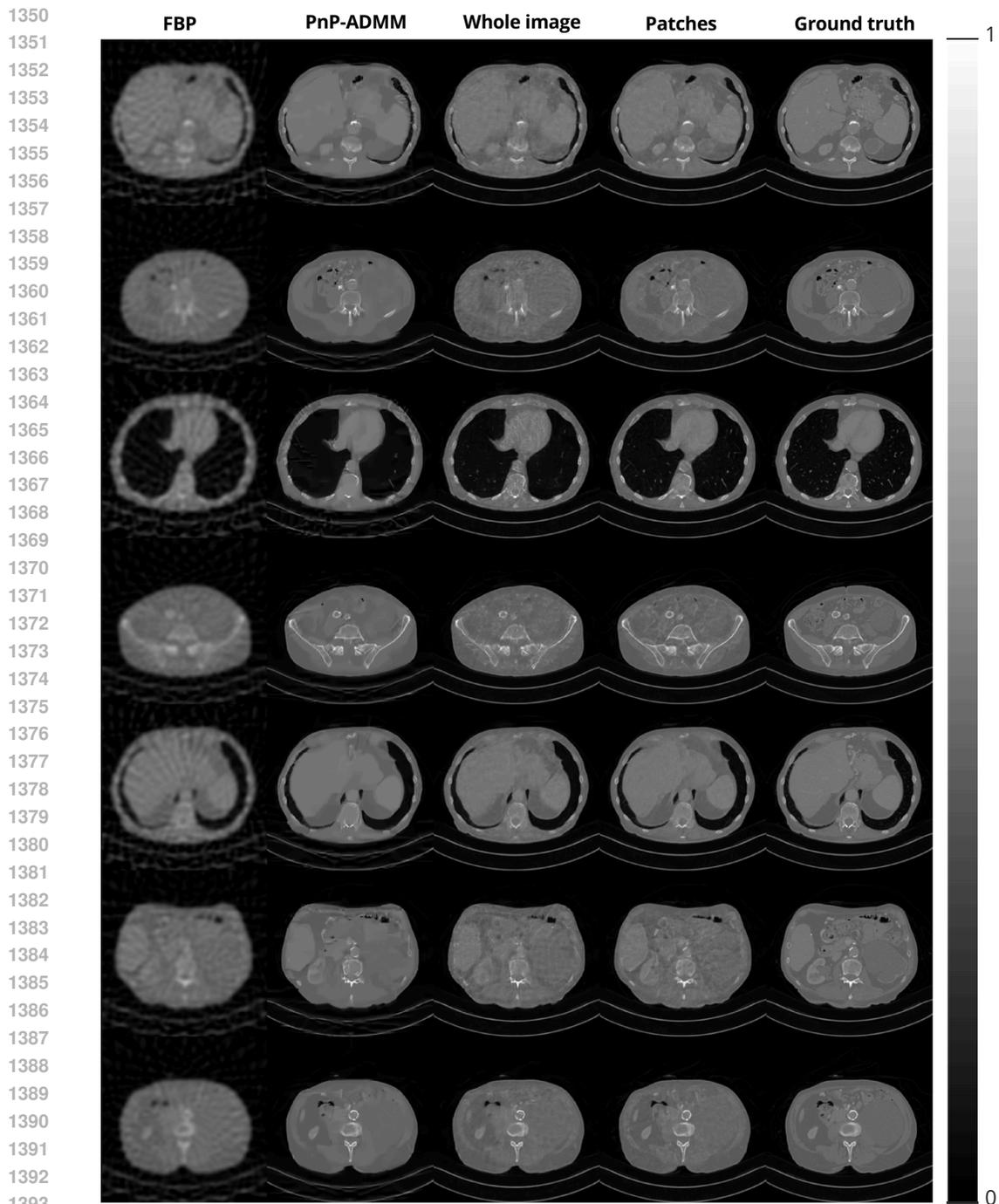
1293

1294

1295



1339 Figure 21: Additional figures for self-supervised 60 view CT recon.
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349



1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

Figure 22: Additional figures for self-supervised 20 view CT recon.

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457



1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511



Figure 24: Additional figures for self-supervised superresolution.

1512 A.7 ADDITIONAL RESULTS

1513 A.7.1 THEORETICAL ANALYSIS OF ALGORITHM 1.

1514 This appendix provides a rough sketch of why Algorithm 1 should work better for the patch-based
1515 model compared to the whole-image model. Firstly, we rewrite (10) in the following form:

$$1516 D(\mathbf{x}) = \sum_c \mathbf{G}'_c D_\theta(\mathbf{G}_c \mathbf{x}, c), \quad (16)$$

1517 where c denotes the location of each patch, \mathbf{G}_c denotes a patch extracting operator that extracts the
1518 patch corresponding to the location c from the whole image \mathbf{x} , and thus \mathbf{G}'_c is an operator that takes
1519 a patch and returns a whole image with the corresponding patch filled in (and the rest of the entries
1520 are zeros). Note that c is input to the patch-based network D_θ through positional encoding.

1521 Now we analyze (11) using this framework. For simplicity we will first drop the conditional expecta-
1522 tion and use the unconditional expectation; later we will analyze the conditional expectation. This
1523 gives us the loss function

$$1524 L(\theta) = \left\| \mathbf{y} - \mathbf{A} \sum_c \mathbf{G}'_c D_\theta(\mathbf{G}_c \mathbf{x}, c) \right\|_2^2. \quad (17)$$

1525 Next, we observe that when \mathbf{A} is the blurring operator or the superresolution operator, \mathbf{A} approx-
1526 imately operates independently on patches of a whole image. For example, particularly when the
1527 patches are not too small, blurring a whole image is approximately equal to blurring the individual
1528 patches and then assembling the patches together. This means that for every patch location c , there
1529 exists a patch operator \mathbf{A}_c such that $\mathbf{G}_c \mathbf{A} \mathbf{x} \approx \mathbf{A}_c \mathbf{G}_c \mathbf{x}$, with \mathbf{x} denoting the whole image. Likewise,
1530 $\mathbf{A} \mathbf{G}'_c \mathbf{x}_c \approx \mathbf{G}'_c \mathbf{A}_c \mathbf{x}_c$, where \mathbf{x}_c denotes the image patch corresponding to the location c .

1531 Hence we replace the loss in (17) with the approximately equal loss

$$1532 L(\theta) = \left\| \mathbf{y} - \sum_c \mathbf{A} \mathbf{G}'_c D_\theta(\mathbf{G}_c \mathbf{x}, c) \right\|_2^2 \approx \left\| \mathbf{y} - \sum_c \mathbf{G}'_c \mathbf{A}_c D_\theta(\mathbf{G}_c \mathbf{x}, c) \right\|_2^2. \quad (18)$$

1533 Furthermore, since the patches in any given iteration of the reconstruction process are nonoverlapping,
1534 we have $\mathbf{y} = \sum_c \mathbf{G}'_c \mathbf{G}_c \mathbf{y}$. Thus the approximate loss is

$$1535 L'(\theta) = \left\| \sum_c \mathbf{G}'_c \mathbf{G}_c \mathbf{y} - \sum_c \mathbf{G}'_c \mathbf{A}_c D_\theta(\mathbf{G}_c \mathbf{x}, c) \right\|_2^2 = \left\| \sum_c \mathbf{G}'_c (\mathbf{G}_c \mathbf{y} - \mathbf{A}_c D_\theta(\mathbf{G}_c \mathbf{x}, c)) \right\|_2^2. \quad (19)$$

1536 Now we have a sum of the form $\|X_1^2 + \dots + X_n^2\|_2^2$, which Song et al. (2024) showed to be upper
1537 bounded by $K(\|X_1\|_2^2 + \dots + \|X_n\|_2^2)$ for a fixed constant K . Applying this inequality gives

$$1538 L'(\theta) \leq K \sum_c \left\| \mathbf{G}'_c (\mathbf{G}_c \mathbf{y} - \mathbf{A}_c D_\theta(\mathbf{G}_c \mathbf{x}, c)) \right\|_2^2 = K \sum_c \left\| \mathbf{G}_c \mathbf{y} - \mathbf{A}_c D_\theta(\mathbf{G}_c \mathbf{x}, c) \right\|_2^2. \quad (20)$$

1539 Table 10 shows that performance is improved by using more backpropagation iterations for the
1540 loss function; hence, although in practice we only perform a fixed number of iterations for speed,
1541 optimally we should aim to reduce the loss $L(\theta)$ to zero. Based on (20), we can instead minimize
1542 the patchwise loss

$$1543 L_c(\theta) = \left\| \mathbf{G}_c \mathbf{y} - \mathbf{A}_c D_\theta(\mathbf{G}_c \mathbf{x}, c) \right\|_2^2. \quad (21)$$

1544 Observe that (21) has the same form as the loss that would be used for refining the network with a
1545 whole image model: $L_w(\theta) = \left\| \mathbf{y} - \mathbf{A} D_\theta(\mathbf{x}) \right\|_2^2$. However, now instead of a loss of a single image,
1546 we now have individual losses of many patches of an image. For example, for the experiments of
1547 Table 1, 25 patches were used to tile each image for each diffusion iteration, so we would have
1548 25 losses. This method of data augmentation provides an explanation for why the patch-based
1549 model obtains better performance than the whole-image model in the single measurement setting.
1550 We additionally note that although the positional encoding input into the network is different for
1551 each patch, the network does not separately learn a distribution for each position, as the weights
1552 are shared across these different positions. This is analogous to the analysis of Dhariwal & Nichol
1553 (2021), where a single diffusion model was trained on the 1000 classes of ImageNet with the class

1566 label of the image being included as an additional input to the network. Since each class only had
 1567 around 1000 images, it would have been very difficult to train a diffusion model on only one of the
 1568 classes, but by training across all the classes at once, a much better network can be trained.

1569 Next, we return to the conditional denoiser in (11). Recall that $D_\theta(\mathbf{x}|\mathbf{y})$ is computed by first com-
 1570 puting $D_\theta(\mathbf{x})$ and then performing $C > 0$ iterations of conjugate gradient descent initialized at
 1571 $D_\theta(\mathbf{x})$. For simplicity, here we analyze the effect of performing gradient descent on $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$.
 1572

1573 From before, we have $\mathbf{y} - \mathbf{A}\mathbf{x} = \sum_c \mathbf{G}'_c(\mathbf{G}_c\mathbf{y} - \mathbf{A}_c\mathbf{G}_c\mathbf{x})$. Making the same assumptions as before
 1574 that the patches are nonoverlapping and \mathbf{A} operates on patches roughly independently, we have

$$1575 \quad \mathbf{A}'\mathbf{y} = \mathbf{A}' \sum_c \mathbf{G}'_c\mathbf{G}_c\mathbf{y} = \sum_c \mathbf{A}'\mathbf{G}'_c\mathbf{G}_c\mathbf{y} \approx \sum_c \mathbf{G}'_c\mathbf{A}'_c\mathbf{G}_c\mathbf{y} \quad (22)$$

1576 Hence, using c and d as patch location indices, the gradient is

$$1577 \quad \mathbf{g} = \mathbf{A}'(\mathbf{y} - \mathbf{A}\mathbf{x}) = \sum_c \mathbf{G}'_c\mathbf{A}'_c\mathbf{G}_c \sum_d \mathbf{G}'_d(\mathbf{G}_d\mathbf{y} - \mathbf{A}_d\mathbf{G}_d\mathbf{x}). \quad (23)$$

1578 Since the patches are nonoverlapping, $\mathbf{G}_c\mathbf{G}'_d = 0$ unless $c = d$, in which case it is the identity.
 1579 Hence, we can combine the double sum into a single sum to get

$$1580 \quad \mathbf{g} = \sum_c \mathbf{G}'_c\mathbf{A}'_c(\mathbf{G}_c\mathbf{y} - \mathbf{A}_c\mathbf{G}_c\mathbf{x}). \quad (24)$$

1581 Note that this takes a very similar form to the full image gradient $\mathbf{g} = \mathbf{A}'(\mathbf{y} - \mathbf{A}\mathbf{x})$: in particular,
 1582 we are simply replacing the full images \mathbf{x} and \mathbf{y} with patches $\mathbf{G}_c\mathbf{x}$ and $\mathbf{G}_c\mathbf{y}$ and then summing over
 1583 all the patches. Thus, $D_\theta(\mathbf{x}|\mathbf{y})$ acts independently on the patches of $D_\theta(\mathbf{x})$. Finally, it is readily
 1584 shown that when using $D_\theta(\mathbf{x}|\mathbf{y})$ in place of $D_\theta(\mathbf{x})$ for (21), we again get a sum of losses over all
 1585 patches.
 1586

1587 Next, we turn to the case of CT reconstruction where \mathbf{A} does not operate approximately indepen-
 1588 dently on patches. Here, we consider the network refining process in Algorithm 1 during the final
 1589 diffusion iteration. By treating the entire conditional denoiser $D_\theta(\mathbf{x}|\mathbf{y})$ as parametrized by one net-
 1590 work $f_\theta(\mathbf{x})$, the loss function simply becomes $L(\theta) = \|\mathbf{y} - \mathbf{A}f_\theta(\mathbf{x})\|_2^2$. This is equivalent to the
 1591 standard DIP formulation. Bager et al. (2020) proved that for this formulation, provided that cer-
 1592 tain technical conditions hold, the minimization problem for θ will converge when gradient descent
 1593 is applied. Most importantly, we need the following definition:

1594 **Definition:** An activation function $\sigma : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is *valid* if it is continuous, monotone, and bounded,
 1595 i.e., there exists $c > 0$ such that $\forall x, \|\sigma(x)\| \leq c\|x\|$.
 1600

1601 For the UNets we used to train the diffusion models, the activation functions were ReLUs and
 1602 sigmoid linear unit (SiLU) which are defined as $\text{SiLU}(x) = x \cdot \sigma(x)$, where $\sigma(x)$ is the sigmoid
 1603 function. Both of these are valid, so the network refining process will converge at the final diffusion
 1604 iteration.
 1605
 1606
 1607
 1608
 1609
 1610
 1611
 1612
 1613
 1614
 1615
 1616
 1617
 1618
 1619

1620 A.7.2 ADDITIONAL EXPERIMENTS

1621
1622 Tables 12 and 13 show LPIPS perception scores for various methods in the single measurement
1623 setting and small dataset setting, respectively. We used the VGG network (Zhang et al., 2018) to
1624 compute these scores and averaged them across all the images in the test dataset. These results show
1625 that our proposed method obtained the images with the best visual image quality.

1626 Table 12: LPIPS score for different methods in single measurement setting.

Method	Deblur	Superresolution
ADMM-TV	0.469	0.542
PNP-ADMM	0.316	0.701
PNP-RED	0.331	0.303
Whole image, naive	0.469	0.469
Patches, naive	0.440	0.465
Self-supervised, whole	0.275	0.339
Self-supervised, patch (Ours)	0.238	0.264
Whole image, correct*	0.194	0.249
Patches, correct*	0.222	0.244

1638
1639 *not available in practice for mismatched distribution inverse problems

1641 Table 13: LPIPS score for different methods in small dataset setting.

Method	Deblur	Superresolution
Best baseline	0.316	0.303
Whole	0.23	0.289
Patch	0.219	0.259

1642
1643
1644
1645
1646
1647
1648
1649 Algorithm 2 provides the pseudocode for the "whole image, correct" method of Table 1.

1650 **Algorithm 2** Whole image recon, no distribution shift

1651
1652 **Require:** $\sigma_1 < \sigma_2 < \dots < \sigma_T, \epsilon > 0, C > 0, \mathbf{y}$
1653 Initialize $\mathbf{x} \sim \mathcal{N}(0, \sigma_T^2 \mathbf{I})$
1654 **for** $i = T : 1$ **do**
1655 Sample $z \sim \mathcal{N}(0, \sigma_i^2 \mathbf{I})$
1656 Set $\alpha_i = \epsilon \cdot \sigma_i^2$
1657 Apply neural network to get $D = D_\theta(\mathbf{x}, \sigma_i)$
1658 Run C iterations of CG for (7) initialized with D
1659 Set $\mathbf{s} = (D - \mathbf{x}) / \sigma_i^2$
1660 Set \mathbf{x} to $\mathbf{x} + \frac{\alpha_i}{2} \mathbf{s} + \sqrt{\alpha_i} \mathbf{z}$
1661 **end for**
1662 **Return** \mathbf{x} .

1663
1664 To demonstrate more statistically significant results, we ran the inverse problems of Table 1 over test
1665 datasets of 25 images in the single measurement setting. Table 14 shows these results: our method
1666 using patches consistently outperforms the baseline and using the whole image model.

1667 Figures 25 and 26 compare the PSNR of the whole-image model and our proposed patch-based
1668 model for each individual image in the 25 image test dataset. The plots show that for all but one
1669 test image in the CT case and for all test images in the deblurring case that the patch-based model
1670 outperformed the whole-image model, illustrating the consistency of our method.

1674
 1675
 1676
 1677
 1678
 1679
 1680
 1681
 1682
 1683
 1684
 1685
 1686
 1687
 1688
 1689
 1690
 1691
 1692
 1693
 1694
 1695
 1696
 1697
 1698
 1699
 1700
 1701
 1702
 1703
 1704
 1705
 1706
 1707
 1708
 1709
 1710
 1711
 1712
 1713
 1714
 1715
 1716
 1717
 1718
 1719
 1720
 1721
 1722
 1723
 1724
 1725
 1726
 1727

Table 14: Extended inverse problem solving results using 25 image test dataset in single measurement setting. Best results are in bold. A two sample t-test shows that for all the experiments, the PSNR and SSIM of the patch-based model exceeds that of the whole image model and the result is statistically significant, all with p-value less than 10^{-7} .

Method	CT, 20 Views		CT, 60 Views		Deblurring		Superresolution	
	PSNR \uparrow	SSIM \uparrow						
PnP-ADMM	30.44	0.843	37.02	0.937	29.93	0.835	26.37	0.771
Whole image, SS	33.51	0.866	40.61	0.959	30.19	0.839	27.78	0.715
Patches, SS (Ours)	34.59	0.891	41.55	0.968	30.95	0.867	28.68	0.839

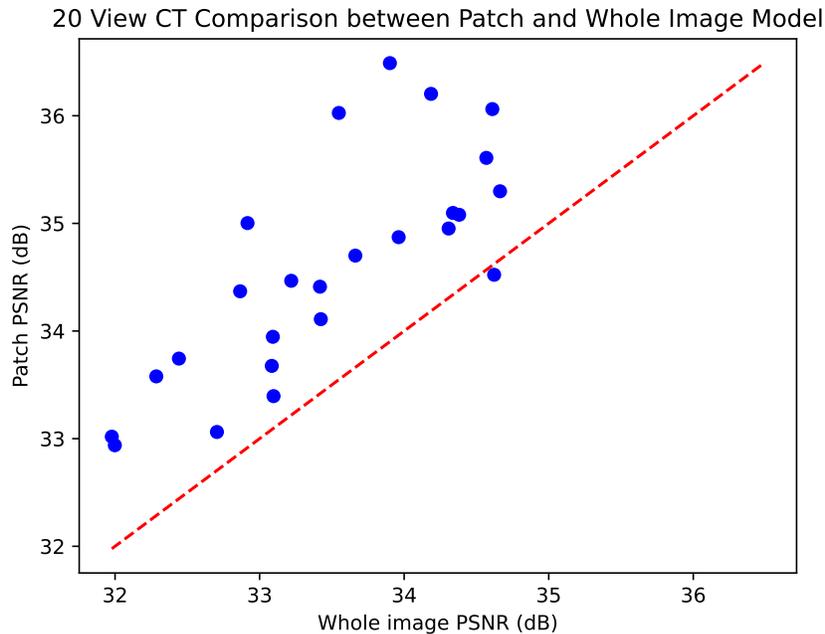


Figure 25: Comparison between PSNR of 20 view CT reconstruction between whole-image model and patch-based model for each image in the test dataset.

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781

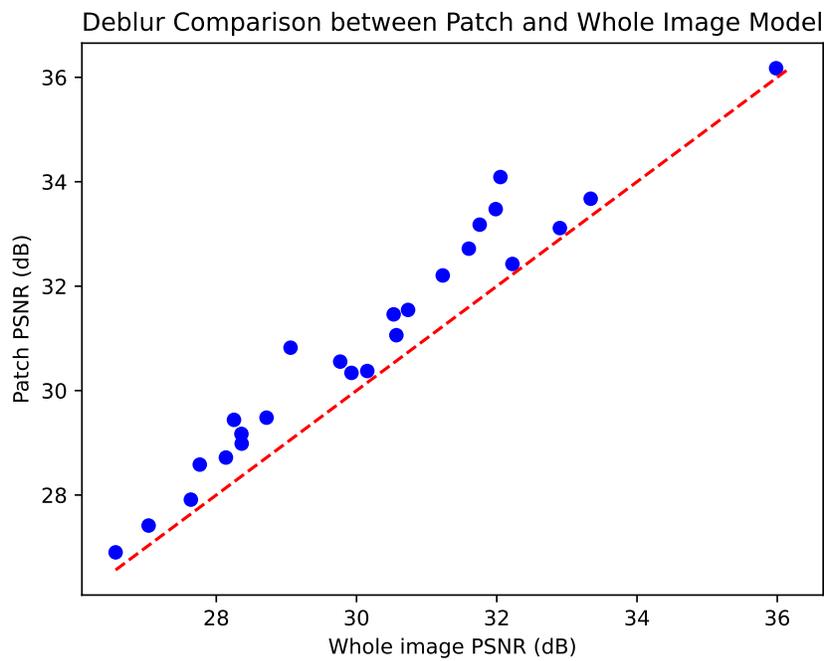


Figure 26: Comparison between PSNR of deblurring between whole-image model and patch-based model for each image in the test dataset.

A.7.3 EXPERIMENTS ON LARGER IMAGES

To show that our method scales to larger images, we ran experiments on 60 view CT reconstruction and deblurring with 512×512 images. For the CT experiments, we still used the AAPM dataset (McCollough et al., 2017) processed in the same way as for Table 1, but kept the slices in their original size of 512×512 . For deblurring, we used the FFHQ dataset (Karras et al., 2019) which contains images of size 512×512 . We scaled each of the RGB channels to between 0 and 1. We used a uniform blur kernel of size 17×17 and added noise with $\sigma = 0.01$. We used the same patch-based networks trained for Table 1 as initializations for these out of distribution experiments. Table 15 shows results of these experiments, where our method obtained the highest quality reconstruction. Although the improvement is modest, note that we trained the patch-based model on phantom images of size 256×256 , which is extremely far out of distribution from the test dataset.

Table 15: Results of inverse problem solving in single measurement setting for 512×512 images.

Method	CT, 60 views		Deblur	
	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow
Baseline	28.33	0.700	24.11	0.649
ADMM-TV	29.36	0.788	28.14	0.760
PnP-ADMM	37.48	0.910	29.77	0.812
Patch, naive	29.32	0.793	26.58	0.749
Patch, SS	37.82	0.919	30.35	0.825

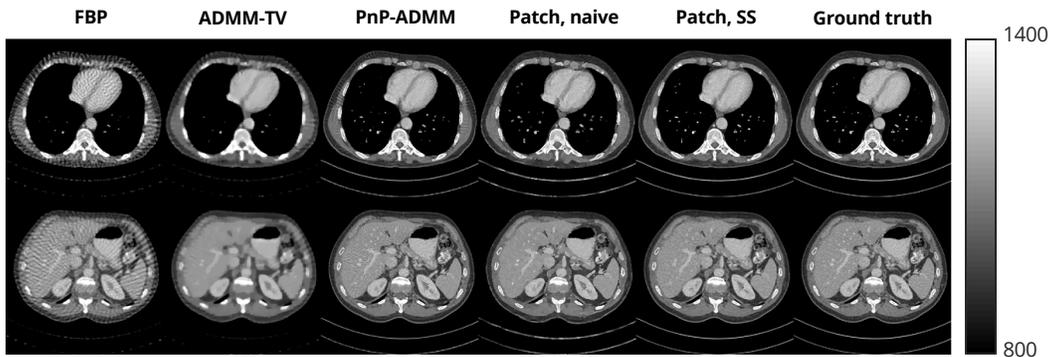


Figure 27: Results of 60 view reconstruction on 512×512 images in single measurement setting displayed in Hounsfield units.

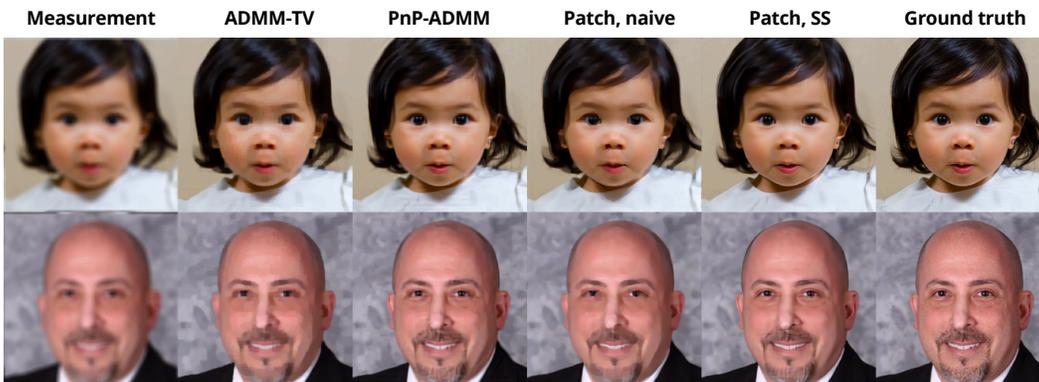


Figure 28: Results of deblurring on 512×512 images in single measurement setting.

1836 A.7.4 ALTERNATE DIFFUSION INVERSE SOLVERS

1837
1838 We ran comparison experiments between various state-of-the-art diffusion inverse solving algo-
1839 rithms including Chung et al. (2024), Wang et al. (2022), and Chung et al. (2023a). Note that
1840 Chung et al. (2024) developed a method targeting 3D inverse problems, but used conjugate gradient
1841 descent to enforce data consistency. Hence, that method closely resembles the methods we used in
1842 our main work (see Tables 1 and 2). Additionally, note that these three methods assume that we
1843 have a network trained on in-distribution data and thus do not refine the network on the fly. Hence,
1844 when applying them in the single measurement setting, we used the network that was trained on the
1845 ellipse phantoms and did not refine the network according to the measurement.

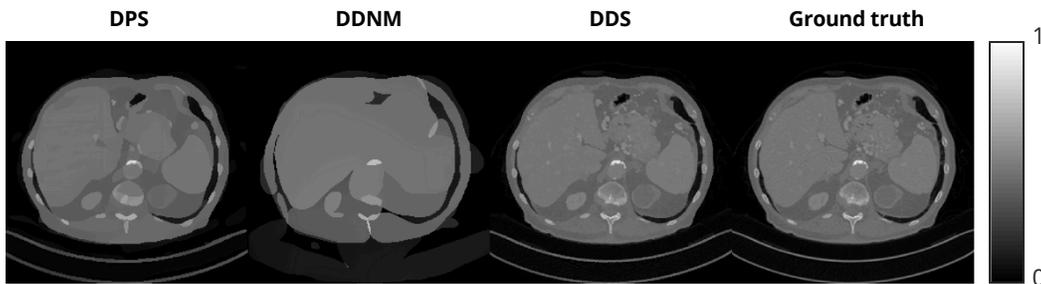
1846 Tables 16 and 17 show the results of using these methods in the single measurement and small
1847 dataset settings, respectively. In the single measurement setting, since these methods do not account
1848 for the mismatched prior, the results are very poor compared to our method that adjusts for the
1849 mismatched prior. Figures 29 and 30 show the failures of these methods visually. Since the network
1850 was trained on ellipse phantoms, the reconstructed images exhibit excessively smooth and rounded
1851 features.

1852 Table 16: Comparison between different diffusion inverse solving methods in single measurement
1853 setting.

Method	CT, 60 views		Deblur	
	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow
DPS (Chung et al., 2023a)	28.28	0.729	27.54	0.864
DDNM (Wang et al., 2022)	23.21	0.572	24.21	0.761
DDS (Chung et al., 2024)	33.97	0.934	26.77	0.782
Ours	41.70	0.967	30.34	0.860

1862 Table 17: Comparison between different diffusion inverse solving methods in small dataset setting.

Method	CT, 60 views		Superresolution	
	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow
DPS (Chung et al., 2023a)	34.24	0.828	28.15	0.814
DDNM (Wang et al., 2022)	26.86	0.860	28.08	0.816
DDS (Chung et al., 2024)	41.21	0.965	28.28	0.830



1882 Figure 29: Visual results of 60 view CT reconstruction using different diffusion inverse solvers in
1883 single measurement setting.

1884
1885
1886
1887
1888
1889

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943



Figure 30: Visual results of deblurring using different diffusion inverse solvers in small dataset setting.

1944 A.7.5 FURTHER BREAKDOWN OF METHODS

1945
1946 It may seem surprising that in Table 2, where a small sample of training data is available, the results
1947 from using the patch-based model are generally worse than those of Table 1, where only a single
1948 measurement is available. This is because for Table 2, after we have fine-tuned the network using
1949 the small dataset, we assume that the network is now in-distribution to the test data, so no network
1950 refining is performed when solving the inverse problem. In contrast, Algorithm 1 is used for Table
1951 1 to refine the network according to the measurement. We illustrate this point further by running
1952 experiments where we applied this network refinement method even after the network was fine-tuned
1953 using the small dataset. Table 18 shows that significant improvements can be made from this further
1954 refinement to fit the measurement. We see that regardless of whether the network was fine-tuned
1955 (FT) with the small dataset, including the self-supervision network refining step results in better
1956 performance.

1957 Table 18: Inverse problem solving results using self-supervision (Algorithm 1) where the initial
1958 network was obtained by fine-tuning on the 10 image dataset.

Method	CT, 20 Views		CT, 60 Views		Deblurring		Superresolution	
	PSNR \uparrow	SSIM \uparrow						
Naive	28.11	0.800	33.10	0.911	25.85	0.742	25.65	0.742
SS only	33.77	0.874	41.45	0.966	30.34	0.860	28.10	0.827
FT only	33.44	0.875	41.21	0.965	29.25	0.840	28.28	0.830
FT and SS	36.43	0.914	42.42	0.971	30.56	0.867	28.38	0.831

1966
1967 We can quantify the degree to which the networks are memorizing the training datasets. We gen-
1968 erated 100 images using both the patch-based model and the whole-image model after fine-tuning
1969 with the small dataset of CT images. For each generation, we compared the NRMSD (normalized
1970 root mean square difference) of the image with each of the 10 CT training images to find the one
1971 with the smallest NRMSD. This was computed by the formula

$$1972 \text{NRMSD}(x, y) = \frac{1}{x_{\max} - x_{\min}} \sqrt{\frac{\sum_{i=1}^n (x_i - y_i)^2}{n}}, \quad (25)$$

1974 where x is the image in the dataset, y is the generated image, and n is the number of pixels in the
1975 image. Table 19 shows the number of images from the 100 generated images where the smallest
1976 NRMSD was less than a certain threshold. This illustrates that the whole image model tended to
1977 memorize the training dataset much more than the patch-based model.

1978
1979 Table 19: Images of the 100 generations that had lowest NRMSD with one of the images in the
1980 small dataset below a certain threshold.

NRMSD (%)	Whole Image	Patches
5	37	0
10	85	4

1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051

A.7.6 ADDITIONAL TRAINING IMAGES

Figures 31 and 32 show some more samples of the ellipse phantoms that were used to train the diffusion models.

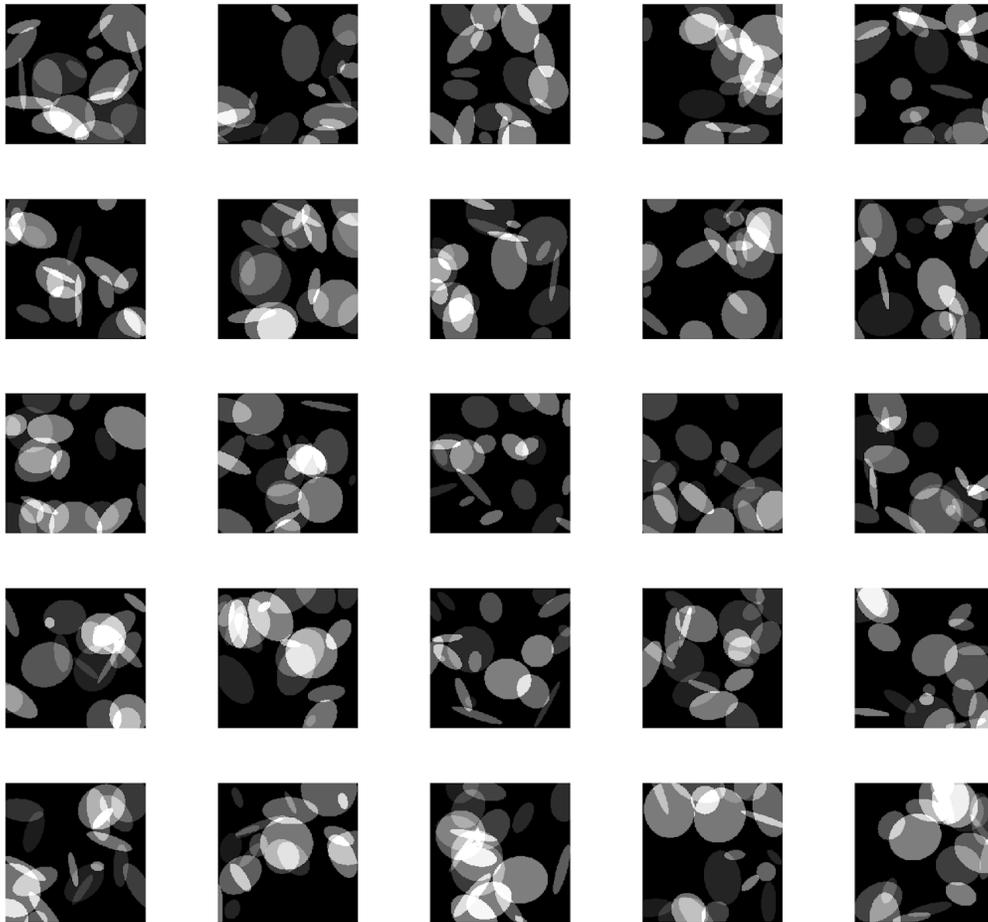


Figure 31: 25 sample grayscale phantoms used to train the mismatched distribution diffusion models.

2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

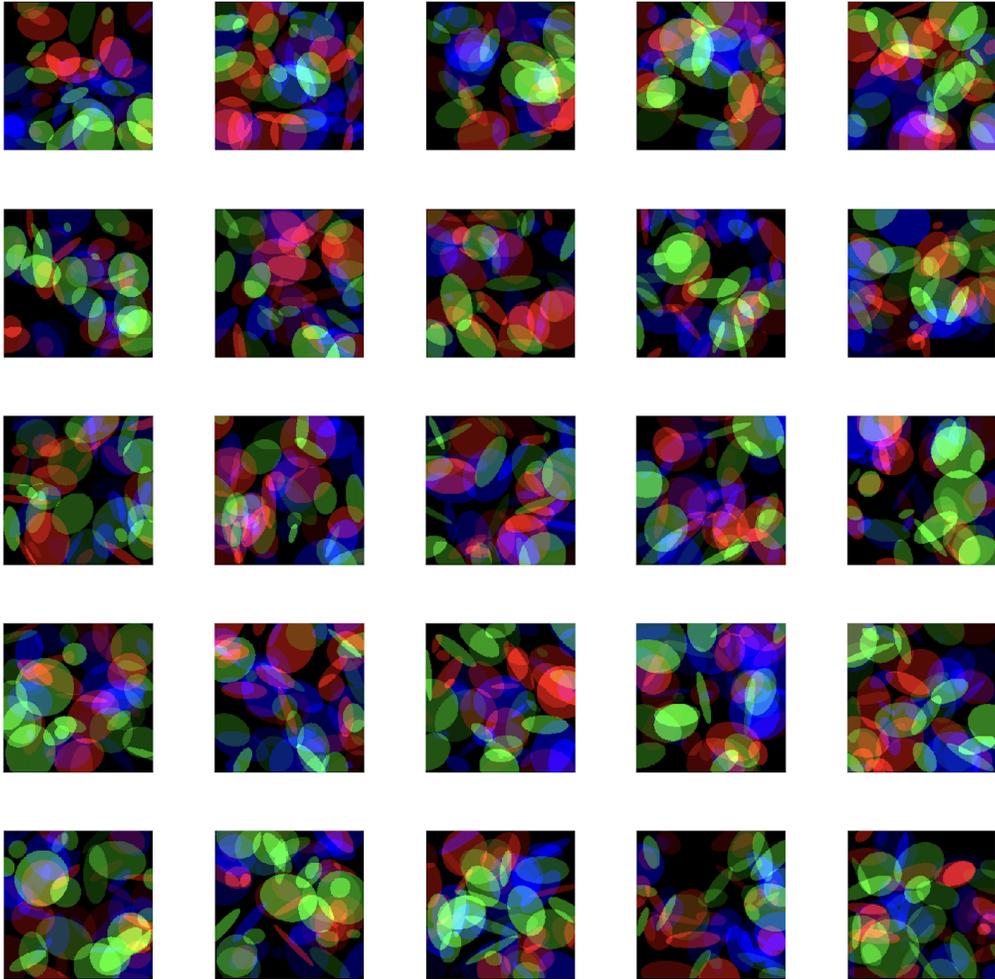


Figure 32: 25 sample color phantoms used to train the mismatched distribution diffusion models.