

Diffusion Models for Large-Scale and 3D Image Reconstruction Problems

by

Jason Hu

A dissertation submitted in partial fulfillment
of the requirements for the degree of
Doctor of Philosophy
(Electrical and Computer Engineering)
in the University of Michigan
2026

Doctoral Committee:

Professor Jeffrey A. Fessler, Co-Chair
Assistant Professor Liyue Shen, Co-Chair
Professor Douglas Noll
Assistant Professor Qing Qu

Jason Hu

jashu@umich.edu

ORCID iD: [0000-0001-5709-1909](https://orcid.org/0000-0001-5709-1909)

© Jason Hu 2026

ACKNOWLEDGMENTS

It would have been impossible for me to have reached this day without the incredible support of many people around me. Their unwavering guidance, encouragement, and insights have impacted the way I think, my goals, and research, and played a significant role in the completion of my PhD thesis.

First of all, I would like to extend my thanks to my advisors Jeff Fessler and Liyue Shen for their continued and supportive mentorship as well as wide domain expertise that have shaped both the direction of my research and my development as an independent thinker. Your constructive criticism and insightful discussions consistently challenged me to refine my ideas, approach problems with greater clarity and rigor, and ultimately helped me develop into the researcher I am today.

Additionally, I would like to thank my fellow researchers and collaborators in both of my research groups. Your discussions, feedback, and shared curiosity created an environment that made tackling difficult problems engaging and rewarding. Thank you for the late night debugging and paper writing sessions as well as coming to celebrate my successes and navigating this journey together. Your support made the day-to-day process of research far more enjoyable.

Finally, my heartfelt gratitude goes out to my friends and family for their constant encouragement, patience, and understanding throughout this journey. Your emotional support was a source of stability and perspective that made it possible for me to persevere through the challenges along the way.

TABLE OF CONTENTS

Acknowledgments	ii
List of Figures	vi
List of Tables	xii
Abstract	xv
Chapter	
1 Introduction	1
1.1 Contributions	2
1.2 Outline	4
2 Background	7
2.1 Inverse Problems	7
2.1.1 Phase Retrieval	8
2.1.2 CT Reconstruction	9
2.1.3 MRI Reconstruction	10
2.2 Diffusion Models	10
2.2.1 Diffusion Inverse Solvers	11
2.2.2 Large-scale Diffusion Models	12
2.2.3 3D Diffusion Models	13
2.3 Distribution Mismatch Problems	13
3 Accelerated Wirtinger Flow with Score-based Image Priors for Holographic Phase Retrieval in Poisson-Gaussian Noise Conditions	16
3.1 Motivation	16
3.1.1 Poisson-Gaussian Log-likelihood for PR	18
3.1.2 Accelerated Wirtinger Flow with Score-based Image Prior	20
3.2 Experiment	22
3.2.1 Experiment Settings	22
3.2.2 Results	24
3.3 Discussion	30
3.4 Conclusion	32
4 Learning Image Priors through Patch-based Diffusion Models for Solving Inverse Problems	34
4.1 Motivation	34

4.2	Methods	37
4.2.1	Patch-wise training	38
4.2.2	Sampling and reconstruction	39
4.3	A Unified Theory of Patch-based Diffusion Models	41
4.3.1	Motivation	41
4.3.2	Methods	41
4.3.3	Examples	46
4.4	Experiments	48
4.5	Conclusion	54
5	DiffusionBlend: Learning 3D Image Prior through Position-aware Diffusion Score Blending for 3D Computed Tomography Reconstruction	55
5.1	Motivation	55
5.2	Methods	57
5.3	Experiments	61
5.4	Conclusion	67
6	Local Patches Meet Global Context: Scalable 3D Diffusion Priors for Computed Tomography Reconstruction	68
6.1	Motivation	68
6.2	Methods	71
6.2.1	Global-Aware Patch Prior	71
6.2.2	Sampling and Reconstruction Algorithm	74
6.3	Results	76
6.3.1	Experimental Setup	76
6.3.2	Unconditional 3D Image Generation	77
6.3.3	Solving Inverse Problems	77
6.3.4	Ablation studies	78
6.4	Conclusion	81
7	Test-Time Adaptation Improves Inverse Problem Solving with Patch-Based Diffusion Models	82
7.1	Motivation	82
7.2	Methods	84
7.2.1	Patch-based prior	84
7.2.2	Diffusion self-supervision	86
7.2.3	Patch-based training	87
7.2.4	Theoretical Analysis	88
7.3	Experiments	90
7.4	Conclusion	98
8	SPAR: Refine a Single Pretrained Diffusion Model to Solve Inverse Problems in Many Modalities	100
8.1	Motivation	100
8.2	Methods	103
8.2.1	Patch-based Prior	103

8.2.2	Model Pretraining	104
8.2.3	Self-Supervised Model Finetuning	105
8.2.4	Theoretical Analysis of Majorizer Loss	107
8.3	Experiments	110
8.4	Conclusion	116
9	Future Work	117
9.1	Improving Existing Works	117
9.1.1	Generalizing Patch-Based Diffusion Models	117
9.1.2	Test-time Adaptation for 3D Problems	118
9.1.3	Diffusion Bridge Network Refining	119
9.1.4	Provable Posterior Sampling	120
9.2	Leveraging Video Diffusion Models	120
	Appendices	122
	Bibliography	185

LIST OF FIGURES

1.1	Dependency flowchart for the chapters of this thesis.	6
3.1	Illustration of Poisson and Gaussian noise statistics in Fourier transform holographic phase retrieval.	17
3.2	Reconstructed images by unregularized methods (Gaussian, Gaussian-Amplitude, Poisson and Poisson-Gaussian) on Histopathology dataset [2], celebA dataset [125] and CT-density dataset. The bottom left/right subfigures correspond to the zoomed in area and the error map for each image. We used $\alpha = 0.035$ and $\sigma = 1$	23
3.3	Reconstructed images on dataset [2]. The bottom left/right subfigures correspond to the zoomed in area and the error map for each image. We used $\alpha = 0.02$ and $\sigma = 1$	24
3.4	Reconstructed images on celebA dataset [125]. The bottom left/right subfigures correspond to the zoomed in area and the error map for each image. We used $\alpha = 0.035$ and $\sigma = 1$	25
3.5	Reconstructed images on CT-density dataset. The bottom left/right subfigures correspond to the zoomed in area and the error map for each image. We used $\alpha = 0.035$ and $\sigma = 1$	25
3.6	Reconstructed images by Gaussian, Poisson and Poisson-Gaussian log-likelihood model with AWFS image prior. Tested on Histopathology dataset [2], celebA dataset [125] and CT-density dataset. The bottom left/right subfigures correspond to the zoomed in area and the error map for each image. α and σ were set to 0.025 and 1, respectively.	26
3.7	Comparison of SSIM and NRMSE varying scaling factor $\alpha \in [0.02, 0.035]$ and STD of Gaussian noise $\sigma \in [0.25, 1.5]$ defined in (3.1).	27
3.8	Reconstructed images by DOLPH [158] and our proposed AWFS method under different σ values, for $\alpha = 0.02$	29
3.9	Comparing AWF vs. WF with NRMSE vs. number of iterations under different noise levels. The curves and shadows represent the mean and standard deviation, respectively.	30
3.10	Reconstructed images by the unregularized Poisson method (the second column) as well as with the AWFS method for different scaling factors α (third to fifth columns). The top and bottom rows show reconstructions from different measurement realizations.	31

4.1	Training the proposed Patch Diffusion Inverse Solver (PaDIS) method. Different sized patches are used in each training iteration. The X and Y position arrays have the same size as the patch and consist of the normalized X and Y coordinates of each pixel of the patch.	35
4.2	Overview of reconstruction process for the proposed Patch Diffusion Inverse Solver (PaDIS) method. Starting at $t = T$, at each iteration we choose a random partition of the zero-padded image and use the neural network trained on patches to get the score function of the entire image. Due to the shifting patch locations, the output image has no boundary artifacts.	35
4.3	Schematic for zero padding and partitioning image into patches	38
4.4	Unconditional generation of CT images. Top row: generation with a network trained on whole image; middle row: patch-only version of [183]; bottom row: proposed PaDIS method.	49
4.5	Results of CT reconstruction. 60 views are used for the top two rows, 20 views are used for the bottom two rows. To better show contrast between organs, we use modified Hounsfield units (HU) in the top figure, while we use the same scale the images were trained on in the bottom figure.	52
4.6	Results of deblurring with Gaussian noise ($\sigma = 0.01$).	53
4.7	Results of superresolution with Gaussian noise ($\sigma = 0.01$).	53
4.8	Results of 60 view CT reconstruction on 512×512 images using modified HU units.	54
5.1	Overview of DiffusionBlend++ compared to previous 3D image reconstruction works. Previous work used a hand-crafted TV term to “regularize” adjacent slices, whereas the proposed approach uses learned diffusion score blending between groups of slices. Here i is the slice index, and t is the reconstruction iteration.	56
5.2	Overview of slice blending process during reconstruction for DiffusionBlend++. At each iteration, we partition the slices of the volume in a different way; slices of the same color are inputted into the network independently. Positional encoding (PE) is also inputted to the network as information about the separation between the slices.	57
5.3	Results of CT reconstruction with 4 views on AAPM dataset, axial view.	62
5.4	Results of DiffusionBlend++ reconstruction with multiple views on AAPM dataset, axial view.	62
5.5	Results of CT reconstruction with 8 views on AAPM dataset, coronal view. DiffusionPatch refers to Algorithm 1 with the same partition for every timestep, and DiffusionBlend+ refers to Algorithm 1 only with partitions of adjacency slices.	66
6.1	Training the 3D patch diffusion model. Noisy 3D patch $\mathbf{G}_c \mathbf{x}_t$ (local details), downsampled volume $\mathbf{D} \mathbf{x}_t$ (global context), and positional encoding are concatenated as inputs to the denoiser network D_θ in each training iteration, to predict the noise ϵ_θ in the input noisy 3D patch. Positional encoding consists of voxel-based coordinates in x, y, z that are normalized to $[-1, 1]$	69

6.2	Schematic illustration for zero padding and partitioning image into 3D patches. Each index i represents one of P^3 possible ways to choose a patch offset tuple.	71
6.3	Unconditional 3D image generation results using the LIDC-IDRI prior. The top row shows axial, coronal, and sagittal slices from a generated volume, and the bottom row shows the corresponding slices from its nearest-neighbor volume in the training dataset. The slice indices for the axial, coronal, and sagittal views are [30, 80, 130, 180, 230], [70, 100, 130, 160, 190], and [60, 130, 160, 190, 210], respectively	77
6.4	Results of our proposed method and comparison methods for 20 view CT recon on LIDC $256 \times 256 \times 256$ dataset. Images are shown in modified Hounsfield units. The top row shows the axial slice and the bottom row shows the sagittal slice from the reconstructed volume.	79
7.1	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.	85
7.2	Dataless setting: Results of 60 view CT reconstruction using self supervised (SS) loss. The display uses modified HU units to show more contrast between organs.	88
7.3	Dataless setting: Results of deblurring using test-time adaptation via the self supervised (SS) loss and several comparison methods.	91
7.4	Results of 60 view reconstruction on 512×512 images in dataless setting displayed in Hounsfield units.	95
7.5	Results of deblurring on 512×512 images in dataless setting.	95
7.6	Comparison of PSNR between patch-based model and whole-image model for overfitting in small dataset setting.	97
7.7	Comparison of SSIM between patch-based model and whole-image model for overfitting in small dataset setting.	97
7.8	Small dataset setting: Results of inverse problem solving. Top row is 60 view CT recon, middle row is deblurring, and bottom row is superresolution.	97
7.9	Unconditional generation of CT images from networks fine-tuned in the small dataset setting. Top two rows were generated with the whole image model; bottom two rows were generated with the patch-based model.	98
8.1	Comparison of traditional domain-specific diffusion models (DM) versus the proposed SPAR method. Left: previous works require separate pre-training for each domain/modality (CT, MRI, CelebA, Places, Dogs). Right: The proposed method enables domain adaptation at testing time through self-supervised learning, avoiding the need for domain-specific pre-training.	101
8.2	Schematic illustration for zero padding and partitioning image into patches. Each index i represents one of P^2 possible ways to choose a patch offset tuple.	103
8.3	Results of our proposed method and comparison methods for various inverse problems. MRI images are complex-valued so we show the pixel magnitudes.	112
8.4	PSNR over time for different loss types in CS-MRI.	115
A.1	Results of PaDIS for 20 view CT reconstruction with different sized patches.	126

A.2	Results for 20 view CT reconstruction with different dataset sizes. Top row shows recon performed by PaDIS; bottom row shows recon performed with the whole-image model.	127
A.3	Results of PaDIS for 20 view CT reconstruction for different positional encoding methods.	128
A.4	Results of PaDIS for 20 view CT reconstruction using different sampling and inverse problem solving algorithms. Top row is with PaDIS and bottom row is with the whole-image model.	129
A.5	Generation of CT images with various acceleration algorithms. Top row shows generation with the EDM sampler [93], middle row uses DDIM [165], bottom row is our proposed method.	136
A.6	Results of extra inverse problem experiments. From top to bottom: 60 view CT, fan beam CT, heavy deblurring. From left to right: baseline, ADMM-TV, whole image diffusion, PaDIS, ground truth.	137
A.7	Unconditionally generated CT images using patch stitching [152] for the top row and patch averaging [142] for the bottom row. Compare this to the images generated by our proposed method in Figure 4.4.	138
A.8	Comparison between PSNR of deblurring between whole-image model and proposed method for each image in the test dataset.	138
A.9	Comparison between PSNR of superresolution between whole-image model and proposed method for each image in the test dataset.	139
A.10	Quantitative results (axial view) of CT reconstruction with 8 views on AAPM dataset for different NFEs and slice blending methods.	145
A.11	Left: Performance of DiffusionBlend++ on more angles, Right: Reconstruction of DiffusionBlend++ with low-dose noise	146
A.12	Comparison of DiffusionBlend++ with classical methods	146
A.13	Results of 3D CT reconstruction with 8 views on LIDC dataset. Top row is axial view, middle row is sagittal view, bottom row is coronal view.	148
A.14	Results of 3D CT reconstruction with 6 views on LIDC dataset. Top row is axial view, middle row is sagittal view, bottom row is coronal view.	148
A.15	Results of 3D CT reconstruction with 4 views on LIDC dataset. Top row is axial view, middle row is sagittal view, bottom row is coronal view.	149
A.16	Results of limited angle 3D CT reconstruction on LIDC dataset. Top row is axial view, middle row is sagittal view, bottom row is coronal view.	149
A.17	Sagittal slices of generated 3D volumes for different numbers of different sampling steps. Generated image quality degrades if the number of steps is reduced too much.	153
A.18	Top row shows the first four axial slices, and the bottom row shows the last four axial slices of the generated volume	154
A.19	Results of our proposed method and comparison methods for 8-view 3D CT reconstruction on AAPM dataset. The top row shows an axial slice and the bottom row shows a sagittal slice from each reconstructed volume.	156

A.20	Results of our proposed method and comparison methods for 60 view CT recon on $256 \times 256 \times 256$ LIDC dataset. Plots are shown in Hounsfield units. The top row shows the axial slice and the bottom row shows the sagittal slice from the reconstructed volume.	157
A.21	Results of our proposed method and comparison methods for 20 view CT recon on $512 \times 512 \times 256$ LIDC dataset. Plots are shown in Hounsfield units. The top row shows the axial slice and the bottom row shows the sagittal slice from the reconstructed volume.	158
A.22	Results of our proposed method for 60 view CT recon on $512 \times 512 \times 256$ LIDC dataset. Plots are shown in Hounsfield units. The top row shows our proposed method and the bottom row shows the ground truth. All slices are axial slices from the same volume.	158
A.23	Unconditional sampling result on $256 \times 256 \times 256$ LIDC-IDRI prior and its nearest neighbor from the training dataset sliced by axial, coronal, and sagittal view.	159
A.24	Unconditional sampling result on $256 \times 256 \times 256$ AAPM prior and its nearest neighbor from the training dataset sliced by axial, coronal, and sagittal view.	159
A.25	Unconditional sampling examples for three different priors (axial slices of 3D volumes).	160
A.26	Unconditional sampling examples for three different priors (coronal slices of 3D volumes).	160
A.27	Unconditional sampling examples for three different priors (sagittal slices of 3D volumes).	161
A.28	Continuous axial slice of a volume sampled from $256 \times 256 \times 256$ LIDC-IDRI prior; every fourth slice is shown	162
A.29	Continuous coronal slice of a volume sampled from $256 \times 256 \times 256$ LIDC-IDRI prior; every fourth slice is shown	180
A.30	Continuous sagittal slice of a volume sampled from $256 \times 256 \times 256$ LIDC-IDRI prior; every fourth slice is shown	181
A.31	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.	182
A.32	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.	182
A.33	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.	182
A.34	Results of 60 view CT recon in the small dataset setting where the size of the small dataset is varied.	183
A.35	Comparison of PSNR between patch-based model and whole-image model for number of network refining iterations in single measurement setting.	183
A.36	Comparison of SSIM between patch-based model and whole-image model for number of network refining iterations in single measurement setting.	183

A.37	Six sample grayscale phantoms and colored phantoms used to train the mismatched distribution diffusion models	183
A.38	Visual results of using various DIS to solve FBCT with 40 views (top) and deblurring (bottom). The same single pretrained network was used in all the experiments. Arrows point to some artifacts in the comparison methods. . . .	184

LIST OF TABLES

3.1	SSIM and NRMSE for Poisson and Poisson-Gaussian likelihoods. Results were averaged across 7 different noise levels by varying $\alpha \in 0.02 : 0.005 : 0.035$ in (3.1), with $\sigma = 1$	27
3.2	SSIM and NRMSE using Poisson Gaussian likelihood with different regularization/image prior approaches. Results were averaged across 7 different noise levels by varying $\alpha \in 0.02 : 0.005 : 0.035$ in (3.1). WFS* runs the same number of iterations as AWFS whereas WFS [†] runs more iterations until convergence.	28
4.1	Comparison of quantitative results on three different inverse problems. Results are averages across all images in the test dataset. Best results are in bold.	51
4.2	Extra inverse problem experiments. Results are averages across all images in the test dataset. Best results are in bold.	52
4.3	Results of 60 view CT reconstruction with 512×512 images. Results are averages across all images in the test dataset. Best results are in bold.	54
5.1	Comprehensive comparison of quantitative results on Sparse-View CT Reconstruction on Axial View for AAPM and LIDC datasets. Best results are in bold.	64
5.2	Comprehensive comparison of quantitative results on Sparse-View CT Reconstruction on Sagittal View for AAPM and LIDC datasets. Best results are in bold.	64
5.3	Comprehensive comparison of quantitative results on Sparse-View CT Reconstruction on Coronal View for AAPM and LIDC datasets. Best results are in bold.	65
5.4	Comprehensive comparison of quantitative results on Limited-Angle CT Reconstruction on All Views for AAPM and LIDC datasets. Best results are in bold.	65
5.5	TV values of different reconstruction algorithms on the AAPM test set	66
5.6	Effectiveness of Blending Modules, Sagittal view performance on AAPM	67
6.1	PSNR (dB) comparison of various methods for CT reconstruction on different datasets and resolutions with best results in bold.	78
6.2	Per volume runtimes of different methods for 8 view CT recon.	79
6.3	FID comparison with and without downsample channel	80
6.4	PSNR in dB of different patch and volume size	80
6.5	PSNR in dB of different K , with best results in bold	81

7.1	Comparison of quantitative results on four different inverse problems in the dataless setting with the self-supervised loss (SS). Results are averages across all images in the test dataset. Best results for practical use are in bold.	93
7.2	Results of inverse problem solving in dataless setting for 512×512 images.	94
7.3	Comparison of results for using self-supervised (SS) diffusion models in small dataset setting. The diffusion model is first trained on ellipse phantoms and then fine-tuned with the small dataset. Best results are in bold.	96
7.4	Comparison of results for using diffusion models trained on thousands of in-distribution images to solve inverse problems. Best results are in bold. Self-supervised refinement of network weights is beneficial even in this setting.	96
7.5	CT reconstruction results where the initial checkpoint was trained on LIDC dataset and refined on the fly with the AAPM measurement.	99
8.1	Comparison of quantitative results on four different medical imaging inverse problems. Results are averages across all images in the test dataset. Best results for practical use are in bold.	111
8.2	Comparison of quantitative results on different natural imaging inverse problems: deblurring and super-resolution. The second column indicates whether the method can handle images from different modalities. Best results for practical use are in bold.	111
8.3	Comparison between diffusion inverse solvers using the same network for CT reconstruction and image deblurring.	114
9.1	Comparison of different diffusion-based methods	120
A.1	Effect of patch size P on CT reconstruction	123
A.2	Dataset size effect on CT reconstruction	123
A.3	Positional encoding effect for CT reconstruction	125
A.4	Dataset size effect on CT reconstruction	125
A.5	Average runtimes of different methods across images in the test dataset for 20 view CT recon.	132
A.6	Effectiveness of Blending Modules, Sagittal view performance on AAPM	144
A.7	Axial PSNR for 8 view SVCT recon for different NFEs	145
A.8	Axial PSNR for 8 view SVCT recon for different slice jump frequencies	145
A.9	Wall times of various methods for 8 view 3D CT reconstruction	146
A.10	Standard Deviation of Performance on Sparse-View CT Reconstruction on Sagittal View for AAPM and LIDC datasets. Best results are in bold.	147
A.11	3D prior modeling methods	150
A.12	Model hyperparameters.	153
A.13	Trade-offs on different patch sizes	155
A.14	Performance of 60 view CT recon using self-supervised network refining with LoRA module. Best results are in bold.	164
A.15	Performance of deblurring using self-supervised network refining with LoRA module. Best results are in bold.	164

A.16	Performance of 60 view CT recon using test-time adaptation with networks of different sizes. Best results are in bold.	165
A.17	Performance of fine-tuning on 60 view CT using checkpoints fine-tuned from different dataset sizes. Best results are in bold.	165
A.18	Performance of Algorithm 5 for 60 view CT reconstruction in single measurement setting with different numbers of backpropagation iterations. Best results are in bold.	166
A.19	Average runtimes of different methods across images in the test dataset for 60 view CT recon.	169
A.20	Comparison of using the proposed SPAR method while training only with CelebA images versus using a mixed dataset. Best results are in bold.	172
A.21	Comparison between different thresholds σ' when using SPAR image reconstruction. Best results are in bold.	173
A.22	Comparison of using the traditional loss versus the proposed majorizer loss for each of the inverse problems.	173
A.23	Comparison of using different K for 60 view parallel beam CT recon. Average runtime for each image is also shown.	174
A.24	Average runtimes of different methods across images in the test dataset for 60 view CT recon.	175
A.25	LPIPS scores for deblurring on LSUN dataset and superresolution on CelebA dataset.	175
A.26	Performance of deblurring on LSUN dataset and superresolution on CelebA dataset with noise level $\sigma = 0.05$. Networks were trained with CelebA images. Best results are in bold.	176
A.27	Results from using DiffPIR to accelerate reconstruction process while applying our network finetuning method for deblurring.	177

ABSTRACT

We live in a world where imaging systems are ubiquitous. From the cell phones in our pockets to our cars and doorbells and on to telescopes and medical scanners, imaging has changed how we share, document, and understand our world. There is an increasing demand to make these systems more efficient by producing higher-quality images with fewer and fewer resources. When we need an image that is of higher quality than can be directly constructed from its measurements, we have an inverse problem and must rely on reconstruction schemes. These reconstruction algorithms fill the gap between the measured samples and the desired high-quality images by leveraging a model that attempts to recover image contents while reducing artifacts.

Diffusion models are trained on large image datasets to learn an underlying prior. This iterative process, which starts with pure noise and is gradually denoised to eventually obtain a clean image, can be modified to generate or restore the clean image that is also enforced to be consistent with the measurement, thereby solving an inverse problem. However, diffusion models are computationally expensive to train and require large, clean training datasets.

To tackle these challenges, this thesis develops methods of using diffusion models to solve inverse problems. In particular, we develop a novel patch-based diffusion prior that allows for image generation and reconstruction using a memory efficient network trained on small patches of an image. We later improve this diffusion prior by including global image context, allowing for better coherence between patches and accelerated sampling. We also develop a 3D diffusion prior that learns the joint distribution between multiple slices of a 3D volume. These methods allow us to solve 3D inverse problems that would otherwise be computationally infeasible to handle with traditional diffusion models. Furthermore, we demonstrate the efficacy of patch-based diffusion models for solving out-of-distribution inverse problems, reducing the need for a large training dataset. This is done by using a self-supervised loss function that refines the network at reconstruction time, ensuring the network output is consistent with the measurement. Motivated by the relative inaccuracy of trained diffusion models for low noise levels, we propose a novel self-supervised majorizer loss function which is theoretically and empirically shown to learn the score function at low noise levels better. Coupled with a channel concatenation operator that allows a network to accept images with differing numbers of channels, our method enables a single diffusion model to solve inverse problems in multiple modalities. The proposed methods are applied to problems in computed tomography, magnetic resonance imaging, and deblurring, among others, showing their wide generalizability with broad scientific and engineering applications.

CHAPTER 1

Introduction

In the real world, many objects cannot be observed directly and must instead be measured through a device or system. A doctor who wishes to view the internal anatomy of a patient can perform a computed tomography (CT) or magnetic resonance imaging (MRI) scan of the patient. A photographer can use a camera to capture the dynamic view of a bustling city. An astronomer studying black holes must use advanced telescopes to capture the desired view. Once the measurement has been obtained by the corresponding imaging system, one must then recover the original desired image from the measurements. This type of problem is known as an *inverse problem*.

In all of these cases, the system that is used to capture the image is inherently noisy. Improving the hardware of the imaging system is one possible solution to mitigate this problem, but is expensive in practice and can lead to slow acquisition times. When we wish to perform fast and cheap imaging, the process will naturally be lossy. Hence, the same measurement will correspond to many possible images. The central problem of computational imaging and this thesis is to develop methods to identify an image close to the ground truth and obtain a high-quality reconstruction.

Traditionally, inverse problems have been solved by reformulating them as an optimization problem consisting of two terms: a *data fidelity* term and a *regularization* term. The data fidelity term ensures that the image is consistent with the measurement, but that term alone is not enough to obtain a high-quality reconstruction due to the lossy nature of the measurement system. Hence, a regularization term, designed to enforce commonly held beliefs about possible clean images, is crucial. This regularization term ranges from simple handcrafted cost functions such as total variation loss or sparsity metrics [29] to complicated data-driven priors that are learned by neural networks [171].

When the measurements are very sparse or the noise level is high, it is especially important to develop a good prior to obtain a high-quality reconstruction. Generative models, and in particular diffusion models, have risen to popularity in recent years, as they are

capable of generating realistic clean images from pure noise. Diffusion models firstly involve adding progressively increasing levels of noise to clean images and training a neural network to denoise them. This process results in a network that has learned the underlying distribution of the training image dataset. Next, to generate a clean image using the trained network, we start with pure noise and use the network to gradually denoise the image, eventually obtaining a clean image belonging to the data distribution. Since diffusion models directly learn the image prior, they are particularly useful for solving inverse problems where the measurements are sparse or noisy, as the prior can be used to fill in missing information. Many methods have been developed for adapting diffusion models from their original purpose of sampling from the prior (generating clean images with no other information) to sampling from the posterior (generating a clean image that is consistent with a measurement, i.e., solving inverse problems), and these methods have shown superior performance over traditional methods.

Despite these successes, diffusion models still face several difficulties when solving inverse problems. Training whole-image diffusion models requires large quantities of clean training data that is often unavailable in many applications. Furthermore, as distributions of images are very high dimensional and complex, deep neural networks with hundreds of millions of parameters are needed to learn the distributions well, making the training process a slow and memory-intensive task. Finally, when solving difficult nonlinear inverse problems, diffusion models exhibit instability as methods of enforcing data fidelity often perturbs the image away from the desired noisy image manifold. Consequently, it is very challenging to apply diffusion models to solve large-scale and 3D imaging problems, applications where training data is limited, and nonlinear inverse problems. This thesis will focus on methods that allow diffusion models to overcome these challenges.

1.1 Contributions

In this work, we develop various novel methods targeted at the aforementioned challenges of diffusion models. These include methods of training and applying different types of diffusion models to solve a large variety of inverse problems, ranging from nonlinear inverse problems to 2D and 3D medical and natural imaging tasks. We also explore different problem settings where memory and training data may be limited. The contributions of this work are summarized as follows:

- **Nonlinear inverse problems.** We develop a novel algorithm to solve the nonlinear inverse problem of phase retrieval with the presence of Poisson and Gaussian noise. We show how diffusion models can be used for this problem and apply an acceleration method to reduce the number of diffusion iterations needed, a crucial

addition to an algorithm that would otherwise be very slow due to the complicated Poisson-plus-Gaussian likelihood. Our simulation experiments show that modeling both types of noise yields better performance than only using one type of noise and the diffusion prior is a better choice than other traditional methods [115].

- **Large-scale diffusion models.** We tackle large-scale imaging problems by developing a novel patch-based diffusion model approach for learning image priors. This patch-based prior can learn the prior of an entire image from solely the patches of the image, so that both at training and reconstruction time, it suffices to input patches of the image into the network. Hence, the method reduces the memory requirement and training time. Extensive experiments show that it exhibits superior performance over traditional whole-image diffusion models for a variety of medical and natural inverse problems [82]. We also develop a framework for general patch-based priors that encompasses several other works and derive convergence guarantees. For 3D image reconstruction, we develop a method that learns a 3D image prior by first learning the joint distributions of multiple slices of the 3D volume and then blending these distributions together. We also improve upon the patch-based prior model by incorporating global volumetric information, leading to a trained network that can unconditionally generate high-resolution volumes. We apply our methods to various problems in CT reconstruction where the measurements are very sparse and demonstrate that this 3D prior yields higher-quality reconstructions than previous state-of-the-art (SOTA) methods [162, 199].
- **Limited training data.** We apply our patch-based diffusion models to the mismatched distribution setting via test-time adaptation. In this setting, we are given a pretrained diffusion model, but the test data belongs to a different distribution from the training distribution, so we must refine the network at test time. Simulations show that our patch-based model is more easily adaptable to the mismatched distribution, is less prone to overfitting in the case of very limited data, and yields better performance at solving inverse problems than whole image models [81]. We then generalize this method using a channel operator and novel majorizer test-time adaptation loss function so that we can use a single pretrained network to solve inverse problems in many modalities [79].
- We raise various avenues of exploration for future work, including extending our unified framework for patch-based diffusion models, applying diffusion bridges for faster image reconstruction, and developing other test-time adaptation methods

that are more generalizable to challenging 3D inverse problems. We also discuss the possibility of applying video diffusion models to solve 3D image reconstruction.

1.2 Outline

To start, Chapter 2 provides background on the mathematical framework used for modeling inverse problems. It also introduces more specifics of the inverse problems that are covered in the rest of the work, including phase retrieval, CT reconstruction, and image enhancement. Next we introduce the relevant background of diffusion models and their relevance to inverse problems before moving onto mismatched distribution problems.

Chapter 3 details our novel algorithm by which we use diffusion models to solve phase retrieval under realistic system and noise models [115]. We demonstrate with exhaustive experiments that the realistic models and score-based image prior yield better reconstructions.

Chapter 4 provides the patch-based diffusion model framework that will be referenced in the subsequent chapters [82]. We develop a diffusion model prior that can be applied to large images by breaking up an image into patches and learning the prior on the patches separately. We then show that while using less memory and training time, in cases of limited training data, our patch-based approach outperforms traditional whole-image diffusion models in various inverse problems.

Chapter 5 further looks at large-scale diffusion models by tackling the 3D CT reconstruction problem with diffusion models [162]. We develop a method to learn 3D diffusion image priors and use it to solve very sparse-view CT reconstruction problems as well as limited-angle CT reconstruction.

Chapter 6 develops a fully 3D prior capable of unconditionally generating high-resolution 3D CT volumes. This is done by incorporating downsampled whole volumes into the patch-based prior developed in Chapter 4 to provide global information into the prior model. We also empirically demonstrate that this improved prior can be leveraged to obtain higher quality CT reconstructions.

In Chapter 7, we move to problems where the training data and test data have different distributions [81]. Specifically, we cover two settings: the first in which we are only given the measurement without any in-distribution training data and must refine the diffusion network on the fly, and the second in which we are given a very small dataset that can be used to fine-tune the network. We show that, when using patch-based diffusion models, desirable performance qualities are achieved compared to whole-image models.

Chapter 8 generalizes the work of the previous chapter even more by pretraining a single diffusion network on a mixed dataset of different modalities that can then be used to solve inverse problems in different modalities such as MRI and natural imaging. Our novel majorized self-supervised loss function, coupled with a solid theoretical foundation, is the driving force behind this method’s flexibility.

Finally, Chapter 9 outlines future work that could be follow-up directions based on this dissertations’s research. Broadly, this consists of extending the theory of patch-based diffusion models by running a more systematic set of experiments and developing the theory for more general (and accelerated) sampling methods. Furthermore, we highlight the potential of diffusion bridges to greatly accelerate inverse problem solving while preserving the generalizability of unconditional diffusion model methods, as well as possible methods of leveraging pretrained video diffusion models to solve 3D image reconstruction problems.

The appendix provides additional results, both quantitative and visual, ablation studies, and theoretical derivations. Summarizing, Figure 1.1 provides an overview of the dependencies of the chapters of this thesis and two of the main topics covered, large-scale imaging and distribution shift problems.

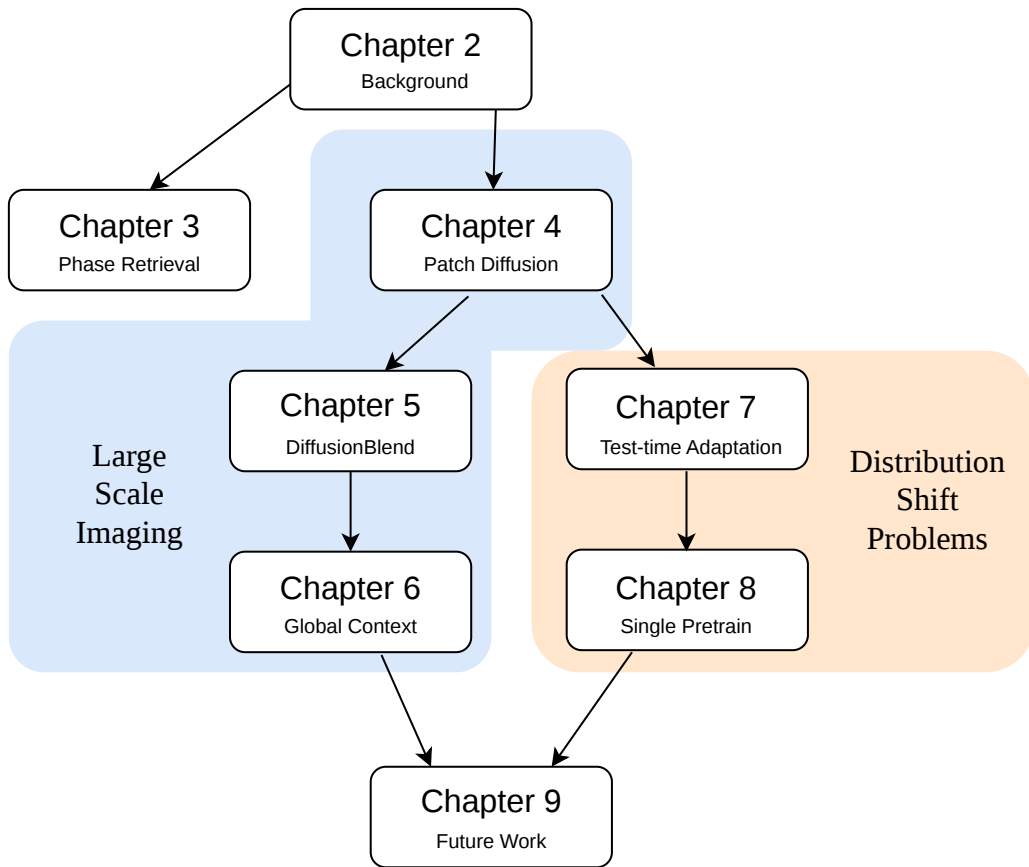


FIG 1.1 – Dependency flowchart for the chapters of this thesis.

CHAPTER 2

Background

This chapter begins with a mathematical formulation of inverse problems, a broad range of specific examples of inverse problems, and general frameworks that are commonly used to solve these types of problems. Next, we highlight specific inverse problems that are addressed in detail in the rest of this thesis and existing methods that can be used to solve them.

2.1 Inverse Problems

Inverse problems are ubiquitous in image processing, and aim to reconstruct an image \boldsymbol{x} from a measurement \boldsymbol{y} , where $\boldsymbol{y} = \mathcal{A}(\boldsymbol{x}) + \epsilon$, \mathcal{A} represents a forward operator, and ϵ represents random unknown noise. In the presence of noise, and particularly when \mathcal{A} is nonlinear or ill-conditioned, the same \boldsymbol{y} can correspond to multiple possible images \boldsymbol{x} . Therefore, to obtain a high-quality reconstruction of \boldsymbol{x} , it is necessary to enforce a prior on \boldsymbol{x} . More precisely, by Bayes' Rule, $p(\boldsymbol{x}|\boldsymbol{y})$ is proportional to $p(\boldsymbol{x}) \cdot p(\boldsymbol{y}|\boldsymbol{x})$. By making assumptions on the noise ϵ , the distribution $p(\boldsymbol{y}|\boldsymbol{x})$ is readily derived, so it remains to establish a prior $p(\boldsymbol{x})$.

Traditionally, total variation (TV) has been used as a regularizer [35]. This regularizer has a denoising effect on the image while preserving sharp boundaries. Similarly, wavelet based methods [52] transform the image to the wavelet domain where it becomes simpler to denoise the image without distorting its main features. More recently, plug and play methods involve alternating between using a denoiser to denoise the image and a data fidelity update step to ensure that the image conforms with the measurement [34]. Block matching 3D (BM3D) is a fast and effective general purpose denoising algorithm, so its use in inverse problems has been studied [136].

In recent years, data-driven methods have risen in popularity in signal and image processing [122, 118, 80, 194, 196]. In particular, for solving inverse problems, when large

amounts of training data is available, a learned prior can be much stronger than the hand-crafted priors used in traditional methods. For instance, plug and play and regularized by denoising methods [172, 197, 123, 83, 121, 33, 156] that involve pretraining a denoiser and applying it at reconstruction time have shown superior performance over using fixed denoisers. These methods have the advantage over supervised deep learning methods such as [89, 104, 169, 187] that the same denoiser may be applied to solve a wide variety of inverse problems.

The next subsections summarize several inverse problems investigated in this thesis.

2.1.1 Phase Retrieval

Phase retrieval (PR) is a nonlinear inverse problem, where the goal is to recover a signal from the (square of) magnitude-only measurements that are corrupted by noise [85]. This problem has applications in astronomy [50], X-ray crystallography [137], optical imaging [157], Fourier ptychography [17, 208, 193, 176] and coherent diffractive imaging (CDI) [106]. For example, in holographic CDI, a coherent beam source illuminates a sample of interest and a reference. When the beam hits the sample, it generates secondary electromagnetic waves that propagate until they reach a detector. By measuring the photon flux, the detector can capture and record a diffraction pattern. This pattern is roughly proportional to the square of Fourier transform magnitude of electric field associated with the illuminated objects [13, 14]. Recovering the structure of the sample from the diffraction pattern generated by the mix of sample and reference is a nonlinear inverse problem known as holographic PR. One approach to this problem is maximum a posteriori (MAP) estimation:

$$\hat{\boldsymbol{x}} = \operatorname{argmax}_{\boldsymbol{x} \in \mathbb{R}^N} p(\boldsymbol{x} | \boldsymbol{y}, \bar{\boldsymbol{b}}, \boldsymbol{A}, \boldsymbol{r}) = \operatorname{argmin}_{\boldsymbol{x} \in \mathbb{R}^N} g(\boldsymbol{x}; \boldsymbol{A}, \boldsymbol{y}, \bar{\boldsymbol{b}}, \boldsymbol{r}) + h(\boldsymbol{x}), \quad (2.1)$$

where \boldsymbol{x} denotes a real latent image to recover, \boldsymbol{y} is the recorded measurement vector, $\bar{\boldsymbol{b}}$ denotes the mean of background measurements, and $\boldsymbol{A} \in \mathbb{C}^{M \times N}$ denotes the system matrix in holographic PR, where M denotes the number of measurements and N denotes the dimension of \boldsymbol{x} . The known reference image \boldsymbol{r} provides additional information to reduce the ambiguity of $\hat{\boldsymbol{x}}$; using an extended reference is a common technique in holographic CDI [155, 70]. Following Bayes' rule, we denote $g(\boldsymbol{x}) = -\log p(\boldsymbol{y}, \boldsymbol{A}, \boldsymbol{r} | \boldsymbol{x})$ and $h(\boldsymbol{x}) = -\log p(\boldsymbol{x})$ as the data fidelity term and the regularization term, respectively. For simplicity we assume \boldsymbol{x} is real, so $\nabla g(\boldsymbol{x})$ denotes the real component of the gradient of log-likelihood. Hence, all priors are trained with real-valued datasets and are applied to the real components of \boldsymbol{x} as well. The method can be extended to complex images.

The term $g(\mathbf{x})$ depends on the noise model, which is generally chosen to be either Gaussian or Poisson. For methods that assume the elements of \mathbf{y} follow independent Gaussian distributions $\mathbf{y} \sim \mathcal{N}(|\mathbf{Ax}|^2 + \bar{\mathbf{b}}, \Sigma)$, where $\Sigma = \sigma^2 \mathbf{I}$, the data fidelity term $g(\mathbf{x})$ in (2.1) becomes $g_{\text{Gau}}(\mathbf{x}) \triangleq \|\mathbf{y} - \bar{\mathbf{b}} - |\mathbf{Ax}|^2\|_2^2$. To solve the corresponding MAP optimization problem, a popular method is Wirtinger flow (WF) [28, 88, 26, 161] using the Wirtinger gradient: $\nabla g_{\text{Gau}}(\mathbf{x}) = 4\mathbf{A}' \text{diag}\{|\mathbf{Ax}|^2 - \mathbf{y} + \bar{\mathbf{b}}\}\mathbf{Ax}$. To determine an appropriate step size for the Wirtinger gradient, one can use its Lipschitz constant or methods such as empirical trial and error, backtracking line search, or observed Fisher information [117]. To further accelerate WF, one can use Nesterov’s momentum methods [139] or optimized gradient methods [97], leading to the accelerated Wirtinger flow (AWF) [193, 21, 58, 63] that is commonly used in solving PR problems. Apart from WF, other methods such as matrix-lifting [30, 27, 157], error reduction (ER) [64], hybrid input-output (HIO) [62], majorize-minimize (MM) [146] and alternating direction method of multipliers (ADMM) [119] have also been proposed.

On the other hand, the Poisson ML model assumes $\mathbf{y} \sim \text{Poisson}(|\mathbf{Ax}|^2 + \bar{\mathbf{b}})$, so that $g(\mathbf{x})$ in (2.1) has the form: $g_{\text{Pois}}(\mathbf{x}) \triangleq \mathbf{1}'(|\mathbf{Ax}|^2 + \bar{\mathbf{b}}) - \mathbf{y}' \log(|\mathbf{Ax}|^2 + \bar{\mathbf{b}})$. Similar to the Gaussian case, one can also apply WF [116] with $\nabla g_{\text{Pois}}(\mathbf{x}) = 2\mathbf{Ax} \odot (\mathbf{1} - \mathbf{y} \oslash (|\mathbf{Ax}|^2 + \bar{\mathbf{b}}))$, where \odot and \oslash denote element-wise multiplication and division, respectively.

2.1.2 CT Reconstruction

Computed tomography (CT) is a medical imaging technique that allows a 3D object to be imaged by shooting X-rays through it [59]. The measurements consist of a set of 2D projection views obtained from setting up the source and detector at different angles around the object. By definition, \mathbf{y} is the (known) set of projection views, \mathcal{A} is the (in most cases assumed to be) linear forward model of the CT measurement system, and \mathbf{x} is the unknown image. The CT reconstruction problem then consists of reconstructing \mathbf{x} given \mathbf{y} . To reduce the radiation dose delivered to the patient, sparse-view CT uses a smaller fraction of X-rays compared to the full-view CT [159]. Additionally, limited-angle CT is useful in cases where patients may have mobility issues and cannot use full-angle CT scans [24]. In these cases, \mathbf{y} contains far less information than \mathbf{x} , so this becomes a challenging inverse problem.

Traditional methods for solving this include regularization-based methods that enforce a previously held belief on \mathbf{x} and likelihood based methods [59, 175, 192, 37]. Data-driven methods have shown tremendous success in signal and image processing in recent years [122, 118, 80, 194]. In particular, for solving inverse problems, when large amounts of

training data is available, a learned prior can be much stronger than the hand-crafted priors used in traditional methods [197, 123]. For past few years, many deep learning-based methods have been proposed for solving the 3D CT reconstruction problem [89, 104, 169, 195]. These methods train a convolutional neural network, such as a U-Net [89], that maps the partial-view filtered backprojection (FBP) reconstructed image to the ground truth image, that is, full-view CT reconstruction. However, these methods often generate blurry images and generalize poorly for out-of-distribution data [4].

2.1.3 MRI Reconstruction

Magnetic resonance imaging (MRI) is another medical imaging technique that uses magnetic fields to obtain images of organs, tissues, and other structures [75]. MRI scanners acquire the Fourier components of the image of interest, called the k-space. However, this acquisition procedure is slow, so acceleration methods such as undersampling the Fourier components have been developed, which can lead to aliasing in reconstructed images. To solve this problem, modern MRI scanners use multiple coils to acquire the Fourier components, providing additional spatial information [145, 67]. Additionally, compressed sensing (CS) MRI [132, 133] improves the quality of the reconstructed images by using suitable sampling patterns. Similar to CT reconstruction and other inverse problems, methods of solving CS MRI reconstruction problems include classical optimization based methods [132], plug and play methods [1], and diffusion model methods [46]. Notably, for deep learning methods, one often represents the complex MRI data using two separate channels consisting of the real and imaginary parts [41].

2.2 Diffusion Models

Diffusion models consist of defining a forward stochastic differential equation (SDE) that adds noise to a clean image [168]: for $t \in [0, T]$, $\mathbf{x}(t) \in \mathbb{R}^d$, we have

$$d\mathbf{x} = -(\beta(t)/2) \mathbf{x} dt + \sqrt{\beta(t)} d\mathbf{w}, \quad (2.2)$$

where $\beta(t)$ is the noise variance schedule of the process. The distribution of $\mathbf{x}(0)$ is the data distribution and the distribution of $\mathbf{x}(T)$ is (approximately) a standard Gaussian. Then, image generation is done through the reverse SDE [3]:

$$d\mathbf{x} = (-\beta(t)/2 - \beta(t)\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)) dt + \sqrt{\beta(t)} d\bar{\mathbf{w}}. \quad (2.3)$$

To train a network to approximate the score function $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)$, it is necessary to first discretize the SDE into a sequence of positive noise scales (for white Gaussian $\mathcal{N}(0, \sigma_k^2)$): $\sigma_1 > \sigma_2 > \dots > \sigma_K$, with σ_K being small enough so that noise of this level does not visibly affect the image, and σ_1 depending on the application. Score matching can be used to train a noise conditional score network (NCSN) [179, 166] as follows:

$$\hat{\theta} = \operatorname{argmin}_{\theta} \sum_{k=1}^K \mathbb{E}_{\mathbf{x}, \tilde{\mathbf{x}}} \left[\left(s_{\theta}(\mathbf{x}, \sigma_k) - \frac{\mathbf{x} - \tilde{\mathbf{x}}}{\sigma_k^2} \right)^2 \right],$$

where $\mathbf{x} \sim p(\mathbf{x})$, $\tilde{\mathbf{x}} \sim \mathbf{x} + \mathcal{N}(0, \sigma_k^2 \mathbf{I})$. (2.4)

With enough data, the neural network $s_{\theta}(\mathbf{x}, \sigma)$ is expected to learn the distribution $p_{\sigma}(\mathbf{x}) = \int p(\mathbf{x}) p_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x}) d\mathbf{x}$ where $p_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x}) = \mathcal{N}(\mathbf{x}, \sigma^2 \mathbf{I})$. Once the neural network has been trained to approximate the score function, one can start with noise and run the reverse SDE to obtain samples from the learned data distribution.

2.2.1 Diffusion Inverse Solvers

Diffusion models serve as a strong image prior as they can generate entire images from pure noise. Most methods that use diffusion models to solve inverse problems involve writing the problem as a conditional generation problem [53, 120, 40] or as a posterior sampling problem [43, 39, 31, 181, 95]. In the former case, the network requires the measurement \mathbf{y} (or an appropriate resized transformation of \mathbf{y}) during training time. Thus, for these task-specific trained models, at reconstruction time, that network is useful only for solving that specific inverse problem. In contrast, for the posterior sampling framework, the network learns an unconditional image prior for \mathbf{x} that can help solve any inverse problem related to \mathbf{x} without retraining.

More specifically, methods that use an unconditionally trained diffusion model to solve inverse problems generally involve alternating between updating the image using the trained score network and updating the image to be consistent with the measurement. These two steps are generally detached, so many choices of sampling methods and SDE solvers are possible for the former step while changing the method by which data fidelity is enforced. For example, [166] uses Langevin dynamics with a simple gradient descent step toward the measurement. Recognizing that the image lies on a noisy data manifold (as opposed to the clean data manifold) during intermediate timesteps, [39] improves on this method by performing gradient descent with respect to the expectation of the clean image at each iteration. As these methods require a large number of iterations, DDIM [165] was developed as an acceleration method that relies on estimating the clean image

at each iteration to guide the noisy image toward the clean image manifold. Building off this method, [181] and [95] rely on hard constraint enforcement via projection onto the manifold where data consistency is satisfied while using DDIM as the sampling method, greatly reducing the number of sampling steps needed when solving inverse problems. Finally, conjugate gradient descent has proven to be useful for enforcing data fidelity when used in conjunction with DDIM when the measurements are assumed to be compressed but noiseless [41]. However, even when the number of sampling steps is large, these methods all struggle to truly sample from the posterior, as they all rely on approximations to the posterior distribution at intermediate timesteps. More recently, [87] proposed a method using resampling and Monte Carlo methods to estimate the posterior distribution at each timestep which was theoretically shown to sample from the posterior distribution when the number of particles is large.

2.2.2 Large-scale Diffusion Models

To reduce the computational burden, latent diffusion models [150] have been proposed, aiming to perform the diffusion process in a much smaller latent space, allowing for faster training and sampling. However, that method requires a pretrained encoder and decoder [57] for a fixed dataset, so it must be retrained for different datasets, and it still requires large amounts of training data. Furthermore, solving inverse problems in the latent space is difficult due to the nonlinearity of the encoder and decoder, so existing methods are slow and unstable, particularly for nonlinear inverse problems [163].

Patch-based diffusion models [183, 55] focus on image generation while training only on patches, reducing the memory requirement by only requiring patch inputs to the score network. However, for high-quality image generation, [183] still requires inputting the entire image into the network at inference time. Furthermore, [55] only allows for unconditional image generation, and it is unclear how to extend this method to solve inverse problems. Supervised patch-based diffusion methods [189, 142] are task specific and do not learn an unconditional image prior that can be applied to all inverse problems. Other patch-based methods [18, 130, 144] learn an unconditional image prior but require the whole image as an input during inference time. Finally, work has been done to perform sampling faster [165, 93, 128], which is unrelated to the training process.

Other strategies for reducing computational costs of diffusion models include multi-level diffusion models that train networks on progressively higher resolution images and employ a similar strategy at inference time to generate high resolution images [69, 9]. Alternative network architectures that reduce the number of parameters have also been

proposed in both the pixel space and latent space [72, 144], but these networks are still difficult to scale up to very large images or 3D images.

2.2.3 3D Diffusion Models

Although diffusion-based methods have shown great performance for solving inverse problems for 2D images in different domains, there are few methods that are able to tackle inverse problems for 3D images because of the infeasible computational and data requirements. Training a network that learns the score function of a 3D volume directly would require a prohibitively large amount of memory and training data. Thus, previous works have focused on learning the distribution of 2D slices of a 3D volume and combining the slices in a coherent way at reconstruction time. Specifically, for 3D CT reconstruction, DiffusionMBIR [42] trains a diffusion model on the axial slices of volumes; at reconstruction time, it uses the total variation (TV) regularizer with a posterior sampling approach to encourage consistency between adjacent slices. Similarly, DDS [41] builds on this work by using accelerated methods of sampling and data consistency to greatly reduce the reconstruction time. However, although the TV regularizer has shown some success in maintaining smoothness across slices, it is not a data-driven method and does not properly learn the 3D prior. TPDM [109] addresses this problem by training a separate prior on the coronal slices of volumes with a conditional sampling approach, which serves as a data-driven method of maintaining slice consistency at reconstruction time, but requires that all the volumes have the same cubic shape. In exchange, this method sacrifices the speed gains made by DDS, requiring alternating updates between the two separate priors, and is also twice as computationally expensive at training time.

2.3 Distribution Mismatch Problems

Training diffusion models well requires vast amounts of clean training data [168, 73], which is infeasible to collect in many applications such as medical imaging [44, 167, 86], black hole imaging [61, 60], and phase retrieval [115, 188]. In particular, for very challenging inverse problems such as black hole imaging [61] and Fresnel phase retrieval [71], no ground truth images are known and one only has a single measurement \mathbf{y} available. In other applications such as dynamic CT reconstruction [147] and single photo emission CT [118], 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 Chapter 7 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.

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) [112] or wavelet transform [52] regularizers that encourage image sparsity. More recently, plug-and-play (PnP) methods have risen in popularity [170, 77, 75, 76]; these methods use a denoiser to solve general inverse problems. Although these methods often use a trained denoiser, [154] found that using an off-the-shelf denoiser such as block matching 3D [49] 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 using the loss function

$$L(\theta) = \|\mathbf{y} - \mathcal{A}(f_\theta(\mathbf{z}))\|_2^2, \quad \mathbf{z} \sim \mathcal{N}(0, \mathbf{I}), \quad (2.5)$$

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 [177]. Various methods have been proposed involving early stopping, regularization, and network initialization [122, 90, 11]. Nevertheless, the method is very sensitive to parameter selection and implementation and can take a long time to train [90].

Most diffusion inverse solving (DIS) methods learn a prior from a large collection of clean in-distribution training images, but recently [10] and [47] 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

$$L(\theta) = \|\mathbf{y} - \mathcal{A}(\text{CG}(\hat{\mathbf{x}}_{0:t}(\mathbf{x}_t; \theta)))\|_2^2 \quad (2.6)$$

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

$$\operatorname{argmin}_{\mathbf{x}} \frac{\gamma}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{x})\|_2^2 + \frac{1}{2} \|\mathbf{x} - \hat{\mathbf{x}}_{0:t}\|_2^2, \quad (2.7)$$

where γ is a tradeoff parameter controlling the strength of the prior versus the measurement. Crucially, these methods introduce an additional LoRA module [78] 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 [47] involving noisy initializations and early stopping to obtain good results and avoid artifacts. Our patch-based model of Chapter 4 avoids this overfitting issue.

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, [138], [207], and [209] 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.

CHAPTER 3

Accelerated Wirtinger Flow with Score-based Image Priors for Holographic Phase Retrieval in Poisson-Gaussian Noise Conditions

3.1 Motivation

Existing algorithms for Poisson-Gaussian holographic phase retrieval are limited in the literature. Handcrafted regularization-based algorithms such as total variation (TV) and higher order generalizations including total generalized variation (TGV) [22, 23] and the L1-norm of coefficients of wavelet transform [52] have been studied extensively. More recently, deep learning (DL)-integrated algorithms for solving inverse problems in computational imaging have been reported to be the state-of-the-art [141]. A trained network can serve as an object prior for regularizing the reconstructed image so it remains near a learned manifold [20]. Incorporating a trained denoising network as a regularizer $h(\cdot)$ led to methods such as plug-and-play (PnP) [34, 204, 92] and regularization by denoising (RED) [149]. In contrast to training a denoiser using clean images, there is growing interest in self-supervised image denoising approaches that do not require clean data as the training target [110, 15, 182].

In addition to training a denoiser as regularizer, generative model-based priors have also been proposed [7, 186]. Recently, diffusion models have gained significant traction for image generation [166, 73, 54, 168]. These probabilistic image generation models start with a clean image and gradually increase the level of noise added to the image, resulting in white Gaussian noise. Then in the reverse process, a neural network is trained to learn the noise in each step to generate or sample a clean image as in the original data distribution.

This chapter based on [115]. Equal contribution with Zongyu Li.

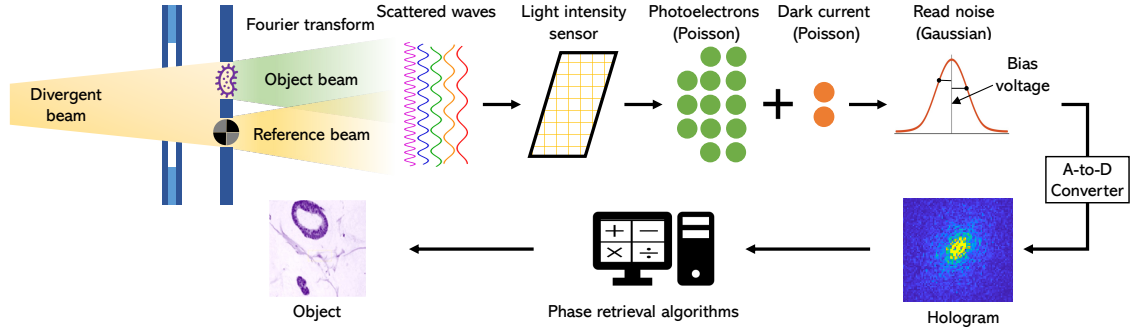


FIG 3.1 – Illustration of Poisson and Gaussian noise statistics in Fourier transform holographic phase retrieval.

The score-based diffusion models estimate the gradients of data distribution and can be used as plug-and-play priors for inverse problems [66] such as image deblurring and MRI and CT reconstruction [108, 86, 46, 48, 127, 167].

However, the realm of using score-based models to perform phase retrieval is relatively unexplored; previous relevant works [158, 66] applied denoising diffusion probabilistic modeling (DDPM) to PR but with less realistic system models and under solely Gaussian or Poisson noise statistics.

In summary, our contributions are as follows:

- We present our proposed algorithm known as accelerated Wirtinger flow with a score-based image prior (*i.e.*, $\nabla h(x)$ in (2.1)) to address the challenge of holographic phase retrieval (PR) problem in the presence of Poisson and Gaussian (PG) noise statistics.
- Theoretically, we derive a Lipschitz constant for the holographic PR's PG log-likelihood and provide a guarantee that sequences generated by the proposed algorithm converge to critical points of the cost function.
- We perform a series of simulation experiments which demonstrate that: 1) Algorithms using the PG likelihood model yield superior reconstructions compared to those relying solely on either the Poisson or Gaussian likelihood models. 2) With the proposed score-based prior as regularization, the proposed approach generates higher quality reconstructions and is more robust to variations of noise levels (without any parameter tuning) than alternative state-of-the-art methods.

3.1.1 Poisson-Gaussian Log-likelihood for PR

Based on the physical model as demonstrated in Fig. 3.1, we model the system matrix \mathbf{A} by the (oversampled and scaled) discrete Fourier transform applied to a concatenation of the sample \mathbf{x} , a blank image (representing the holographic separation condition [107]) and a known reference image \mathbf{r} , so that \mathbf{y} follows the Poisson plus Gaussian distribution:

$$\begin{aligned} \mathbf{y} &\sim \mathcal{N}(\text{Poisson}(|\mathbf{A}(\mathbf{x})|^2 + \bar{\mathbf{b}}), \sigma^2 \mathbf{I}), \\ \mathbf{A}(\mathbf{x}) &\triangleq \alpha \mathcal{F}\{\mathbf{x}, \mathbf{0}, \mathbf{r}\}. \end{aligned} \quad (3.1)$$

Here σ^2 denotes the variance of Gaussian noise, and α denotes a scaling factor (quantum efficiency, conversion gain, etc.) after applying the Fourier transform. Plugging the negative log-likelihood of (3.1) into (2.1) leads to

$$g_{\text{PG}}(\mathbf{x}) = \sum_{i=1}^M g_i(\mathbf{x}), \quad g_i(\mathbf{x}) \triangleq -\log \left(\frac{\sum_{n=0}^{\infty} e^{-(|\mathbf{a}'_i \mathbf{x}|^2 + \bar{b}_i)} \cdot (|\mathbf{a}'_i \mathbf{x}|^2 + \bar{b}_i)^n \cdot \frac{e^{-\left(\frac{y_i - n}{\sqrt{2}\sigma}\right)^2}}{\sqrt{2\pi\sigma^2}}}{n!} \right).$$

Here \mathbf{a}'_i denotes the i th row of \mathbf{A} (since \mathbf{A} is linear). We opt to use WF for estimating \mathbf{x} because it is commonly used in practice due to its simplicity and efficiency [28]. The WF algorithm is based on the gradient of (3.2):

$$\begin{aligned} \nabla g_{\text{PG}}(\mathbf{x}) &= 2\mathbf{A}' \text{diag}\{\phi_i(|\mathbf{a}'_i \mathbf{x}|^2 + \bar{b}_i; y_i)\} \mathbf{A} \mathbf{x}, \\ \phi(u; v) &\triangleq 1 - \frac{s(u, v-1)}{s(u, v)}, \quad s(a, b) \triangleq \sum_{n=0}^{\infty} \frac{a^n}{n!} e^{-\left(\frac{b-n}{\sqrt{2}\sigma}\right)^2}. \end{aligned} \quad (3.2)$$

Lemma 1. The function $\phi(u)$ is Lipschitz differentiable and the Lipschitz constant for its derivative $\dot{\phi}(u)$ is:

$$\max\{|\ddot{\phi}(u)|\} \triangleq \mu = \left(1 - e^{-1/\sigma^2}\right) e^{\frac{2y_{\max}-1}{\sigma^2}}, \quad (3.3)$$

where $y_{\max} \triangleq \max_{i \in \{1, \dots, M\}} y_i$.

The proof of Lemma 1 is given in [38].

To facilitate choosing step sizes, the following Theorem provides a Lipschitz constant for the gradient of the log-likelihood (3.2).

Theorem 1 Assume $|x_j|$ is bounded above by C for each j , a Lipschitz constant of $\nabla g_{\text{PG}}(\mathbf{x})$ is

$$\begin{aligned} \mathcal{L}(\nabla g_{\text{PG}}) \triangleq & 4C^2 \|\mathbf{A}\|_2^2 \|\mathbf{A}\|_\infty^2 \left(1 - e^{-\frac{1}{\sigma^2}}\right) e^{\frac{2y_{\max}-1}{\sigma^2}} \\ & + 2\|\mathbf{A}\|_2^2 \left|1 - C^2 \|\mathbf{A}\|_\infty^2 \left(1 - e^{-\frac{1}{\sigma^2}}\right) e^{\frac{2y_{\max}-1}{\sigma^2}}\right|. \end{aligned} \quad (3.4)$$

where $y_{\max} \triangleq \|\mathbf{y}\|_\infty$.

Proof: Let $g_{\text{PG}}(\mathbf{x})$ denote a function that maps a vector $\mathbf{x} \in \mathbb{R}^N$ to a scalar; it is the sum of each $g_i(\mathbf{x}) \triangleq \phi_i(|\mathbf{a}'_i \mathbf{x}|^2 + b_i; y_i)$ over $i = 1, \dots, M$. Let $\mathbf{g}(\mathbf{x})$ denote a function that maps a vector $\mathbf{x} \in \mathbb{R}^N$ to the measurement space \mathbb{R}^M ; it is the concatenation of each $g_i(\mathbf{x})$. So $\nabla g_{\text{PG}}(\mathbf{x}) \in \mathbb{R}^N$, $\nabla^2 g_{\text{PG}}(\mathbf{x}) \in \mathbb{R}^{N \times N}$, and $\nabla \mathbf{g}(\mathbf{x}) \in \mathbb{R}^{M \times N}$.

By the chain rule, the Hessian of g_{PG} is

$$\nabla^2 g_{\text{PG}}(\mathbf{x}) = 2\mathbf{A}' (\text{diag}\{\mathbf{A}\mathbf{x}\} \nabla \mathbf{g}(\mathbf{x}) + \text{diag}\{\mathbf{g}(\mathbf{x})\} \mathbf{A}). \quad (3.5)$$

Assume $|x_j| \leq C$ for each j . Then it follows that $\|\text{diag}\{\mathbf{A}\mathbf{x}\}\|_2 \leq C\|\mathbf{A}\|_\infty$ by the construction of matrix-vector multiplication, leading to a Lipschitz constant for $\nabla g_{\text{PG}}(\mathbf{x})$:

$$\mathcal{L}(\nabla g_{\text{PG}}) = 2C\|\mathbf{A}\|_2 \|\mathbf{A}\|_\infty \|\nabla \mathbf{g}(\mathbf{x})\|_2 + 2\|\mathbf{A}\|_2^2 \|\text{diag}\{\mathbf{g}(\mathbf{x})\}\|_2. \quad (3.6)$$

Here $\mathcal{L}(\nabla g_{\text{PG}})$ denotes a Lipschitz constant for ∇g_{PG} , not necessarily the best one. To compute $\|\nabla \mathbf{g}(\mathbf{x})\|_2$, we substitute the Lipschitz constant of $\dot{\phi}(u)$ into (3.2) and apply Lemma 1, leading to

$$\|\nabla \mathbf{g}(\mathbf{x})\|_2 \leq 2C\|\mathbf{A}\|_2 \|\mathbf{A}\|_\infty \left(1 - e^{-\frac{1}{\sigma^2}}\right) e^{\frac{2y_{\max}-1}{\sigma^2}}. \quad (3.7)$$

To compute $\|\text{diag}\{\mathbf{g}(\mathbf{x})\}\|_2$, let

$$t \in [b, \max_i\{|\mathbf{a}'_i \mathbf{x}|^2\} + b] \subseteq \mathcal{T} \triangleq [b, C^2\|\mathbf{A}\|_\infty^2 + b]. \quad (3.8)$$

From the fact that $\dot{\phi}(t) \leq 1$ by its construction, one can show:

$$\begin{aligned} \|\text{diag}\{\mathbf{g}(\mathbf{x})\}\|_2 = \|\mathbf{g}(\mathbf{x})\|_\infty & \leq \max_{t \in \mathcal{T}} \{|\dot{\phi}(t)|\} \\ & \leq \left|1 - C^2\|\mathbf{A}\|_\infty^2 \max\{|\ddot{\phi}(t)|\}\right|. \end{aligned} \quad (3.9)$$

Combining (3.6), (3.7) and (3.9) completes the proof of Theorem 1. \blacksquare

Theorem 1 extends [38] by considering a nonlinear transformation ($\mathbf{A}\mathbf{x} \rightarrow |\mathbf{A}\mathbf{x}|^2$) and a different system matrix \mathbf{A} .

In the practical implementation we approximate (3.2), with a finite sum following [38]:

$$s(a, b) \approx \sum_{n=0}^{n^+} \frac{a^n}{n!} e^{-\left(\frac{b-n}{\sqrt{2}\sigma}\right)^2}, \quad n^+ = \lceil n^* + \delta\sigma \rceil, \quad (3.10)$$

with n^* given by

$$\begin{aligned} n^* &= \sigma \mathcal{W}\left(\frac{a}{\sigma^2} e^{b/\sigma^2}\right) \\ &\approx \sigma \left(\frac{b}{\sigma^2} \log\left(\frac{a}{\sigma^2}\right) - \log\left(\frac{b}{\sigma^2} \log\left(\frac{a}{\sigma^2}\right)\right) \right) \\ &= \frac{b}{\sigma} \log\left(\frac{a}{\sigma^2}\right) - \sigma \log\left(\frac{b}{\sigma^2} \log\left(\frac{a}{\sigma^2}\right)\right), \end{aligned} \quad (3.11)$$

where $\mathcal{W}(\cdot)$ denotes the Lambert function. The accuracy of this approximation is controlled by δ . Reference [38] provides a comprehensive analysis on the maximum error value of the truncated sum (3.10), finding that the infinite sum is well-approximated by taking a manageable number of terms. In particular, we found that for our problem specifications, taking 100 terms was sufficient to ensure an approximation with error less than 0.1 percent.

3.1.2 Accelerated Wirtinger Flow with Score-based Image Prior

For accelerating the WF algorithm, we followed the implementation of [111] as its convergence guarantee was proved. Assuming that the true score function can be learned properly, when we have a trained score function $s_\theta(\mathbf{x}, \sigma)$ by applying (2.4), the gradient descent algorithm for MAP estimation (2.1) has the form: $\mathbf{x}_{t+1} = \mathbf{x}_t - \mu(\nabla g(\mathbf{x}_t) + s_\theta(\mathbf{x}_t, \sigma_k))$. Algorithm 1 summarizes our proposed AWFS algorithm. (Supplement shows the vanilla version without acceleration.) Similar to Langevin dynamics, we choose σ_k to be a descending sequence of noise levels. In practice, we generally use each noise level a fixed number of times, with geometrically spaced noise levels between some lower and upper bound. The step size factor β in Algorithm 1 can be selected empirically, but we show that the Lipschitz constant of the gradient $\nabla g_{\text{PG}}(\mathbf{x}_t) + s_\theta(\mathbf{x}_t, \sigma_k)$ exists as demonstrated in Theorem 2 (see proof in Supplement); hence with sufficiently small step size β , the sequence generated by Algorithm 1 will converge to a critical point of the posterior distribution in (2.1).

We assume that the data allows the neural network to learn the score function well, *i.e.*, $s_\theta(\mathbf{x}, \sigma) \approx \nabla \log(p_\sigma(\mathbf{x}))$, where $p_\sigma(\mathbf{x}) = p(\mathbf{x}) \otimes \mathcal{N}(0, \sigma^2)$, where \otimes denotes (circular) convolution. Supplement shows that $\nabla \log(p_\sigma(\mathbf{x}))$ is Lipschitz continuous on $[-C, C]^N$.

Algorithm 1 Proposed AWFS method for PR: accelerated WF with score-based image prior.

Require: Measurement \mathbf{y} , system matrix \mathbf{A} , momentum factor $\eta_0 = 1$, step size factor β , truncation operator $\mathcal{P}_C(\cdot) \rightarrow [0, C]$; initial image \mathbf{x}_0 , initial auxiliary variables $\mathbf{z}_0 = \mathbf{w}_0 = \mathbf{v}_0 = \mathbf{x}_0$, initialize $\sigma_1 > \sigma_2 > \dots > \sigma_K$.

for $k = 1 : K$ **do**

for $t = 1 : T$ **do**

 Set step size $\mu = \beta\sigma_k^2$.

 Set $\Delta z_{t,k} = \frac{\eta_{t-1,k}}{\eta_{t,k}}(\mathbf{z}_{t,k} - \mathbf{x}_{t,k})$.

 Set $\Delta x_{t,k} = \frac{\eta_{t-1,k}-1}{\eta_{t,k}}(\mathbf{x}_{t,k} - \mathbf{x}_{t-1,k})$.

 Set $\mathbf{w}_{t,k} = \mathcal{P}_C(\mathbf{x}_{t,k} + \Delta z_{t,k} + \Delta x_{t,k})$.

 Compute $s_\theta(\mathbf{x}_{t,k}, \sigma_k)$ and $s_\theta(\mathbf{w}_{t,k}, \sigma_k)$.

 Set $\mathbf{z}_{t+1,k} = \mathbf{w}_{t,k} - \mu(\nabla g_{\text{PG}}(\mathbf{w}_{t,k}) + s_\theta(\mathbf{w}_{t,k}, \sigma_k))$.

 Set $\mathbf{v}_{t+1,k} = \mathbf{x}_{t,k} - \mu(\nabla g_{\text{PG}}(\mathbf{x}_{t,k}) + s_\theta(\mathbf{x}_{t,k}, \sigma_k))$.

 Set $\eta_{t+1,k} = \frac{1}{2} \left(1 + \sqrt{1 + 4\eta_{t,k}^2} \right)$.

 Choose weight factor $\gamma_{t,k}$ (see Theorem 2).

 Set $\mathbf{x}_{t+1,k} = \mathcal{P}_C(\gamma_{t,k}\mathbf{z}_{t+1,k} + (1 - \gamma_{t,k})\mathbf{v}_{t+1,k})$.

end for

end for

Return $\mathbf{x}_{T,K}$.

Using $p_\sigma(\mathbf{x})$, we define the smoothed posterior as

$$p_\sigma(\mathbf{x}|\mathbf{A}, \mathbf{y}, \bar{\mathbf{b}}, \mathbf{r}) \propto p(\mathbf{y}|\mathbf{A}, \mathbf{x}, \bar{\mathbf{b}}, \mathbf{r})p_\sigma(\mathbf{x}). \quad (3.12)$$

Theorem 2 For a smooth density function $p_{\sigma_k}(\mathbf{x})$ that has finite expectation with $\sigma_k > 0$, the Lipschitz constant of $\nabla g_{\text{PG}}(\mathbf{x}_{t,k}) + s_\theta(\mathbf{x}_{t,k}, \sigma_k)$ exists when each element in $\mathbf{x}_{t,k}$ satisfies $0 \leq |x_j| \leq C$ for each j . Furthermore, if the weighting factor $\gamma_{t,k}$ is set to 0 if $p_{\sigma_k}(\mathbf{v}_{t+1}|\mathbf{y}, \mathbf{A}, \bar{\mathbf{b}}) \leq p_{\sigma_k}(\mathbf{z}_{t+1}|\mathbf{y}, \mathbf{A}, \bar{\mathbf{b}})$ and 1 otherwise, then with sufficiently small β , the inner iteration sequence $\{\mathbf{x}_{t,k}\}_{t=1}^T$ generated by Algorithm 1 is bounded, and any accumulation point of that sequence as $t \rightarrow \infty$ is a critical point of the smoothed posterior distribution $p_{\sigma_k}(\mathbf{x}|\mathbf{y}, \mathbf{A}, \bar{\mathbf{b}}, \mathbf{r})$ in (3.12).

Sketch of Proof: By Lipschitz continuity of $\log(p_\sigma(\mathbf{x}))$, and from the design of Algorithm 1, $\mathbf{x}_{t,k}$ and $\mathbf{w}_{t,k}$ are both bounded between $[-C, C]$ for all t, k , so the Lipschitz constant \mathcal{L}^* of $\nabla g_{\text{PG}}(\cdot) + s_\theta(\cdot)$ exists. With the stepsize μ satisfying $0 < \mu < \frac{1}{\mathcal{L}^*}$, and the weighting factor $\gamma_{t,k} \in \{0, 1\}$ being chosen according to the higher posterior probability between $p_{\sigma_k}(\mathbf{z}|\mathbf{A}, \mathbf{y}, \bar{\mathbf{b}}, \mathbf{r})$ and $p_{\sigma_k}(\mathbf{v}|\mathbf{A}, \mathbf{y}, \bar{\mathbf{b}}, \mathbf{r})$ (see [111]), we satisfy all conditions in

Theorem 1 of [111], which establishes the critical-point convergence of the sequence $\mathbf{x}_{t,k}$ generated by Algorithm 1 for any $\sigma_k, k = 1, \dots, K$. Hence the sequence $\mathbf{x}_{t,k}$ generated by Algorithm 1 converges as $t \rightarrow \infty$ to a critical-point of the posterior $p_{\sigma_k}(\mathbf{x}|\mathbf{A}, \mathbf{y}, \bar{\mathbf{b}}, \mathbf{r})$ for any σ_k . Supplement shows the full proof of Theorem 2.

Theorem 2 states that at each iteration, we should choose $\gamma_{t,k} \in \{0, 1\}$ according to the higher posterior probability to ensure convergence. However, we found empirically that setting $\gamma_{t,k} = 0.5$ also led to convergent sequences (see Fig. 3.9), so we use this practical version for all the results in Section 3.2. An alternative to our score-based image prior approach would be to alternate data fidelity updates with the sampling step of a diffusion model based on the denoising diffusion probabilistic models (DDPM) framework [158]. Here, we implement and present the results of a slightly altered version (details shown in the Supplement). There are two main changes for this version: first, we initialize with the result of a baseline reconstruction method (in this case, the unregularized Poisson method), rather than pure noise; second, we use the PG likelihood to enforce data fidelity as opposed to the Gaussian likelihood. The former change enables the latter change to be feasible, as it reduces the required number of sampling steps by a factor of 30, allowing for the expensive computation of the PG likelihood each step.

3.2 Experiment

3.2.1 Experiment Settings

Data. We tested all algorithms on three datasets: 162 histopathology images related to breast cancer [2] (train/val/test is 122/20/20); 920 images from CelebA dataset [125] (train/val/test is 800/100/20); and 720 images from a homemade CT-density dataset (train/val/test is 600/100/20). The CT-density dataset was generated from single-photon emission computerized tomography (SPECT)/CT images for Yttrium-90 radionuclide therapy after applying the CT-to-density calibration curve [114]. Although the size of training datasets are relatively small compared to typical datasets such as ImageNet or LSUN [73, 168] that have millions of images, we do not require the score functions to learn image priors strong enough to generate realistic images from white Gaussian noise; rather, it suffices for the priors to be able to denoise moderately noisy images.

System Model. The system matrix is based on the discrete Fourier transform of the concatenation of the true image \mathbf{x} , a blank image $\mathbf{0}$ and a reference image \mathbf{r} with scaling and oversampling. We set the scaling factor α to be in the range $[0.02, 0.035]$ so that the

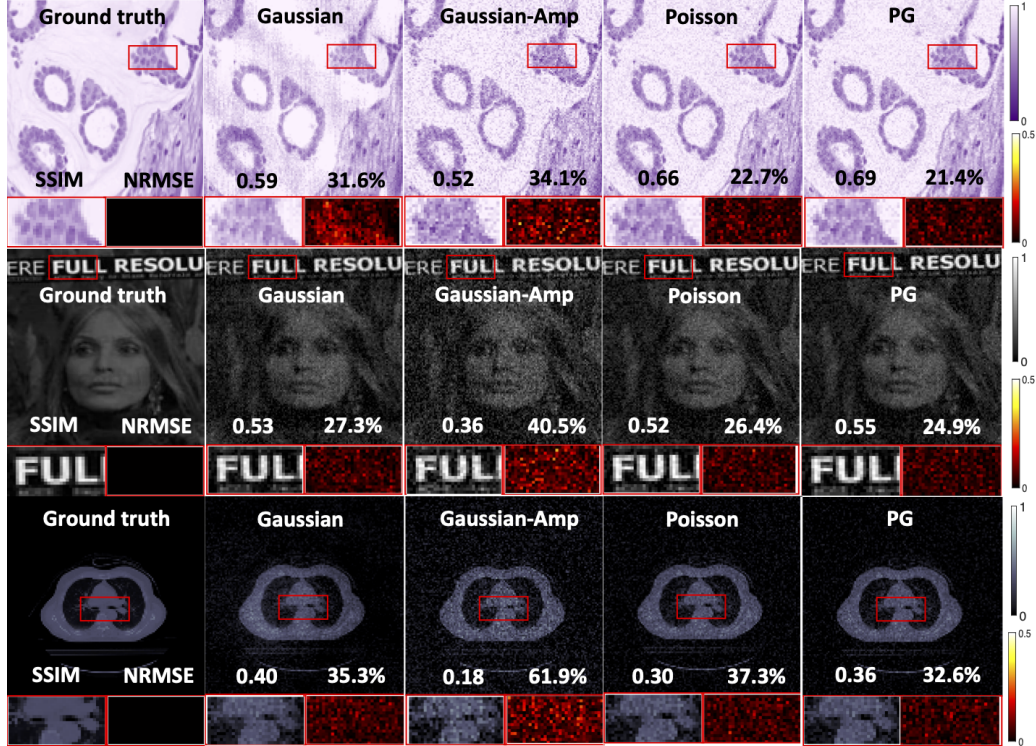


FIG 3.2 – Reconstructed images by unregularized methods (Gaussian, Gaussian-Amplitude, Poisson and Poisson-Gaussian) on Histopathology dataset [2], celebA dataset [125] and CT-density dataset. The bottom left/right subfigures correspond to the zoomed in area and the error map for each image. We used $\alpha = 0.035$ and $\sigma = 1$.

average counts per pixel range from 6 to 25; the oversampled ratio is set to 2. We set \mathbf{r} to be a binary random image similar to what was used in [107]. The standard deviation of the Gaussian read noise added to the measurements \mathbf{y} was set as $\sigma \in [0.5, 1.5]$.

Algorithms. For unregularized algorithms, we implemented Gaussian WF [28], Poisson WF [12] and Poisson-Gaussian WF. For regularized algorithms, we implemented smoothed total variation (TV) based on the Huber function [84, p. 184] and PnP/RED methods with the DnCNN denoiser [205]: PnP-ADMM [178], PnP-PGM [91], and RED-SD [149]. We also implemented the RED-SD algorithm with “Noise2Self” zero-shot image denoising network [15] (RED-SD-SELF). For diffusion models, we implemented DOLPH [158] and our proposed AWFS. Supplement shows the implementation details of each algorithm. We used spectral initialization [131] for the Gaussian PR and Poisson PR methods; we then used the output results from Poisson PR to initialize other algorithms. We ran all algorithms until the normalized root mean squared error (NRMSE) between consecutive iterations differed by less than 0.01 percent or reached the maximum number of iterations (e.g., 50).

To evaluate the robustness and limitation of these algorithms, we first tuned the parameters for each algorithm at the noise level when $\alpha = 0.030$ and $\sigma = 1$, and then held them fixed throughout all experiments (Table 3.1, Table 3.2, Fig. 3.7 and Fig. 3.8). In practice the ground truths are unknown, so oracle tuning of test datasets is infeasible (though some form of cross validation may be possible). Though the numbers reported could fluctuate after careful refinement, *e.g.*, by performing grid search on tuning parameters, such techniques would potentially impede the algorithm’s practical use.

Network Training. For PnP denoising networks, we trained all denoisers on different noise levels $\sigma \in \{9, 11, 13, 15\}$ and found that $\sigma = 15$ worked the best on our data. We also used the denoiser scaling technique from [196] to dynamically adjust the performance of all PnP methods. To perform score matching, we applied 20 geometrically spaced noise levels between 0.005 and 0.1 on each of the training images. All networks were implemented in PyTorch and trained on an NVIDIA Quadro RTX 5000 GPU using the ADAM optimizer [99] for 1000 epochs with the best one being selected based on the mean squared error (MSE) validation loss.

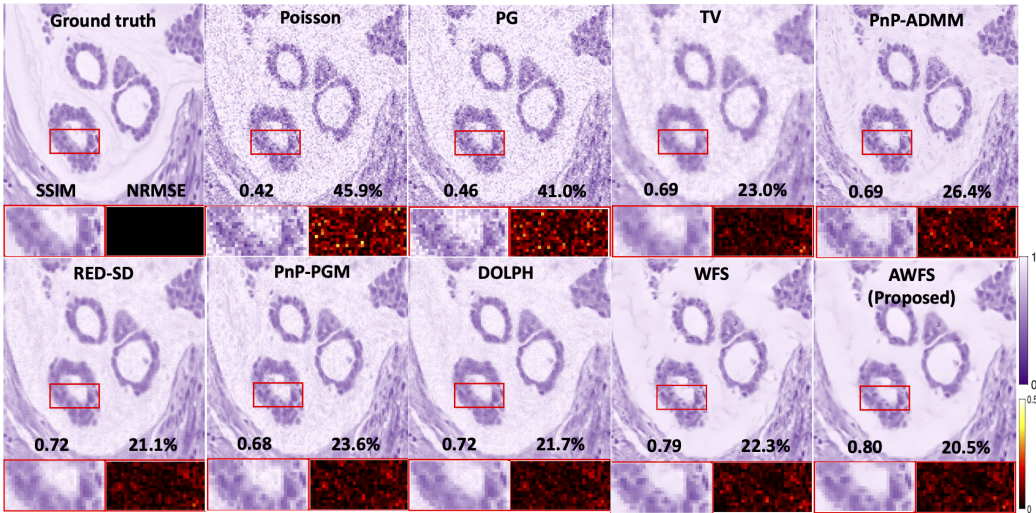


FIG 3.3 – Reconstructed images on dataset [2]. The bottom left/right subfigures correspond to the zoomed in area and the error map for each image. We used $\alpha = 0.02$ and $\sigma = 1$.

3.2.2 Results

We compared all implemented algorithms both qualitatively, by visualizing the reconstructed images and residual errors, and quantitatively, by computing the NRMSE and

structural similarity index measure (SSIM) [184]. Due to the global phase ambiguity, *i.e.*, all the algorithms can recover the signal only to within a constant phase shift due to the loss of global phase information, we corrected the phase of \hat{x} by $\hat{x}_{\text{corrected}} \triangleq \text{sign}(\langle \hat{x}, x_{\text{true}} \rangle) \hat{x}$.

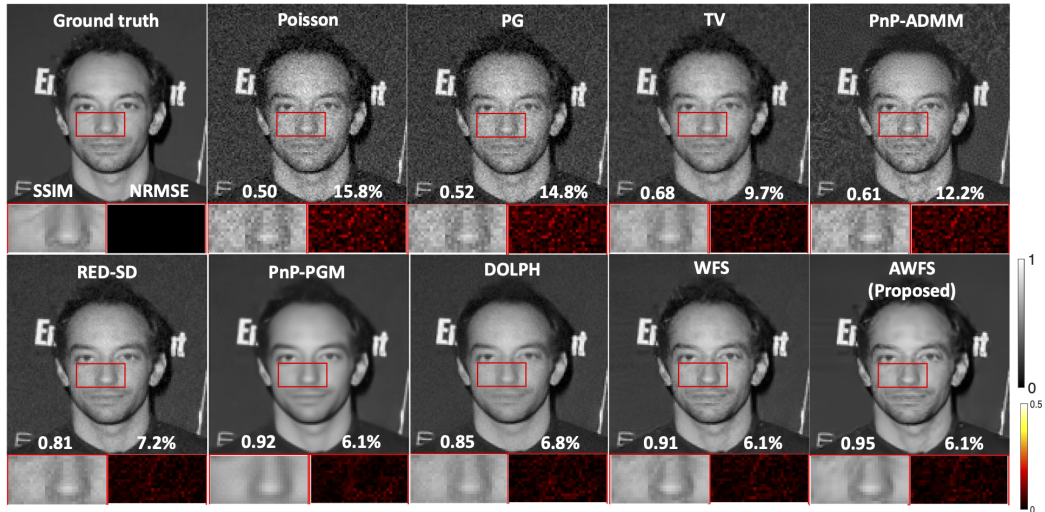


FIG 3.4 – Reconstructed images on celebA dataset [125]. The bottom left/right subfigures correspond to the zoomed in area and the error map for each image. We used $\alpha = 0.035$ and $\sigma = 1$.

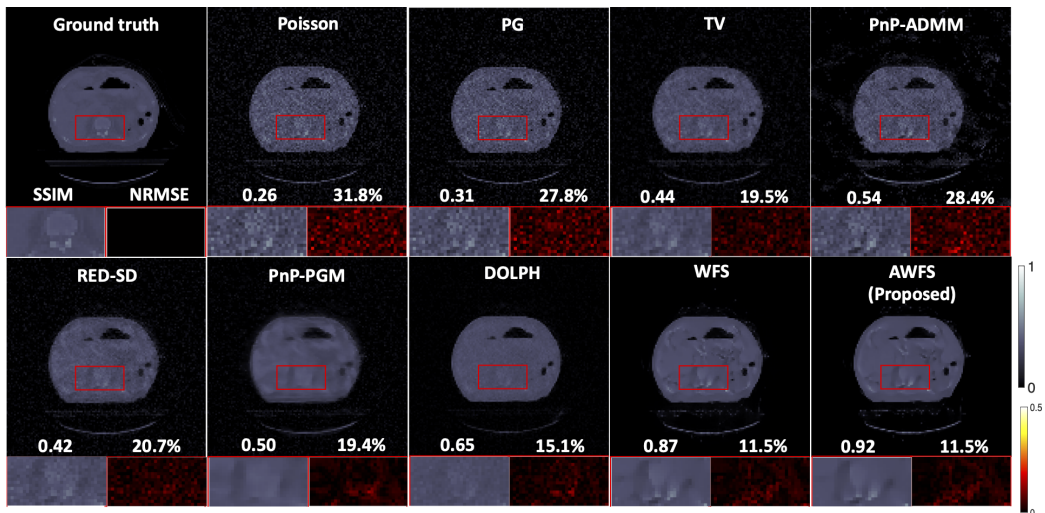


FIG 3.5 – Reconstructed images on CT-density dataset. The bottom left/right subfigures correspond to the zoomed in area and the error map for each image. We used $\alpha = 0.035$ and $\sigma = 1$.

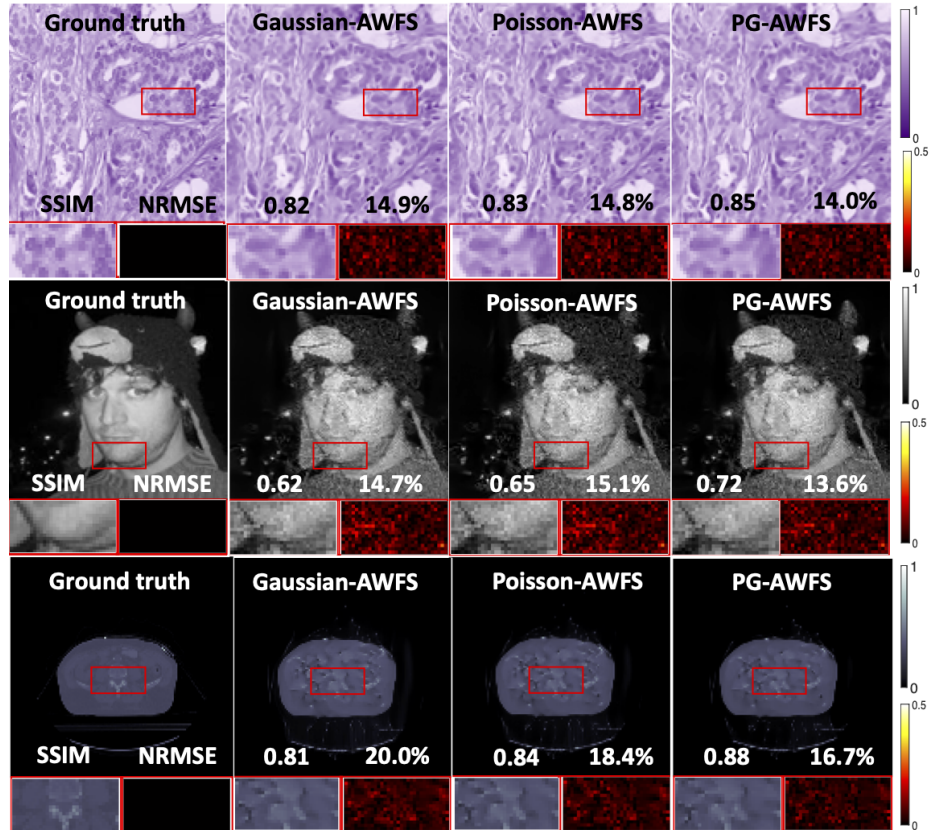


FIG 3.6 – Reconstructed images by Gaussian, Poisson and Poisson-Gaussian log-likelihood model with AWFS image prior. Tested on Histopathology dataset [2], celebA dataset [125] and CT-density dataset. The bottom left/right subfigures correspond to the zoomed in area and the error map for each image. α and σ were set to 0.025 and 1, respectively.

Fig. 3.2 shows experiments of running unregularized methods based on different noise models on the histopathology, CelebA, and CT density datasets. For comparison, we ran the unregularized methods with a Gaussian only noise model, Poisson only, and PG noise model.

Fig. 3.3, Fig. 3.4 and Fig. 3.5 visualize reconstructed images generated by algorithms mentioned in the previous section. The WF with PG likelihood outperforms WF with Poisson likelihood with a consistently higher SSIM and lower NRMSE. Of the regularized algorithms with PG likelihood, our proposed AWFS had less visual noise and achieved greater detail recovery compared to other methods, as evidenced by the zoomed-in area in these figures. Fig. 3.6 shows that for a variety of datasets, when combined with the AWFS method, while the Poisson only and Gaussian only models lead to reasonable reconstructions, the PG noise model leads to the highest quality image. For all three datasets

TBL 3.1 – SSIM and NRMSE for Poisson and Poisson-Gaussian likelihoods. Results were averaged across 7 different noise levels by varying $\alpha \in 0.02 : 0.005 : 0.035$ in (3.1), with $\sigma = 1$.

Likelihood	Unregularized (SSIM/NRMSE)	DOLPH (SSIM/NRMSE)	AWFS (SSIM/NRMSE)	RED-SD	PnP-PGM	TV
DataSet: Histopathology [2]						
Gaussian [28]	0.52 ± 0.18	41.2 ± 25.3	0.76 ± 0.07	18.0 ± 3.0	0.84 ± 0.06	16.2 ± 3.7
Poisson [12]	0.54 ± 0.18	31.7 ± 10.2	0.72 ± 0.13	19.5 ± 6.1	0.83 ± 0.06	16.2 ± 3.7
Poisson-Gaussian	0.57 ± 0.18	28.9 ± 9.0	0.80 ± 0.06	16.0 ± 2.9	0.85 ± 0.05	15.4 ± 3.7
DataSet: CelebA [125]						
Gaussian [28]	0.31 ± 0.09	55.6 ± 13.9	0.70 ± 0.12	14.5 ± 17.4	0.72 ± 0.16	15.3 ± 11.8
Poisson [12]	0.39 ± 0.10	24.5 ± 11.4	0.61 ± 0.12	15.6 ± 10.6	0.72 ± 0.16	15.2 ± 11.8
Poisson-Gaussian	0.42 ± 0.10	21.8 ± 9.1	0.71 ± 0.11	13.7 ± 11.1	0.74 ± 0.15	14.8 ± 11.9
DataSet: CT-Density						
Gaussian [28]	0.29 ± 0.09	50.5 ± 8.0	0.51 ± 0.12	22.4 ± 3.9	0.82 ± 0.11	19.1 ± 4.8
Poisson [12]	0.19 ± 0.06	48.9 ± 13.1	0.38 ± 0.11	25.6 ± 7.5	0.84 ± 0.08	17.8 ± 4.3
Poisson-Gaussian	0.24 ± 0.06	40.8 ± 9.5	0.55 ± 0.08	20.0 ± 3.3	0.88 ± 0.05	16.4 ± 3.7

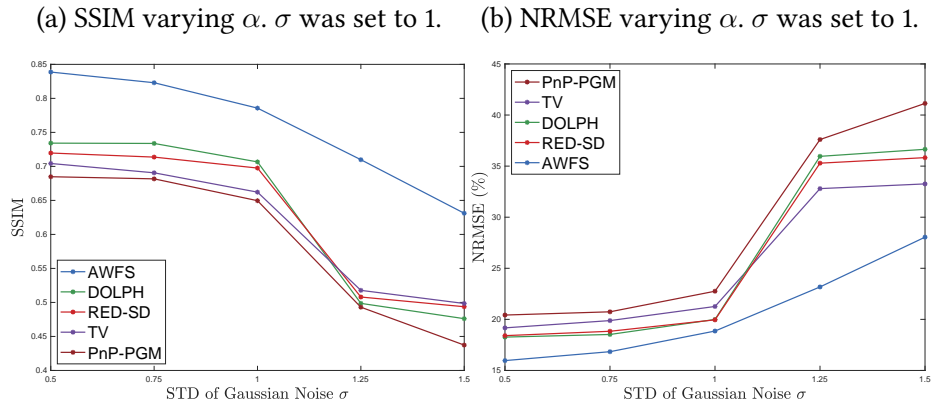
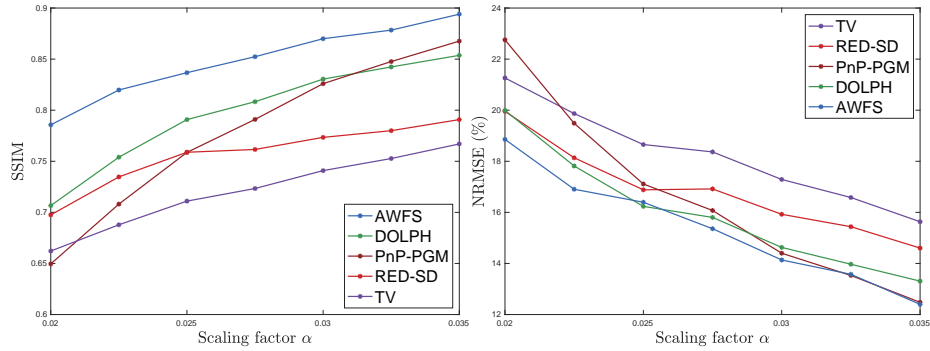


FIG 3.7 – Comparison of SSIM and NRMSE varying scaling factor $\alpha \in [0.02, 0.035]$ and STD of Gaussian noise $\sigma \in [0.25, 1.5]$ defined in (3.1).

TBL 3.2 – SSIM and NRMSE using Poisson Gaussian likelihood with different regularization/image prior approaches. Results were averaged across 7 different noise levels by varying $\alpha \in 0.02 : 0.005 : 0.035$ in (3.1). WFS* runs the same number of iterations as AWFS whereas WFS[†] runs more iterations until convergence.

Dataset	Histopathology [2]		CelebA [125]		CT-Density	
Methods	SSIM	NRMSE (%)	SSIM	NRMSE (%)	SSIM	NRMSE (%)
Unregularized	0.57 ± 0.18	28.9 ± 9.0	0.42 ± 0.10	21.8 ± 9.1	0.24 ± 0.06	40.8 ± 9.5
RED-SD-SELF [15]	0.66 ± 0.13	21.9 ± 4.5	0.60 ± 0.09	15.9 ± 10.6	0.34 ± 0.04	28.1 ± 4.1
PnP-ADMM [178]	0.71 ± 0.11	20.7 ± 4.2	0.56 ± 0.08	16.7 ± 8.1	0.55 ± 0.03	31.2 ± 2.7
TV regularizer	0.72 ± 0.11	18.2 ± 3.9	0.64 ± 0.07	14.4 ± 8.6	0.41 ± 0.03	23.7 ± 2.8
RED-SD [149]	0.76 ± 0.09	16.8 ± 3.6	0.69 ± 0.11	13.9 ± 10.9	0.38 ± 0.04	25.9 ± 4.0
PnP-PGM [91]	0.78 ± 0.11	16.5 ± 4.5	0.74 ± 0.14	13.5 ± 11.3	0.42 ± 0.07	24.6 ± 4.4
DOLPH [158]	0.80 ± 0.06	16.0 ± 2.9	0.71 ± 0.11	13.7 ± 11.1	0.55 ± 0.08	20.0 ± 3.3
WFS*	0.76 ± 0.12	18.2 ± 5.5	0.63 ± 0.16	16.9 ± 11.8	0.53 ± 0.17	21.3 ± 7.6
WFS [†]	0.83 ± 0.06	16.2 ± 4.0	0.70 ± 0.16	15.7 ± 11.8	0.74 ± 0.13	17.3 ± 4.8
AWFS (Proposed)	0.85 ± 0.05	15.4 ± 3.7	0.74 ± 0.15	14.8 ± 11.9	0.88 ± 0.05	16.4 ± 3.7

shown, when used in conjunction with our AWFS method, including both Poisson and Gaussian likelihoods results in the highest quality reconstruction both in terms of quantitative metrics as well as visually. Thus, although the score function provides a useful prior for recovering an image when the measurement is very noisy, a proper noise model is also crucial to a high quality reconstruction.

For quantitative evaluations, Table 3.1 illustrates the effect of using our proposed PG likelihood as compared to the simpler Poisson likelihoods. We did not run the Gaussian likelihood with DOLPH or AWFS due to the abysmal performance with this likelihood. In all cases, usage of the PG likelihood results in improved image quality in terms of both metrics. Table 3.2 consists of experiments using the PG likelihood and shows the efficacy of the proposed AWFS method over other methods. In particular, our AWFS had superior quantitative performance over all other compared methods on the histopathology and CT-density datasets; in contrast, the PnP-Prox showed the lowest NRMSE on celebA dataset. This is likely due to higher randomness in celebrity faces because the effectiveness of generative models can vary depending on the dataset used. Thus, when provided with a small amount of training data with high randomness, image denoising models (DnCNN) may be more effective than generative models.

We also tested the robustness of the leading algorithms in Table 3.2, by varying both scaling factor α and STD of Gaussian noise σ . Fig. 3.7 and Fig. 3.8 illustrate results, where

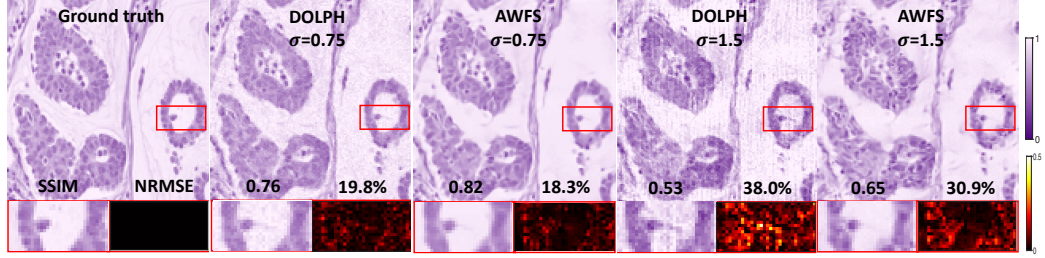


FIG 3.8 – Reconstructed images by DOLPH [158] and our proposed AWFS method under different σ values, for $\alpha = 0.02$.

our AWFS algorithm had the highest SSIM and lowest NRMSE. In Fig. 3.8, AWFS demonstrated minimal variations in SSIM and NRMSE metrics than DOLPH as evidenced by the smaller discrepancies in SSIM (0.17 vs. 0.23) and NRMSE (12.6% vs. 18.2%) when σ varies from 0.75 to 1.5. Fig. 3.9 compares the convergence rate of AWF vs. WF for the Poisson and PG likelihood, respectively. Under a variety of noise levels, AWF consistently converged faster than WF in terms of number of iterations.

It is a known property of diffusion models that they can produce images with hallucinated features if the measurements are insufficiently informative. In the case of low-count phase retrieval with serious corruptions of both Poisson and Gaussian noise, as is investigated here, the measurement is highly corrupted and contains magnitude-only measurements of the original signal. Thus, it may be difficult for the diffusion models to avoid some otherwise realistic hallucinations if the data consistency is not strong enough to guide the model away from such hallucinations. On the other hand, if the measurements are less corrupted, then the data consistency should be strong enough to avoid such hallucinations. Fig. 3.10 provides examples of this for the CT image dataset via a comparison of the reconstruction quality of the AWFS method over a range of count levels. With the lowest scaling factor, *e.g.*, $\alpha = 0.02$, the measurements were seriously corrupted with noise, and the method may hallucinate some features. However, at higher count level, *e.g.*, $\alpha = 0.05$, there is enough information in the measurement to enforce consistency and avoid noticeable hallucinations. We performed the same experiment twice with different noisy initializations and all other parameters held equal to demonstrate robustness of the method under different initializations.

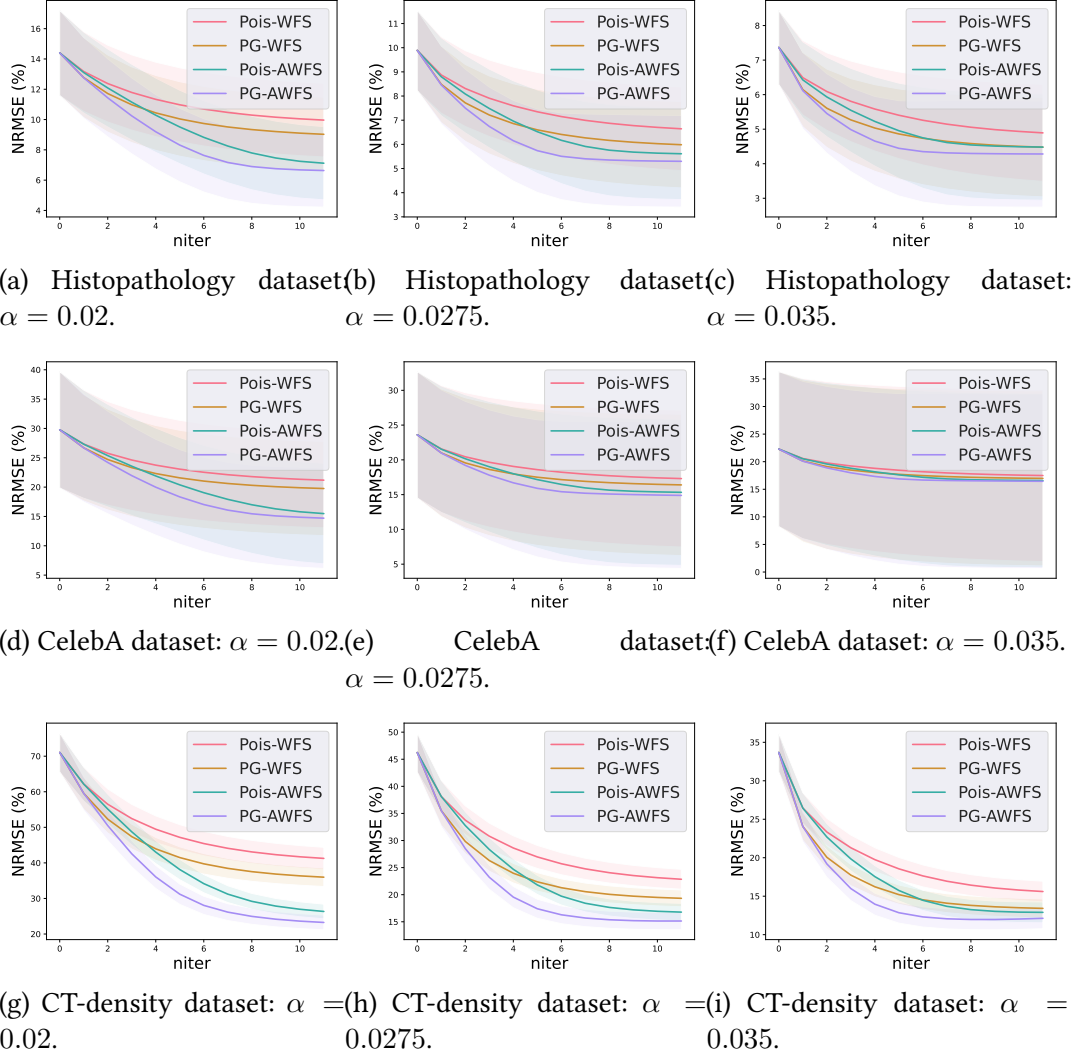


FIG 3.9 – Comparing AWF vs. WF with NRMSE vs. number of iterations under different noise levels. The curves and shadows represent the mean and standard deviation, respectively.

3.3 Discussion

PR has a long-standing history in the field of signal processing and imaging. Pioneering works such as the error reduction (ER) and hybrid input-output (HIO) algorithms by Gerchberg Saxton [64] and Fienup [62] have been proposed to address this problem. These iterative algorithms involve constraints imposed on evaluations between the image domain and frequency domain. However, these methods have limitations in terms of the quality of reconstructed images and their convergence remains uncertain [201]. Another approach to solving PR problems is through compressed sensing and optimization techniques like

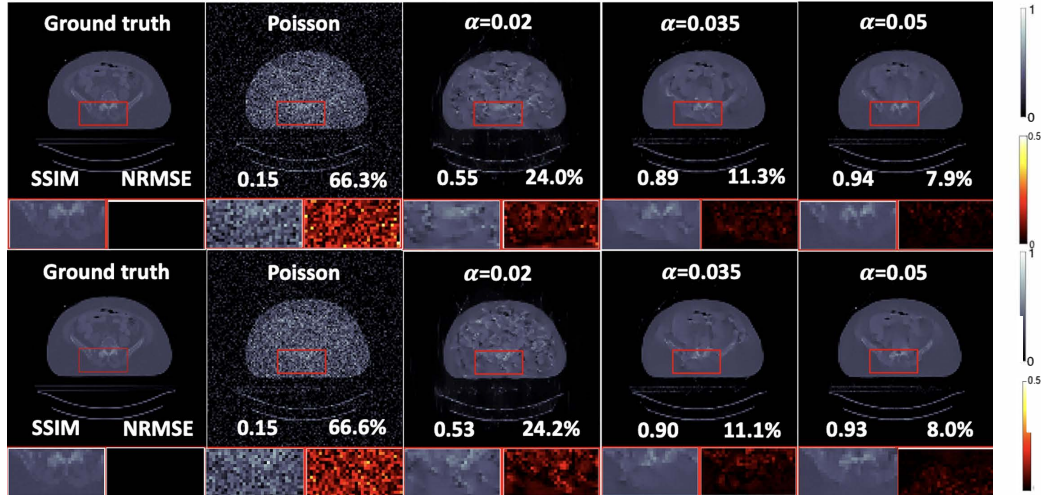


FIG 3.10 – Reconstructed images by the unregularized Poisson method (the second column) as well as with the AWFS method for different scaling factors α (third to fifth columns). The top and bottom rows show reconstructions from different measurement realizations.

Wirtinger flow (WF) [28], matrix lifting [30, 27, 157], MM [146] and ADMM [119]. This chapter focuses on the WF algorithm due to being straightforward to incorporate with the DL regularizer for the image prior. The likelihood modelling of the noise statistics existing in the measurement is also critical. Previous studies have primarily focused on modelling either Gaussian or Poisson likelihood only, but in practical scenarios, both types of noise are often encountered. Therefore, this chapter contributes to a more practical perspective of addressing the holographic PR problem by using a PG likelihood and incorporating state-of-the-art deep learning image priors. In the case where the measurement is contaminated with Poisson and Gaussian noise, the speedup in reconstruction is crucial, as the bottleneck of our algorithm is in computing the PG likelihood. Additionally, though it is viable to perform a large number of neural network evaluations to perform image reconstruction, it is unrealistic to compute a similarly large number of PG likelihoods. Thus, we perform acceleration in WF algorithm following [111], which guarantees convergence to a critical point for the Holographic PR problems.

In our evaluation of three datasets, we consistently observed that the use of PG likelihood yielded superior performance compared to using either Poisson or Gaussian likelihood alone, as expected. Additionally, the results obtained from the CT-density dataset were generally of lower quality than those from the other two datasets. This can be attributed to lower average counts per pixel (many zero pixels near the image borders). Using a DL image prior can be considered from two perspectives: training a denoiser

or training to learn the density distribution of images. In our work, we applied both approaches and observed that the effectiveness of these methods differed depending on the dataset tested. Specifically, in the Histopathology dataset [2] and the CT-density dataset, where the images share similar structures, the generative models performs better even when trained with limited data. In the case of the CelebA dataset [125], which includes a wide variety of celebrity faces, generative models did not exhibit as strong performance as denoiser methods when trained on limited data. This is likely due to the fact that generating high-quality images is generally more challenging than removing noise from existing images and may necessitate a larger training dataset. The effectiveness of accelerated WF compared to vanilla WF is due to the non-convexity of the PR problem. Although recent advances in geometric landscape analysis of PR can guarantee that all local minimizers are global even with random initialization [25], in practice the measurements are contaminated by noise so that many more measurements are required for the cost function to have a benign geometric landscape.

Despite the promising results achieved with our proposed AWFS approach, there are several limitations of our work. First, the approximate calculation of the infinite sum in (3.2) is accurate but computationally expensive. Future work should seek ways to accelerate this calculation while maintaining accuracy. Second, we did not implement and test the accelerated WF applied on the diffusion posterior sampling (DPS) method [39], for which the network is fine-tuned from a pretrained state-of-the-art diffusion model. This approach has the potential to advance current methods in PR problem and we will investigate it in the future. Another limitation of the proposed method is that it has been demonstrated on measurements that are based on simulations. To further demonstrate the efficacy of the method in a real-world setting, future work should consist of evaluating the accuracy of the methods when run on real measurement data. Finally, our experiments are limited to real-valued images, however, our method can be extended to handle complex-valued images by splitting real and imaginary components into separate reconstruction routines with different pretrained neural networks [211]. Addressing these limitations will be the future direction of this work.

3.4 Conclusion

We proposed a novel algorithm based on Accelerated Wirtinger Flow and Score-based image prior (AWFS) for Poisson-Gaussian holographic phase retrieval. Simulation experiments demonstrate that the proposed AWFS method produced the best reconstruction quality both qualitatively and quantitatively and was more robust to various noise levels,

compared to other state-of-the-art methods. Furthermore, we proved that our proposed algorithm has a critical-point convergence guarantee. Therefore, after the method is extended to accommodate complex-valued images, it should have much promise for real-world PR applications.

CHAPTER 4

Learning Image Priors through Patch-based Diffusion Models for Solving Inverse Problems

4.1 Motivation

Diffusion models require large quantities of training data and vast computational power to be able to generate high resolution images; Song et al. [168] and Ho et al. [73] used several days to weeks of training on over a million training images in the ImageNet [153] and LSUN [200] datasets to generate 256×256 images. This high cost motivates the research on improved training efficiency for diffusion models, such as fine-tuning an existing checkpoint on a different dataset [138, 42] to reduce the required training time and data. However, this strategy restricts the range of usable network architectures and requires the existence of a pretrained network, which limits the range of applications. Besides the demanding training data and computational cost, diffusion models also struggle in very large-scale problems, such as very high-resolution or 3D images.

Our proposed method tackles these challenges in a unified way by training diffusion models on patches of images, as opposed to whole images (see Fig. 4.1). We provide the location of the randomly extracted patch to the network to help it learn global image properties. Since each training image contains many patches, the required size of the training

This chapter based on [82]. The related code is at <https://github.com/jasonhu4/PaDIS>.

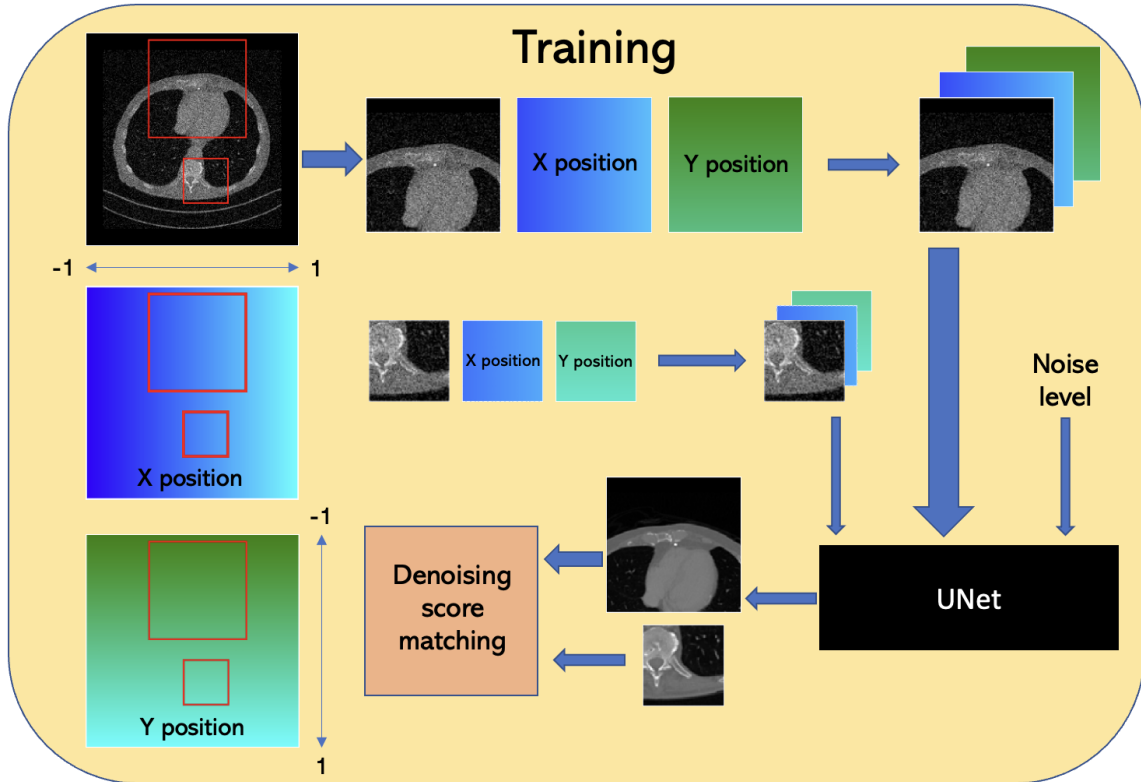


FIG 4.1 – Training the proposed Patch Diffusion Inverse Solver (PaDIS) method. Different sized patches are used in each training iteration. The X and Y position arrays have the same size as the patch and consist of the normalized X and Y coordinates of each pixel of the patch.

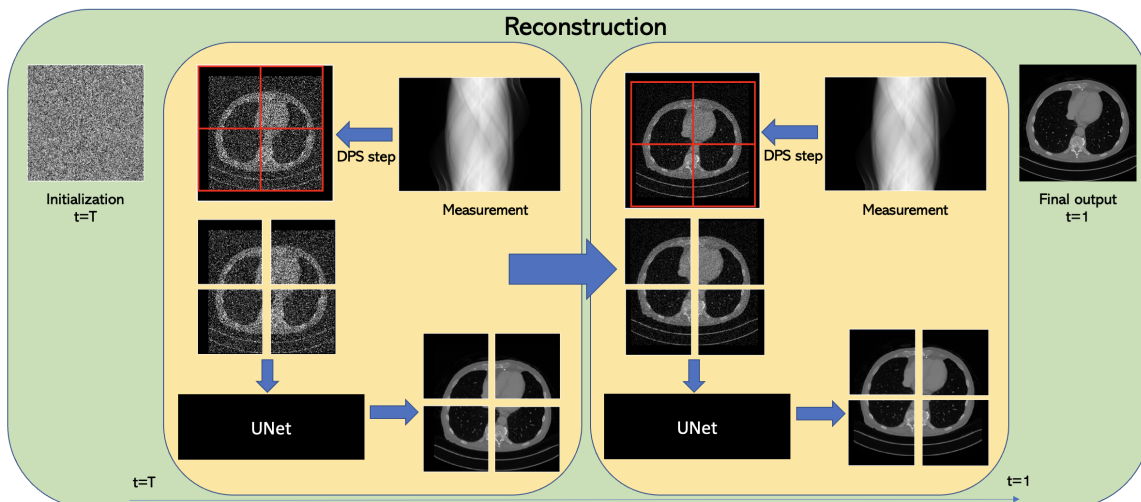


FIG 4.2 – Overview of reconstruction process for the proposed Patch Diffusion Inverse Solver (PaDIS) method. Starting at $t = T$, at each iteration we choose a random partition of the zero-padded image and use the neural network trained on patches to get the score function of the entire image. Due to the shifting patch locations, the output image has no boundary artifacts.

dataset is greatly reduced, from the millions usually needed to generate high-quality images to only a couple thousand or even several hundred. The required memory and training time is also reduced because it is never necessary to backpropagate the whole image through the network. Our proposed method allows for a flexible network architecture and only requires it to accept images of any size, a property true of many UNets [73], so there is much more flexibility in the architecture design than fine-tuning methods.

At inference time (see Fig. 4.2), by first zero padding the image, the proposed approach partitions it into patches in many different ways (see Fig. 4.3), eliminating the boundary artifacts between patches that would appear if non-overlapping patches were used. We develop a method to express the distribution of the whole image in terms of the patch distribution that is learned by the proposed patch-based network. By incorporating positional information of patches, this framework allows us to compute the score function of the whole image without ever inputting the whole image into the network. Unlike previous patch-based works that may be task-specific [190, 191, 142], the prior defined by our approach may be treated in a black box manner and then paired with any other stochastic differential equation (SDE) solver to sample from the prior, or with any general-purpose inverse problem solver to perform image reconstruction. We conduct experiments on multiple datasets and different inverse problems and demonstrate that the proposed method is able to synthesize the patches to produce reasonably realistic images and very accurate reconstructions for inverse problems. Furthermore, PaDIS provides a promising avenue for which generation and inverse problem solving of very large and high-dimensional images may be tackled in the future. (See Ch. 6.)

In summary, our main contributions are as follows:

- We provide a theoretical framework whereby a score function of a high-resolution high-dimensional image is learned purely through the score function of its patches.
- The proposed method greatly reduces the amount of memory and training data needed compared to traditional diffusion models.
- The trained network has great flexibility and can be used with many existing sampling algorithms and is the first patch-based model that can solve inverse problems in an unsupervised manner.
- We perform experiments on a variety of inverse problems to show superior image quality over whole image diffusion model methods while being far less resource heavy.

4.2 Methods

To be able to solve large 2D imaging problems as well as 3D and 4D inverse problems, our goal is to develop a model for $p(\mathbf{x})$ that does not require inputting the entire image \mathbf{x} into any neural network. We would like to simply partition \mathbf{x} into nonoverlapping patches, learn the distribution of each of the patches, and then piece them together to obtain the distribution of \mathbf{x} . However, this would result in boundary artifacts between the patches. Directly using overlapping patches would result in sections of the image covered by multiple patches to be updated multiple times, which is inconsistent with the theory of diffusion models. Ideally, we would like to use nonoverlapping patches to update \mathbf{x} , but with a variety of different patch tiling schemes so that boundaries between patches do not remain the same through the entire diffusion process.

To accomplish this task, we zero pad \mathbf{x}_0 and learn the distribution of the resulting zero padded image. More precisely, consider an $N \times N$ image \mathbf{x}_0 and let $1 \leq P < N$ be an integer denoting the patch size and let $k = \lfloor N/P \rfloor$. Then \mathbf{x}_0 could be covered with a $(k+1) \times (k+1)$ nonoverlapping patch grid but that would result in $(k+1)P - N$ additional rows and columns for the patches. Hence, we zero pad \mathbf{x}_0 on all four sides by $M = (k+1)P - N$ to form a new image with $N+2M$ rows and columns. With slight abuse of notation, we also denote this larger image by \mathbf{x} . Let i, j be positive integers between 1 and M inclusive. Fig. 4.3 illustrates that we may partition \mathbf{x} into $(k+1)^2 + 1$ regions as follows: $(k+1)^2$ of the regions consist of evenly chopping up the square consisting of rows i through $i+N+P$ and columns j through $j+N+P$ into a $k+1$ by $k+1$ grid, with each such partition being $P \times P$, and the last region consists of the remaining bordering part of \mathbf{x} that is not included in the first $(k+1)^2$ regions. This last region will always be entirely zeros, and the $(k+1)^2$ square patches fully cover the central $N \times N$ image.

Each pair of integers i and j correspond to one possible partition, so when we let i and j range through all the possible values, our proposal is to model the distribution of \mathbf{x} as the following product distribution:

$$p(\mathbf{x}) = \prod_{i,j=1}^{i,j=M} p_{i,j,B}(\mathbf{x}_{i,j,B}) \prod_{r=1}^{(k+1)^2} p_{i,j,r}(\mathbf{x}_{i,j,r}) / Z, \quad (4.1)$$

where $\mathbf{x}_{i,j,B}$ represents the aforementioned bordering region of \mathbf{x} that depends on the specific values of i and j , $p_{i,j,B}$ is the probability distribution of that region, $\mathbf{x}_{i,j,r}$ is the r th $P \times P$ patch when using the partitioning scheme corresponding to the values of i and j , $p_{i,j,r}$ is the probability distribution of that region, and Z is an appropriate scaling

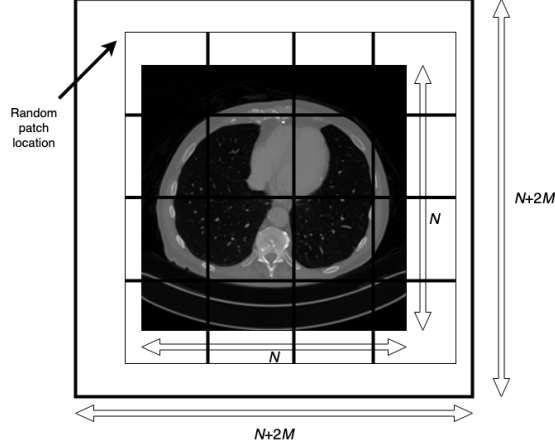


FIG 4.3 – Schematic for zero padding and partitioning image into patches

factor. Generative models based on products of patch probabilities have a long history in the Markov random field literature; see **S** A.1.5.

The score function of the entire image follows directly from (4.1):

$$\nabla \log p(\mathbf{x}) = \sum_{i,j=1}^{i,j=M} \left(\mathbf{s}_{i,j,B}(\mathbf{x}_{i,j,B}) + \sum_{r=1}^{(k+1)^2} \mathbf{s}_{i,j,r}(\mathbf{x}_{i,j,r}) \right). \quad (4.2)$$

Thus, we have expressed the score function of the whole image \mathbf{x} as sums of scores of patches $\mathbf{x}_{i,j,r}$ and the border $\mathbf{x}_{i,j,B}$. The former can be learned through score matching as in **S** 4.2.1. For the latter, by construction, if \mathbf{x} is a zero padded image then $\mathbf{x}_{i,j,B}$ is everywhere zero. Thus, for all \mathbf{x} , $D(\mathbf{x}_{i,j,B}) = \mathbb{E}[\mathbf{x}_{i,j,B}, t = 0 | \mathbf{x}_{i,j,B}, t = T]$ is everywhere zero, where D represents the denoiser function. Then the computation of its score function is trivial by Tweedie’s formula [56].

Importantly, unlike previous papers like [183] and [55], our method can compute the score function of the entire image through only inputs of patches to the neural network. This makes it possible to learn a black box function for score functions of large images, where directly training a diffusion model would be infeasible due to memory and time constraints. Furthermore, **S** 4.4 shows that in the case of limited data, the large number of patches present in each training image allows the patch-based model to learn a model for the underlying distribution that performs better than whole image models.

4.2.1 Patch-wise training

Following the work in [183] and [93], we train a denoiser using the score matching approach. Our neural network $D_\theta(\mathbf{x}, \sigma_t)$ accepts the noisy image \mathbf{x} and a noise level σ_t , and

is trained through the following loss function:

$$\mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma_t^2 I)} \|D_\theta(\mathbf{x} + \boldsymbol{\epsilon}, \sigma_t) - \mathbf{x}\|_2^2. \quad (4.3)$$

By Tweedie’s formula [56], the score function is given by $s_\theta(\mathbf{x}, \sigma_t) = (D_\theta(\mathbf{x}, \sigma_t) - \mathbf{x})/\sigma_t^2$. Here, we want to learn the score function of patches $\mathbf{x}_{i,j,r}$, so we apply (4.3) to patches of \mathbf{x} instead of the entire image. Following [183], we extract patches randomly from the zero-padded image \mathbf{x} . To incorporate the positional information of the patch, we first define the x positional array as the size $N + 2M$ 2D array consisting of the x positions of each pixel of the image scaled between -1 and 1; the y positional array is similarly defined. We then extract the corresponding patches of these two positional arrays and concatenate them along the channel dimension of the noisy image patch as inputs into the network, nothing that noise is not added to the positional arrays.

When directly applying (4.3), it would suffice to fix the patch size P and then train using size P patches exclusively. However, [183] found that by training with patches of varying sizes, training time can be reduced and the neural network learns cross-region dependencies better. Hence, we train both with patches of size P and also patches of *smaller* size, where the size is chosen according to a stochastic scheduling scheme. By using the UNet architecture in [73], the same network can take images of different sizes as input. The Appendix provides details of the experiments.

4.2.2 Sampling and reconstruction

The proposed patch-based method for learning $p(\mathbf{x})$ may be paired with any method that would otherwise be used for sampling with a full image diffusion model, such as Langevin dynamics [166] and DDPM [73], as well as acceleration methods such as second-order solvers [93] and DDIM [165]. At training time, the network only receives patches of the image as input, along with the positional information of the patch. Nevertheless, we show that when the number of sampling iterations is sufficiently large, the proposed method is still capable of generating reasonably realistic appearing entire images. However, our main goal is solving large-scale inverse problems, not image generation.

Computing $\mathbf{s}(\mathbf{x})$ via (4.2) would require summing over the score functions of all $(k + 1)^2$ patches a total of M^2 times (corresponding to the M^2 different ways of selecting i and j). This method would be prohibitively slow due to the size of M^2 ; instead, for each iteration we randomly choose integers i and j between 1 and M and estimate $\mathbf{s}(\mathbf{x})$ using just that corresponding term of the outer summation.

Algorithm 2 Patch Diffusion Inverse Solver (PaDIS)

Require: $\sigma_1 < \sigma_2 < \dots < \sigma_T$, $\epsilon > 0$, $\zeta_i > 0$, P, M, \mathbf{y}

Initialize $\mathbf{x} \sim \mathcal{N}(0, \sigma_T^2 \mathbf{I})$

for $t = T : 1$ **do**

 Sample $z \sim \mathcal{N}(0, \sigma_t^2 \mathbf{I})$

 Set $\alpha_t = \epsilon \cdot \sigma_t^2$

 Randomly select integers $i, j \in [1, M]$

 For all $1 \leq r \leq (k + 1)^2$, extract patch $\mathbf{x}_{i,j,r}$

 Compute $D_{i,j,r} = D_\theta(\mathbf{x}_{i,j,r}, \sigma_t)$

 Set $\mathbf{s}_{i,j,r} = (D_{i,j,r} - \mathbf{x}_{i,j,r}) / \sigma_t^2$

 Apply (4.2) to get $\mathbf{s} = \mathbf{s}(\mathbf{x}, \sigma_t)$

 Set \mathbf{x} to $\mathbf{x} - \zeta_t \nabla_{\mathbf{x}_t} \|\mathbf{y} - \mathcal{A}(D)\|_2^2$

 Set \mathbf{x} to $\mathbf{x} + \frac{\alpha_t}{2} \mathbf{s} + \sqrt{\alpha_t} \mathbf{z}$

end for

For solving inverse problems with diffusion models, there are general algorithms e.g., [46] and [39], as well as more model-specific algorithms, e.g., [95] and [181]. Here, we demonstrate that our method applies to a broad range of inverse problems, and opt to use generalizable methods that do not rely on any special properties (such as the singular value decomposition of the system matrix as in [95], [181]) of the forward operator. We found that DPS [39] in conjunction with Langevin dynamics generally yielded the most stable and high quality results, so we use this approach as our central algorithm. Similar to [39], we chose the stepsize to be $\zeta_i = \zeta / \|\mathbf{y} - \mathcal{A}(D(\mathbf{x}))\|_2$. To the best of our knowledge, this is the first work that learns a fully patch-based diffusion prior and applies it to solve inverse problems; we call our method **Patch Diffusion Inverse Solver (PaDIS)**. Computing the score functions of the patches at each iteration is easily parallelized, allowing for fast sampling and reconstruction. Alg. 2 shows the pseudocode for the main image reconstruction algorithm; the appendix shows the pseudocode for the other implemented algorithms.

Finally, we comment on some high-level similarities between our proposed method and [109]; in both cases, the image in question is partitioned into smaller parts in multiple ways. In [109], one of the partitions consists of 2D slices in the x-y direction, and the other partition consists of 2D slices in the x-z direction, whereas for our method, each of the partitions consists of $(k + 1)^2$ square patches and the outer border region. For both methods, the score functions of each of the parts are learned independently during training. Then for each sampling iteration, both methods involve choosing one of the partitions, computing the score functions for each of the parts independently, and then

updating the entire image by updating the parts separately. For our approach, the zero-padding method allows for many possible partitions of the image and eliminates boundary artifacts between patches.

4.3 A Unified Theory of Patch-based Diffusion Models

4.3.1 Motivation

Although patch-based diffusion models have demonstrated success in solving inverse problems [81, 82], certain aspects of the training and reconstruction process lack theoretical justification. Furthermore, multiple other works that appear to use a different framework turn out to share the same framework as these previous two works. This section¹ provides a theoretical justification for much of the training and reconstruction process of several works in the field, as well as a series of experiments that deepen our understanding of the impacts of different types of patch-based diffusion models on inverse problem solving.

4.3.2 Methods

We begin by formulating a unified patch framework that encapsulates various methods for combining patch scores into the score of an entire image. To start, let the entire (possibly zero padded) image be denoted by \mathbf{x} . Then a general patch-based model for a denoiser of \mathbf{x} is given by

$$\mathbf{D}(\mathbf{x}) = \left(\sum_c \mathbf{G}'_c \mathbf{D}(\mathbf{G}_c \mathbf{x}) \right) \oslash \left(\sum_c \mathbf{G}'_c \mathbf{G}_c \mathbf{1} \right), \quad (4.4)$$

where the sums iterate over all (possibly overlapping) patches c , \oslash denotes element-wise division, and the $\mathbf{1}$ represents an image of all ones. This model states that to denoise an entire image, we can denoise the patches individually, add together the contributions from each of the patches to the whole image, and then perform pixel-wise division to account for the number of patches that contribute to each pixel. We will show later how PaDIS and [142] fall under this framework.

We start by showing that there is a straightforward equivalence between this equation for denoisers and a similar equation for score functions. For convenience, define $\mathbf{M} = \sum_c \mathbf{G}'_c \mathbf{G}_c \mathbf{1}$. By Tweedie’s formula, $\mathbf{s}(\mathbf{x}) = (\mathbf{D}(\mathbf{x}) - \mathbf{x})/\sigma^2$. Then by using the fact that

¹This section is not in the paper [82] that forms the basis for the rest of this chapter.

$(\sum_c \mathbf{G}'_c \mathbf{G}_c \mathbf{x}) \oslash \mathbf{M} = \mathbf{x}$, one can readily use algebra to derive the relationship

$$\mathbf{s}(\mathbf{x}) = \left(\sum_c \mathbf{G}'_c \mathbf{s}_V(\mathbf{G}_c \mathbf{x}) \right) \oslash \mathbf{M}, \quad (4.5)$$

where \mathbf{s}_V is the score function of a patch.

Following this, if \mathbf{M} is a constant M times the matrix of all ones, then the probability distribution of the entire image takes the following form:

$$p(\mathbf{x}) = \frac{1}{Z} \left(\prod_c p_c(\mathbf{G}_c \mathbf{x}) \right)^{1/M}, \quad (4.6)$$

where Z is a scaling constant.

Unfortunately, in the general case where \mathbf{M} is not a constant matrix, this “score” function need not correspond to any probability distribution, as illustrated by the following example. Firstly, score functions are the gradient of the (log of a) function, so to simplify things we rewrite (4.5) as

$$\nabla f(\mathbf{x}) = \left(\sum_c \mathbf{G}'_c \nabla f_c(\mathbf{G}_c \mathbf{x}) \right) \oslash \mathbf{M}, \quad (4.7)$$

for the patch functions f_c and overall function f . Now we take the specific example of \mathbf{x} containing just three pixels in a row, which we represent as $\mathbf{x} = [u \ v \ w]$. Then we define the first patch as the first two pixels and the second patch as the next two pixels. The functions that operate on these two patches are defined as $f_1([u \ v]) = uv$ and $f_2([v \ w]) = v^2 + 2w^2$; their gradients are $\nabla f_1 = (v, u)$ and $\nabla f_2 = (2v, 4w)$. Only the second pixel in \mathbf{x} is overlapping, so $\mathbf{M} = [1 \ 2 \ 1]$. Putting it all together gives $\nabla f(\mathbf{x}) = (v, \frac{1}{2}u + v, 4w)$. However, the curl of this vector field is readily computed to be $(0, 0, -1/2)$. As this is nonzero, this is not a conservative vector field, meaning it cannot be equal to the gradient of a function. Hence, a “score function” written in the form (4.7) need not have a corresponding probability distribution.

In the general case, this “score” function may not correspond to a probability distribution; nevertheless, we show how the training process can be simplified from the loss from entire image training. By incorporating the positional encoding c^* of patches as an input

to the network, the denoising score matching loss is given by

$$L(\theta) = \mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma_t^2 \mathbf{I})} \|\mathbf{s}_\theta(\mathbf{x} + \boldsymbol{\epsilon}, \sigma_t) + \boldsymbol{\epsilon}/\sigma_t\|_2^2 \quad (4.8)$$

$$= \mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma_t^2 \mathbf{I})} \left\| \sum_c \mathbf{G}'_c \mathbf{s}_V(\mathbf{G}_c \mathbf{x}, c^*; \theta) \odot \mathbf{M} + \boldsymbol{\epsilon}/\sigma_t \right\|_2^2. \quad (4.9)$$

Note that $(\sum_c \mathbf{G}'_c \mathbf{G}_c \boldsymbol{\epsilon}) \odot \mathbf{M} = \boldsymbol{\epsilon}$. So we rewrite the norm inside the expectation as

$$L = \left\| \sum_c \mathbf{G}'_c \mathbf{s}_V(\mathbf{G}_c \mathbf{x}, c^*; \theta) \odot \mathbf{M} + \sum_c \mathbf{G}'_c \mathbf{G}_c \boldsymbol{\epsilon} \odot \mathbf{M} / \sigma_t \right\|_2^2 \quad (4.10)$$

$$= \left\| \left(\sum_c \mathbf{G}'_c (\mathbf{s}_V(\mathbf{G}_c \mathbf{x}, c^*; \theta) + \mathbf{G}_c \boldsymbol{\epsilon} / \sigma^2) \right) \odot \mathbf{M} \right\|_2^2. \quad (4.11)$$

Assuming that the patches collectively cover the entire image, each element of \mathbf{M} is at least 1. Thus we have

$$L \leq \left\| \sum_c \mathbf{G}'_c (\mathbf{s}_V(\mathbf{G}_c \mathbf{x}, c^*; \theta) + \mathbf{G}_c \boldsymbol{\epsilon} / \sigma^2) \right\|_2^2. \quad (4.12)$$

Now we have a term of the form $\|X_1 + X_2 + \dots + X_n\|_2^2$, which we can bound using a similar derivation as (6.6) in Ch. 6 (derived in the appendix of Ch. 5):

$$L \leq K \sum_c \left\| \mathbf{G}'_c (\mathbf{s}_V(\mathbf{G}_c \mathbf{x}, c^*; \theta) + \mathbf{G}_c \boldsymbol{\epsilon} / \sigma^2) \right\|_2^2, \quad (4.13)$$

where K is a constant independent of \mathbf{x} . Finally, since \mathbf{G}'_c is an operator restoring a patch to its location in a full image, it is invariant under the Euclidean norm. Putting everything together, we have

$$L(\theta) \leq K \cdot \mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma_t^2 \mathbf{I})} \mathbb{E}_{\text{random } c} \left\| \mathbf{s}_V(\mathbf{G}_c \mathbf{x}, c^*; \theta) + \mathbf{G}_c \boldsymbol{\epsilon} / \sigma_t^2 \right\|_2^2. \quad (4.14)$$

This loss function indicates that (perhaps surprisingly!) regardless of the method by which the patches are used to tile the image, the method of training the patch-based network is the same: extract patches randomly from whole images in the training dataset, and then perform score matching on the resulting patches. This principle extends to overlapping patches and even patches of different sizes. The following subsections demonstrate how various patch-based diffusion models fall under this framework.

Stochasticity Method Evaluation. For certain probability models, it is too computationally expensive to evaluate the score functions of all the patches involved, so we instead stochastically select some patches and perform an update according to them. We show next that using stochastic updates leads to convergence to the same distribution as if we used the entire score function.

Theorem. Suppose $\mathbf{s}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \mathbf{s}_i(\mathbf{x})$. When sampling using Langevin dynamics [166], instead of using $\mathbf{s}(\mathbf{x})$, we use $\mathbf{s}_i(\mathbf{x})$ for a randomly chosen i . Then under certain regularity conditions, the sampling algorithm will lead to an image that is drawn from the distribution $p(\mathbf{x})$.

Proof. The update step that uses the full score function is

$$\Delta \mathbf{x}_t = \frac{\epsilon_t}{2} \cdot \frac{1}{N} \sum_{i=1}^N \mathbf{s}_i(\mathbf{x}_t) + \sqrt{\epsilon_t} \mathbf{z}, \quad (4.15)$$

whereas the update step in the stochastic algorithm is

$$\Delta \mathbf{x}_t = \frac{\epsilon_t}{2} \cdot \mathbf{s}_j(\mathbf{x}_t) + \sqrt{\epsilon_t} \mathbf{z}. \quad (4.16)$$

We wish to show that the distribution of random vectors in sequence $\{\mathbf{x}_t\}$ generated by the stochastic algorithm converge to samples drawn from the same distribution as the algorithm using the full score function.

The crux of the proof is the following idea. We first fix some $\epsilon_0 \in (0, 1)$. By the infinite sum properties on ϵ_t , there is a subsequence $t_1 < t_2 < \dots$ such that $\sum_{t=t_s+1}^{t_{s+1}} \epsilon_t$ approaches ϵ_0 as s approaches infinity. Now we collectively consider the $t_{i+1} - t_i$ update steps belonging to one term of the subsequence. Across these steps, there are mainly two sources of randomness: the stochasticity of choosing which score function to use, and the added Gaussian noise. We will show that the amount of deviation induced by the stochastic score functions is $O(\epsilon_0)$ while that of the Gaussian noise is $O(\sqrt{\epsilon_0})$. Thus, taking the limit as ϵ_0 approaches 0 shows that almost all the variance will be due to the Gaussian noise. Hence, these $t_{i+1} - t_i$ update steps are equivalent to a single Langevin update step using the full score function with step size ϵ_0 . Thus, the stochastic algorithm converges to the same distribution as the algorithm using the full score function.

To show this, first, define $g(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \mathbf{s}_i(\mathbf{x})$ and $h_t(\mathbf{x}) = \mathbf{s}_{j_t}(\mathbf{x}) - g(\mathbf{x})$. Then the stochastic gradient at time t is $g(\mathbf{x}) + h_t(\mathbf{x})$. The total update due to gradients between steps t_s and t_{s+1} is

$$\sum_{t=t_s+1}^{t_{s+1}} \frac{\epsilon_t}{2} (g(\mathbf{x}_t) + h_t(\mathbf{x}_t)). \quad (4.17)$$

We analyze the two terms of this sum separately. First,

$$\sum_{t=t_s+1}^{t_{s+1}} \frac{\epsilon_t}{2} g(\mathbf{x}_t) = \sum_{t=t_s+1}^{t_{s+1}} \frac{\epsilon_t}{2} g(\mathbf{x}_{t_s}) + \sum_{t=t_s+1}^{t_{s+1}} \frac{\epsilon_t}{2} (g(\mathbf{x}_t) - g(\mathbf{x}_{t_s})). \quad (4.18)$$

The first term here is just $\frac{\epsilon_0}{2} g(\mathbf{x}_{t_s})$, which is an update step using the gradient of the full score function. Therefore we need to bound the standard deviation of the error term:

$$E = \left| \sum_{t=t_s+1}^{t_{s+1}} \frac{\epsilon_t}{2} (g(\mathbf{x}_t) - g(\mathbf{x}_{t_s})) \right| \leq \sum_{t=t_s+1}^{t_{s+1}} \frac{\epsilon_t}{2} |g(\mathbf{x}_t) - g(\mathbf{x}_{t_s})|. \quad (4.19)$$

Assuming that g has a local Lipschitz constant L ,

$$E \leq \sum_{t=t_s+1}^{t_{s+1}} \frac{\epsilon_t L}{2} \|\mathbf{x}_t - \mathbf{x}_{t_s}\|_2 \leq \frac{\epsilon_0 L}{2} \max_{t_s < t \leq t_{s+1}} \|\mathbf{x}_t - \mathbf{x}_{t_s}\|_2. \quad (4.20)$$

Now it suffices to note that as ϵ_0 goes to 0, so do all the step sizes within the subsequence. Hence, $\|\mathbf{x}_t - \mathbf{x}_{t_s}\|_2$ goes to 0 (at a rate that is unknown but not crucial to the analysis). Therefore, we can write

$$E \leq \frac{\epsilon_0 L}{2} o(1) \leq O(\epsilon_0). \quad (4.21)$$

In summary we have shown that

$$\sum_{t=t_s+1}^{t_{s+1}} \frac{\epsilon_t}{2} g(\mathbf{x}_t) = \frac{\epsilon_0}{2} g(\mathbf{x}_{t_s}) + O(\epsilon_0). \quad (4.22)$$

Now we analyze the second term which is:

$$\sum_{t=t_s+1}^{t_{s+1}} \frac{\epsilon_t}{2} h_t(\mathbf{x}_t) = \sum_{t=t_s+1}^{t_{s+1}} \frac{\epsilon_t}{2} h_t(\mathbf{x}_{t_s}) + \sum_{t=t_s+1}^{t_{s+1}} \frac{\epsilon_t}{2} (h_t(\mathbf{x}_t) - h_t(\mathbf{x}_{t_s})). \quad (4.23)$$

Similar to the derivation for g , the second term is readily bounded by $O(\epsilon_0)$ assuming that h_t is Lipschitz. For the first term, note that each $h_t(\mathbf{x}_{t_s})$ is an iid random variable. Furthermore, the variance of each such random variable is bounded by the maximum values taken by $s_i(\mathbf{x})$ along the trajectory of \mathbf{x} , which in turn is assumed to be bounded. Hence, we get a weighted average of iid random variables with bounded variance. In particular, the sum of the weights is $O(\epsilon_0)$, so the average deviation due to the entire term is $O(\epsilon_0)$.

Putting everything together,

$$\sum_{t=t_s+1}^{t_{s+1}} \frac{\epsilon_t}{2} (g(\mathbf{x}_t) + h_t(\mathbf{x}_t)) = \frac{\epsilon_0}{2} g(\mathbf{x}_{t_s}) + O(\epsilon_0). \quad (4.24)$$

Now the variance due to the added Gaussian noise is simple: at each step, the noise has variance ϵ_t and it is independent. Thus the total noise has variance ϵ_0 , which means the average deviation due to this noise is $O(\sqrt{\epsilon_0})$. Hence, as ϵ_0 approaches 0, the random error from the stochastic gradient is negligible compared to the added Gaussian noise. This shows that across the steps t_s to t_{s+1} , the expected value of the update is $\frac{\epsilon_0}{2} g(\mathbf{x}_{t_s})$ and the variance of the added Gaussian noise is ϵ_0 . This variance aligns exactly with a Langevin update step using the full score function with a step size of ϵ_0 . Hence, we have shown that the samples generated by the stochastic algorithm converges to the same probability distribution as using the full score function.

4.3.3 Examples

Next, we provide specific examples of works that fall under this framework.

Image restoration in adverse weather conditions. The work of [142] uses a supervised diffusion model to perform image restoration for images taken in adverse weather conditions. In this paper, the learned patch prior is conditioned on a degraded version of the clean image. We can readily apply this method for learning an unsupervised prior by removing the condition during training. Then the restoration algorithm consists of the step $\Omega_t = \Omega_t + P_d \cdot \epsilon_\theta(\mathbf{x}_t^{(d)}, t)$, where Ω_t is the noise of the whole image at the t th timestep, d is an iterator over the overlapping patches, P_d is an operator that maps a patch to its location in the whole image, and $\mathbf{x}_t^{(d)}$ is an image patch. This step is followed by the step $M = M + P_d \cdot \mathbf{1}$. At the end of each diffusion timestep, the noise is scaled accordingly with the step $\Omega_t = \Omega_t \odot M$. This definition of M is consistent with our definition of M and the update steps follow (4.4). Finally, the training process outlined in [142] involves selecting a random patch of a clean image and performing denoising score matching, consistent with (4.14).

PaDIS. Our proposed method of using patch shifts with nonoverlapping patches falls under this general framework [82]. We start with the score function model in (4.2) and rewrite it using denoisers as follows:

$$\mathbf{D}(\mathbf{x}) = \frac{1}{M^2} \sum_{i,j=1}^{i,j=M} \left(\mathbf{D}_{i,j,B}(\mathbf{x}_{i,j,B}) + \sum_{r=1}^{(k+1)^2} \mathbf{D}_{i,j,r}(\mathbf{x}_{i,j,r}) \right), \quad (4.25)$$

where we include an additional $1/M^2$ factor for consistency with later derivations. We drop the \mathbf{G}'_c terms that map the denoised patches back to the whole image for convenience. By definition, the denoised border $\mathbf{D}_{i,j,B}(\mathbf{x}_{i,j,B})$ is everywhere zero. Note that for each of the M^2 random patch shifts corresponding to a pair of indices i and j , the resulting partition of the zero padded image consists of nonoverlapping patches. Interpreting each such shift as containing $(k+1)^2 + 1$ patches, the sum in total contains $((k+1)^2 + 1)M^2$ patches. Furthermore, each pixel is included in exactly M^2 of the patches. Hence, we may rewrite (4.25), this time with all the patch extracting operators as

$$\mathbf{D}(\mathbf{x}) = \left(\sum_c \mathbf{G}'_c \mathbf{D}(\mathbf{G}_c \mathbf{x}) \right) \oslash (M^2 \cdot \mathbf{1}), \quad (4.26)$$

where $\mathbf{1}$ is the matrix of all ones. Now we are iterating over all $((k+1)^2 + 1)M^2$ patches using a single iterator c . This sum encapsulates both of the sums in (4.25) so that M^2 of the terms consists of denoising the image border while the remaining $(k+1)^2 M^2$ terms consist of denoising square patches. Finally, with this tiling scheme, it is clear that $\sum_c \mathbf{G}'_c \mathbf{G}_c \mathbf{1} = M^2 \mathbf{1}$. Hence, we have shown how the score model for PaDIS follows (4.4).

Strictly speaking, when using (4.4) to update the image, it would be necessary to compute the score function of every patch in the sum.

TPDM. Used for solving 3D inverse problems, the probability distribution model used in TPDM [109] is given by

$$p(\mathbf{x}) = \left(\prod_{i=1}^N q_\theta(\mathbf{x}[:, :, i])^\alpha \right) \left(\prod_{j=1}^N q_\phi(\mathbf{x}[j, :, :])^\beta \right) / Z, \quad (4.27)$$

where α and β are fixed constants and Z is a scaling constant. This model uses a total of $2N$ patches used to cover the volume, where N of them are along the first indexing dimension and N of them are along the third indexing dimension. Each voxel in the volume is covered by exactly two such patches, so $\mathbf{M} = 2 \cdot \mathbf{1}$. Therefore, by setting $\alpha = \beta = 1/2$, we recover the form of (4.6). Furthermore, training is done separately on the 2D slices in the two perpendicular directions of the volumes. Similar to PaDIS, while it would be necessary to use both the distributions q_θ and q_ϕ to update the volume at each step of the diffusion process during reconstruction, the TPDM approach instead takes a stochastic approach of choosing one of the two possible ways.

DiffusionBlend++. The probability distribution model used in DiffusionBlend++ (see Ch. 5) is given by

$$\mathbf{s}(\mathbf{x}) = |S|^{-1} \sum_{\mathcal{S}=\mathcal{S}_1 \cup \dots \cup \mathcal{S}_r} \sum_{i=1}^r \nabla \log p(\mathbf{x}[:, :, \mathcal{S}_i]), \quad (4.28)$$

where the outer sum runs over (possibly a subset of) all possible ways to partition the N slices of the volume into r sets of N/r elements each. This is very similar to PaDIS: \mathbf{M} is simply a constant times the matrix of all ones, assuming that each slice is included in the same number of partitions. Furthermore, at reconstruction time, DiffusionBlend++ again uses a stochastic approach to avoid computing the entire score function in (4.28).

In conclusion, we have provided a general framework for patch-based diffusion models and identified several recent works that fall under this framework. Furthermore, we have shown some theoretical justifications and insights that enhance our understanding of these models. In the future, this section would benefit from a series of carefully designed experiments that verify the theoretical results and examine differences in results between possible design choices for patches.

4.4 Experiments

Experimental setup. For the main CT experiments, we used the AAPM 2016 CT challenge data [135] that consists of 10 3D volumes. We cropped the data in the Z-direction to select 256 slices and then rescaled the images in the XY-plane to have size 256×256 . Finally, we used the XY slices from 9 of the volumes to define 2304 training images, and used 25 of the slices from the tenth volume as test data. For the deblurring and superresolution experiments, we used a 3000 image subset of the CelebA-HQ dataset [126] (with each image being of size 256×256) for training to demonstrate that the proposed method works well in cases with limited training data. We preprocessed the data by dividing all the RGB values by 255 so that all the pixel values were between 0 and 1. The test data was a randomly selected subset of 25 of the remaining images. In all cases, we report the average metrics across the test images: peak SNR (PSNR) in dB, and structural similarity metric (SSIM) [185]. For the colored images, these metrics were computed in the RGB domain.

For the main patch-based networks, we trained mostly with a patch size of 56×56 to allow the target image to be completely covered by 5 patches in each direction while minimizing the amount of zero padding needed. We used the network architecture in [93] for both the patch-based networks and whole image networks. All networks were trained

on PyTorch using the Adam optimizer with 2 A40 GPUs. The Appendix provides full details of the hyperparameters.

Image generation. Our proposed method is able to learn a reasonable prior for whole images, despite never being trained on any whole images. Fig. 4.4 shows generation results for the CT dataset using three different methods. The top row used the network trained on whole images; the middle row used the method of [183] except that the entire image is never supplied to the network either at training or sampling time; the bottom row used the proposed method. The middle row shows that the positional encoding does ensure reasonably appropriate generated patches at each location. However, simply generating each of the patches independently and then naively assembling them together leads to obvious boundary artifacts due to lack of consistency between patches. Our proposed method solves this problem by using overlapping patches via random patch grid shifts, leading to generated images with continuity throughout. However, some of these images still appear asymmetric or have otherwise unrealistic global structure. This is expected as the network is not trained on regions of the image that are far apart and thus cannot learn long range correlations; in a later chapter, we will develop an improved method that can mitigate this issue.

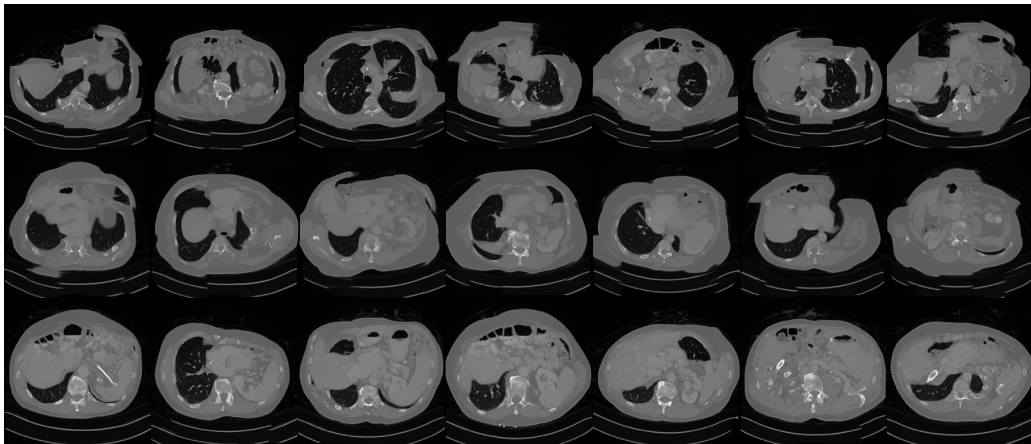


FIG 4.4 – Unconditional generation of CT images. Top row: generation with a network trained on whole image; middle row: patch-only version of [183]; bottom row: proposed PaDIS method.

Inverse problems. We tested the proposed method on a variety of different inverse problems: CT reconstruction, deblurring, and superresolution. For the forward and backward projectors in CT reconstruction, we used the implementation provided by [173]. We performed two sparse view CT (SVCT) experiments: one using 8 projection views, and

one using 20 projection views. Both of these were done using a parallel-beam forward projector where the detector size was 512 pixels. For the deblurring experiments, we used a uniform blur kernel of size 9×9 and added white Gaussian noise with $\sigma = 0.01$ where the clean image was scaled between 0 and 1. For the superresolution experiments, we used a scaling factor of 4 with downsampling by averaging and added white Gaussian noise with $\sigma = 0.01$. DPS has shown to benefit significantly from using a higher number of neural function evaluations (NFEs) [39], so we use 1000 NFEs for all of the diffusion model experiments.

For the comparison methods, we trained a diffusion model on entire images using the same denoising score matching method shown in Section 4.2.1. The inference process was identical to that of the patch-based method, with the exception that the score function of the image at each iteration was computed directly using the neural network, as opposed to needing to break up the zero-padded image into patches. We also compared with two traditional methods: applying a simple baseline and reconstructing via the total variation regularizer (ADMM-TV). 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 [74]. We also implemented two plug and play (PnP) methods: PnP-ADMM [196] and regularization by denoising (RED) [83]. We trained denoising CNNs on each of the datasets following [151] and then used them to solve the inverse problems.

For further comparison of diffusion model methods, we implemented different sampling and data consistency algorithms and applied them in conjunction with our patch-based prior. In particular, we compared with Langevin dynamics [166], predictor-corrector sampling [46], and variation exploding DDNM (VE-DDNM) [181]. For all these sampling methods, we used the variation exploding framework for consistency and to avoid retraining the network. We also compared with two other ways of assembling priors of patches to form the prior of an entire image: patch averaging [142] and patch stitching [152].

Table 4.1 shows the main inverse problem solving results. The best results were obtained after training the patch-based models for around 12 hours, while the whole image models needed to be trained for 24-36 hours, demonstrating a significant improvement in training time. Furthermore, when averaging the metrics across the test dataset, our proposed method outperformed the whole image method in terms of PSNR and SSIM for all the inverse problems. The score functions of all the patches can be computed in parallel for each iteration, so the reconstruction times for these methods were very similar (both

around 5 minutes per image). The whole-image results could be more favorable if more training data were used.

We also ran ablation studies examining the effect of various parameters of the proposed method. Namely, we studied the usage of different patch sizes during reconstruction, varying dataset sizes, importance of positional encoding for patches, and different sampling and inverse problem solving algorithms. The results of these studies are in App. A.1.1.

TBL 4.1 – Comparison of quantitative results on three different inverse problems. Results are averages across all images in the test dataset. Best results are in bold.

Method	CT, 20 Views		CT, 8 Views		Deblurring		Superresolution	
	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow
Baseline	24.93	0.595	21.39	0.415	24.54	0.688	25.86	0.739
ADMM-TV	26.82	0.724	23.09	0.555	28.22	0.792	25.66	0.745
PnP-ADMM [196]	26.86	0.607	22.39	0.489	28.82	0.818	26.61	0.785
PnP-RED [83]	27.99	0.622	23.08	0.441	29.91	0.867	26.36	0.766
Whole image diffusion	32.84	0.835	25.74	0.706	30.19	0.853	29.17	0.827
Langevin dynamics [166]	33.03	0.846	27.03	0.689	30.60	0.867	26.83	0.744
Predictor-corrector [46]	32.35	0.820	23.65	0.546	28.42	0.724	26.97	0.685
VE-DDNM [181]	31.98	0.861	27.71	0.759	-	-	26.01	0.727
Patch Averaging [142]	33.35	0.850	28.43	0.765	29.41	0.847	27.67	0.802
Patch Stitching [152]	32.87	0.837	26.71	0.710	29.69	0.849	27.50	0.780
PaDIS (Ours)	33.57	0.854	29.48	0.767	30.80	0.870	29.47	0.846

In addition to the main inverse problems shown in Table 4.1, we also ran experiments on three additional inverse problems to demonstrate the effectiveness of our method on a wide variety of inverse problems: 60 view CT reconstruction, fan beam reconstruction using 180 views, and deblurring with a uniform kernel of size 17×17 . For the deblurring experiment, we again added Gaussian noise with level $\sigma = 0.01$ to the measurement. Table 4.2 shows the quantitative results of these experiments and Figure A.6 shows the visual results.

In the bottom of Figure 4.5, some artifacts are present in the reconstructions obtained by the diffusion model methods, although they are more apparent in the whole image model than with PaDIS. The measurements are very compressed in this case, so it is very difficult for any model to obtain diagnostic-quality reconstructions; the baselines perform

TBL 4.2 – Extra inverse problem experiments. Results are averages across all images in the test dataset. Best results are in bold.

Method	CT, 60 Views		CT, Fan Beam		Heavy Deblurring	
	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow
Baseline	25.89	0.746	20.07	0.521	21.14	0.569
ADMM-TV	30.93	0.833	25.78	0.719	26.03	0.724
Whole image diffusion	35.83	0.894	26.89	0.835	28.35	0.808
PaDIS (Ours)	39.28	0.941	29.91	0.932	28.91	0.818

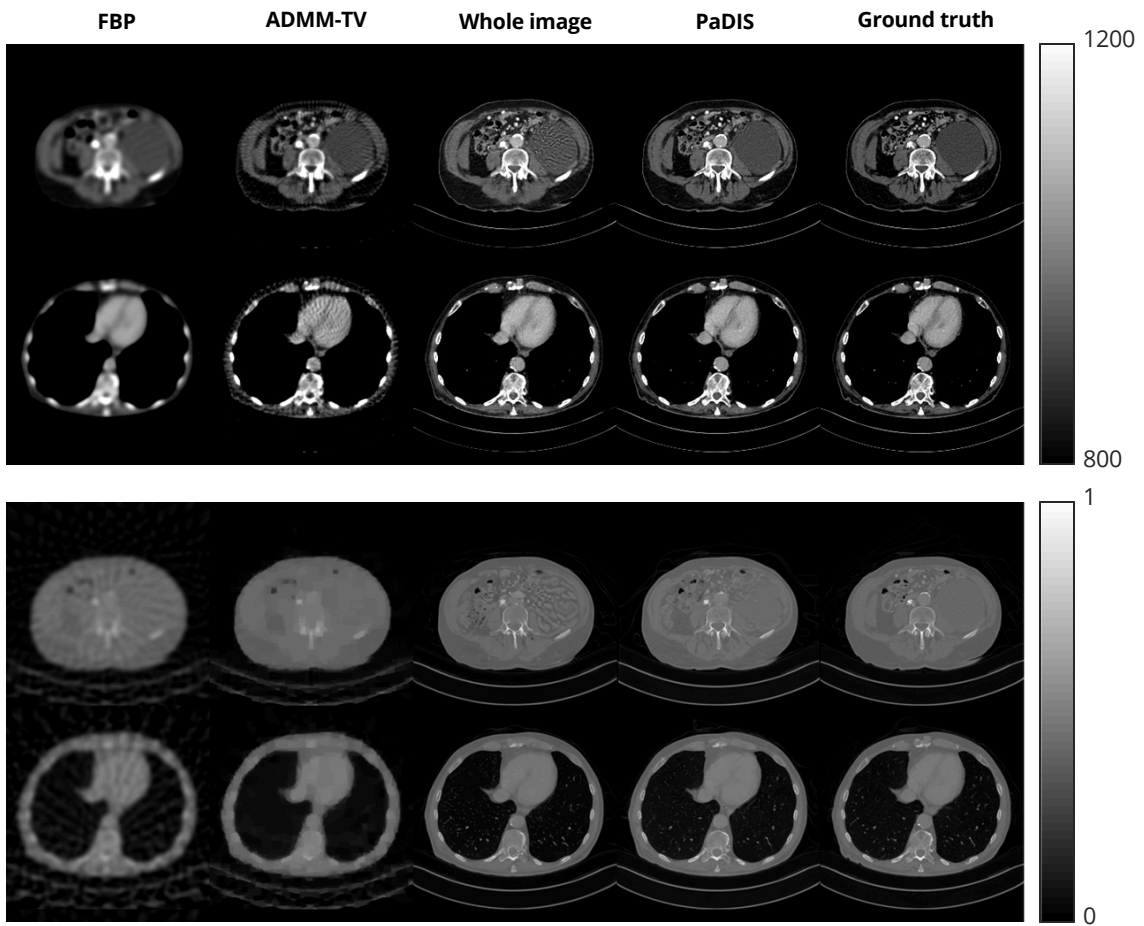


FIG 4.5 – Results of CT reconstruction. 60 views are used for the top two rows, 20 views are used for the bottom two rows. To better show contrast between organs, we use modified Hounsfield units (HU) in the top figure, while we use the same scale the images were trained on in the bottom figure.

significantly worse in terms of quantitative metrics and exhibit severe blurring. In clinical settings, patient diagnosis are typically performed with CT scans consisting of hundreds

of views. The top of Figure 4.5 shows that when 60 views are used, our proposed method yields a much better reconstruction without artifacts. Nevertheless, we show the potential of our proposed methods to reconstruct images from very sparse views with a decent image quality, which could be potentially used for applications such as patient positioning.

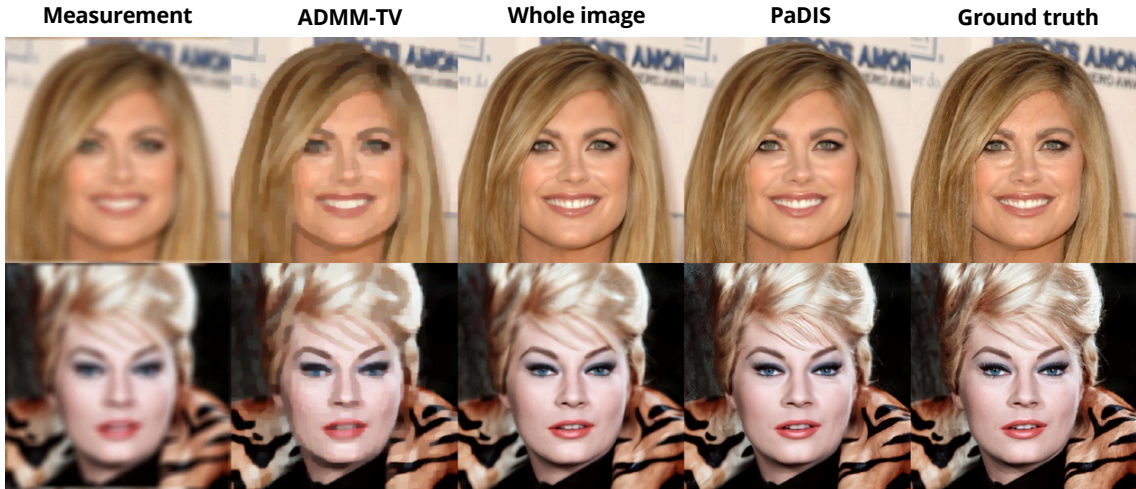


FIG 4.6 – Results of deblurring with Gaussian noise ($\sigma = 0.01$).

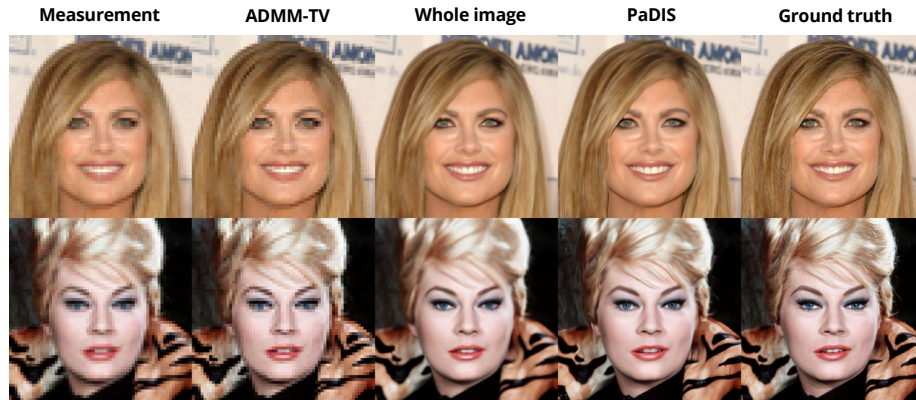


FIG 4.7 – Results of superresolution with Gaussian noise ($\sigma = 0.01$).

Finally, since the original AAPM dataset contained CT images of resolution 512×512 , we ran 60 view CT reconstruction experiments with these higher resolution images. Due to the lack of data, we did not train a whole-image model on these higher resolution images. We used a zero padding width of 64 and largest patch size of 64×64 for training. Figure 4.8 shows the visual results of these experiments. Hence, our proposed methods

can obtain high quality reconstructions for both different inverse problems and also for different image sizes.

TBL 4.3 – Results of 60 view CT reconstruction with 512×512 images. Results are averages across all images in the test dataset. Best results are in bold.

Method	FBP	ADMM-TV	PaDIS
PSNR \uparrow	28.38	29.48	36.93
SSIM \uparrow	0.699	0.788	0.899

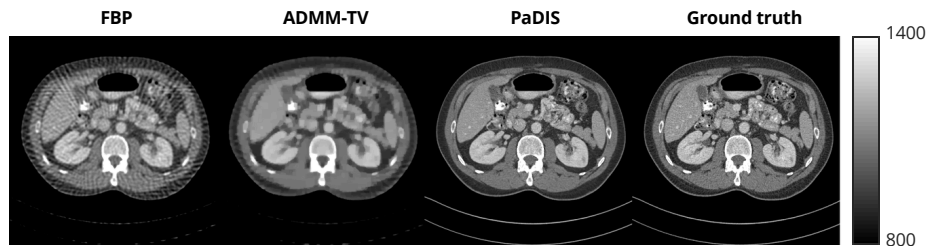


FIG 4.8 – Results of 60 view CT reconstruction on 512×512 images using modified HU units.

4.5 Conclusion

In this work, we presented a method of using score-based diffusion models to learn image priors through solely the patches of the image, combined with suitable position encoding. Simulation results demonstrated how the method can be used to solve a variety of inverse problems. Extensive experiments showed that under conditions of limited training data, the proposed method outperforms methods involving whole image diffusion models. In the future, more work could be done on higher quality image generation using a multi-scaled resolution approach [68, 72], using acceleration methods for faster reconstruction, and higher dimensional image reconstruction. Image priors like those proposed in this chapter have the potential to benefit society by reducing X-ray dose in CT scans. Generative models have the risk of inducing hallucinations and being used for disinformation.

CHAPTER 5

DiffusionBlend: Learning 3D Image Prior through Position-aware Diffusion Score Blending for 3D Computed Tomography Reconstruction

5.1 Motivation

Learning efficient 2D image priors using diffusion models is already computationally expensive, which requires large-scale of training data, training time, and GPU memory. For example, previous works [168, 73] require training for several days to weeks on over a million training images in the ImageNet [153] and LSUN [200] datasets to generate high-quality 2D natural images of size 256×256 . Hence, directly learning a 3D diffusion prior on the entire CT volume would be practically infeasible or prohibitively expensive due to the demanding requirements of training data and computational cost. In addition, real clinical CT data is usually limited and scarce and often has a resolution larger than $256 \times 256 \times 400$, which makes directly training the data prior very challenging. The problem of tackling 3D image inverse problems, especially for medical imaging remains a challenging open research question.

A few recent works [42, 109, 41] have discussed and proposed to solve 3D image reconstruction problems either through employing some hand-crafted regularization to enforce consistency between 2D slices when reconstructing 3D volumetric images [41, 42], or through training several diffusion models for 2D images on each plane (axial, coronal,

This chapter based on [162]. Equal contribution with Bowen Song.

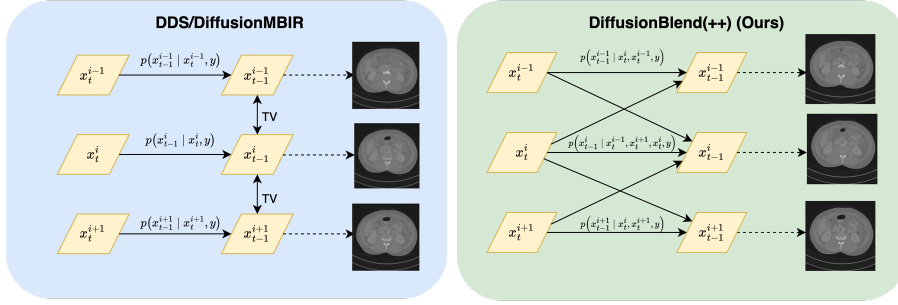


FIG 5.1 – Overview of DiffusionBlend++ compared to previous 3D image reconstruction works. Previous work used a hand-crafted TV term to “regularize” adjacent slices, whereas the proposed approach uses learned diffusion score blending between groups of slices. Here i is the slice index, and t is the reconstruction iteration.

and sagittal), and performing reverse sampling with each model alternatively [109]. However, all of these works only learn the distribution of a single 2D slice via the diffusion model, while having not yet explored the dependency between slices that is required to better model the real 3D image prior.

To overcome these limitations, we propose a novel method, called DiffusionBlend, that enables learning the distribution of 3D image patches (a batch of nearby 2D slices), and blends the scores of patches to model the entire 3D volume distribution for image reconstruction. Specifically, we firstly propose to train a diffusion model that models the joint distribution of 3D image patches (nearby 2D slices) in the axial plane conditioning on the slice thickness. Then, we introduce a random blending algorithm that approximates the score function of the entire 3D volume by using the trained 3D-patch score function. Moreover, we can either directly use the trained model to predict the noise of a single 2D slice by taking its corresponding 3D patch as input, or applying a random blending algorithm that firstly randomly partitions the volume into different 3D patches at each time step and then computes the score of each 3D patch during reverse sampling. Through either way, we can output the predicted noise of the entire 3D volume. In this way, our proposed method is able to enforce cross-slice consistency without any hand-crafted regularizer. Our method has the advantage of being fully data-driven and can enforce slice consistency without the TV regularizer as demonstrated in Fig. 5.1. Through exhaustive experiments of ultra-sparse-view and limited-angle 3D CT reconstruction on different datasets, we validate that our proposed method achieves superior reconstruction results for 3D volumetric imaging, outperforming previous state-of-the-art (SOTA) methods. Furthermore, our method achieves better or comparable inference time than SOTA methods, and requires minimum hyperparameter tuning for different tasks and settings.

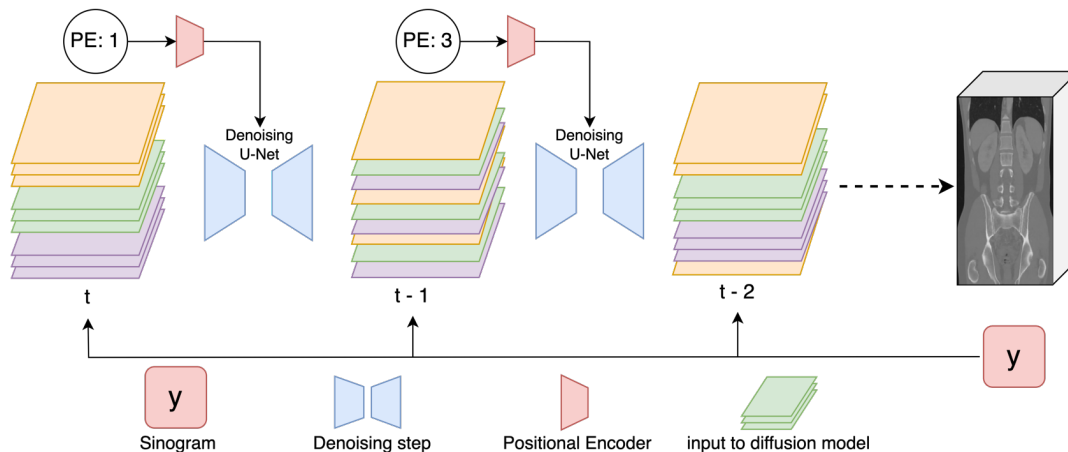


FIG 5.2 – Overview of slice blending process during reconstruction for DiffusionBlend++. At each iteration, we partition the slices of the volume in a different way; slices of the same color are inputted into the network independently. Positional encoding (PE) is also inputted to the network as information about the separation between the slices.

In summary, our main contributions are as follows:

- We propose DiffusionBlend(++): a novel method for 3D medical image reconstruction through 3D diffusion priors. To the best of our knowledge, our method is the first diffusion-based method that learns the 3D-patch image prior incorporating the cross-slice dependency, so as to enforce the consistency for the entire 3D volume without any external regularization.
- Specifically, instead of independently training a diffusion model only on separated 2D slices, we propose a novel method that first trains a diffusion model on 3D image patches (a batch of nearby 2D slices) with positional encoding, and at inference time, employs a new approach of random partitioning and diffusion score blending to generate an isotropically smooth 3D volume.
- Extensive experiments validate our proposed method achieves **state-of-the-art** reconstruction results for 3D volumetric imaging for the task of ultra-sparse-view and limited-angle 3D CT reconstruction on different datasets, with improved inference time efficiency and minimal hyperparameter tuning.

5.2 Methods

Instead of modeling the 2D slices of the 3D volume as independent data samples during training time, and then applying regularization between slices at reconstruction time, we

propose incorporating information from neighboring slices at training time to enforce consistency between slices. More precisely, our first approach models the data distribution of a 3D volume with H slices in the z dimension as follows:

$$p(\mathbf{x}) = \prod_{i=1}^H p(\mathbf{x}[:, :, i] | \mathbf{x}[:, :, i-j:i-1], \mathbf{x}[:, :, i+1:i+j])/Z, \quad (5.1)$$

where j is a positive integer indicating the number of neighboring slices above and below the target slice that are being used as conditions to predict the target slice, and Z is a normalizing constant. To deal with boundary conditions where the third index may exceed the bounds of the original volume, we apply repetition padding above and below the main volume.

For training, we simply concatenate each of the conditioned slices with the target slice along the channel dimension to serve as an input to the neural network. Then we apply denoising score matching to predict the noise of the target slice as the loss function of the neural network:

$$\mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\mathbf{y} \sim \mathcal{N}(\mathbf{x}, \sigma_t^2 I)} \mathbb{E}_{i \in [1, H]} \left\| \frac{D_\theta(\mathbf{y}[:, :, i-j:i+j], \sigma_t) - \mathbf{x}[:, :, i]}{\sigma_t^2} \right\|_2^2. \quad (5.2)$$

At reconstruction time, the score function of the entire volume decomposes as a sum of score functions of each of the slices:

$$\nabla \log p(\mathbf{x}) = \sum_{i=1}^H \nabla \log p(\mathbf{x}[:, :, i] | \mathbf{x}[:, :, i-j:i-1], \mathbf{x}[:, :, i+1:i+j]). \quad (5.3)$$

In this way, we have rewritten the score of the 3D volume as sums of the scores of the 2D slices learned by the network. This means that we can now apply any algorithm that uses diffusion models to solve inverse problems to solve the 3D CT reconstruction problem. Furthermore, this method of blending together information from different slices allows us to learn a prior for the entire volume that combines information from different slices. We call this method **DiffusionBlend**.

To learn an even better 3D image prior, instead of learning the conditional distribution of individual target slices, we can learn the **joint distribution** of several neighboring slices at once, which we call a 3D patch. Letting k be the number of slices in each patch, we can partition the volume into 3D patches and approximate the distribution of the volume as

$$p(\mathbf{x}) = \left(\prod_{i=1}^{H/k} p(\mathbf{x}[:, :, (i-1)k+1:i*k]) \right) / Z, \quad (5.4)$$

where Z is a normalizing constant. Comparing this with (5.1), the main difference is instead of conditioning on neighboring slices, we are now incorporating the neighboring slices as a joint distribution. This allows for much faster reconstruction, as k slices are updated simultaneously according to their score function. However, this method faces similar slice consistency issues as in [42], since certain pairs of adjacent slices (namely, pairs whose slice indices are congruent to 0 and 1 modulo k) are never updated simultaneously by the network.

To deal with this issue, we propose two additional changes. Firstly, instead of using the same partition (updating the same k slices) at once for each iteration, we can use a different partition so that the previous border slices can be included in another partition. For example, we can randomly sample the end index of the first 3D patch for adjacency slices. Let m be uniformly sampled from $1, 2, \dots, k$, we can use the partition

$$\mathcal{S} = \{1, 2, \dots, H\} = \{1, \dots, m\} \cup \{m+1, \dots, m+1+k\} \cup \dots \cup \{H-k+1, \dots, H\}, \quad (5.5)$$

instead of $\mathcal{S} = \{1, 2, \dots, H\} = \{1, \dots, k\} \cup \{k+1, \dots, 2k\} \cup \dots \cup \{H-k+1, \dots, H\}$, where m is the offset index number in the new partition. We can then compute the score on the new partition. More generally, we can choose an arbitrary partition of \mathcal{S} into H/k sets, each containing k elements for each iteration, updating each slice in the small set simultaneously for that iteration.

Secondly, to better capture information between nonadjacent slices, we apply relative positional encoding as an input to the network. More precisely, if a 3D patch has a slice thickness (the distance between two slices) of p , then we let p be input of the positional encoding for that 3D patch. The positional encoding block consists of a sinusoidal encoding module and several dense connection modules, which has the same architecture as the timestamp embedding module of the same diffusion model. In this manner, the network is able to learn how to incorporate information from nonadjacent slices and captures more global information about the entire volume. Recall that for 3D patches of adjacent slices, the border between patches may have inconsistencies. To address this, we can **concatenate each border** as a new 3D patch, and then compute the score from it. If there are k slices in an adjacency-slice 3D patch, then the new 3D patch has the relative positional encoding of k , and also has a size of k . For instance, if the previous partition is $(1,2,3),(4,5,6),(7,8,9)$, the new partition is $(1,4,7),(2,5,8),(3,6,9)$. Here we are forming a new partition with jumping slices. In practice, since we need a pretrained natural image checkpoint due to scarcity of medical image data, we set $k = 3$ for facilitating fine tuning from natural image checkpoints.

We call the partitioning by 3D patch with adjacent slices as **Adjacency Partition**, and the partitioning by 3D patch with jumping slices as **Cross Partition**. Letting $r = H/k$ be the number of 3D patches, with a random partition, this method is stochastically averaging the different estimations of the $\nabla \log p(\mathbf{x})$ by different partitions. Specifically, the estimation of score by a single partition $\mathcal{S}_1 \cup \dots \cup \mathcal{S}_r$ is given by $\sum_{i=1}^r \nabla \log p(\mathbf{x}[:, :, \mathcal{S}_i])$. Ideally, we want to compute

$$|\mathcal{S}|^{-1} \sum_{\mathcal{S}=\mathcal{S}_1 \cup \dots \cup \mathcal{S}_r} \sum_{i=1}^r \nabla \log p(\mathbf{x}[:, :, \mathcal{S}_i]). \quad (5.6)$$

Similar to [42, 39, 109], we can share the summation in (5.6) across different diffusion steps since the difference between two adjacent iterations \mathbf{x}_i and \mathbf{x}_{i+1} is minimal.

In summary, we have shown how the score function of the entire volume can be written in terms of scores of the slices of the volume. Hence, similar to DiffusionBlend, this method can be coupled with any inverse problem solving algorithm. The scores of the slices can be approximated using a neural network. Training this network consists of randomly selecting k slices from a volume and concatenating them along the channel dimension to get the input to the network (along with the positional encoding of the slices), and then using denoising score matching as in (5.2) as the loss function.

Sampling and reconstruction. With Eq. 5.6, each reconstruction step would require computing the score functions corresponding to each of the partitions of \mathcal{S} , and then summing them to get the score function $s(\mathbf{x})$. We propose the variable sharing technique for this method, and only need to compute the score of one partition per time step. Hence, each iteration, we instead randomly choose one of the partitions of \mathcal{S} and update the volume of intermediate samples by the score function. Finally, we use repetition padding if H is not a multiple of k . This method incorporates a similar slice blending strategy as DiffusionBlend, but allows for significant acceleration at reconstruction time as k slices are updated at once. Furthermore, it allows the network to learn joint information between slices that are farther apart without requiring the increase in computational cost associated with increasing k . We call this method **DiffusionBlend++**. The pseudocode of the algorithm can be found in Alg. 3.

In practice, we choose not to select from all possible partitions, but instead select from those where the indices in each \mathcal{S}_i are not too far apart, as the joint information between slices that are very far apart is hard to capture. Table A.11 summarizes the different 3D image prior models. The appendix provides more details about the partition selection scheme.

Krylov subspace methods. Following the work of [41], we apply Krylov subspace methods to enforce data consistency with the measurement. At each timestep t , by using Tweedie’s formula [56], we compute $\hat{\mathbf{x}}_t = \mathbb{E}[\mathbf{x}_0|\mathbf{x}_t]$, and then apply the conjugate gradient method

$$\hat{\mathbf{x}}'_t = \text{CG}(\mathbf{A}^* \mathbf{A}, \mathbf{A}^* \mathbf{y}, \hat{\mathbf{x}}_t, M), \quad (5.7)$$

where in practice, the CG operator involves running M CG steps for the normal equation $\mathbf{A}^* \mathbf{y} = \mathbf{A}^* \mathbf{A} \mathbf{x}$. We combine this method with the DDIM sampling algorithm [165] to decrease reconstruction time. To summarize, we provide the algorithm for DiffusionBlend++ below. The Appendix provides the training algorithms for our proposed method as well as the reconstruction algorithm for DiffusionBlend.

Algorithm 3 DiffusionBlend++

Require: Forward model \mathbf{A} , sinogram \mathbf{y} , hyperparameter k

Initialize $\mathbf{x}_T \sim \mathcal{N}(0, \sigma_T^2 \mathbf{I})$

for $t = T : 1$ **do**

 Randomly select a partition $\mathcal{S} = \mathcal{S}_1 \cup \dots \cup \mathcal{S}_r$

if $t \bmod k = 0$ **then** use cross partition, **else** use random adjacency partitions

 Compute the relative positional encoding PE_t

 For each i , compute $\epsilon_\theta(\mathbf{x}_t[:, :, \mathcal{S}_i], PE_t)$

 Compute $\mathbf{s} = \nabla \log p(\mathbf{x}_t)$ using (5.6)

 Compute $\hat{\mathbf{x}}_t = \mathbb{E}[\mathbf{x}_0|\mathbf{x}_t]$ using Tweedie’s formula

 Set $\hat{\mathbf{x}}'_t = \text{CG}(\mathbf{A}^* \mathbf{A}, \mathbf{A}^* \mathbf{y}, \hat{\mathbf{x}}_t)$

 Sample \mathbf{x}_{t-1} using $\hat{\mathbf{x}}'_t$ and \mathbf{s} via DDIM sampling

end for

Return \mathbf{x}

5.3 Experiments

Experimental setup. We used the public CT dataset from the AAPM 2016 CT challenge [135] that consists of 10 volumes. We rescaled the images in the XY-plane to have size 256×256 without altering the data in the Z-direction and used 9 of the volumes for training data and the tenth volume as test data. The training data consisted of approximately 5000 2D slices and the test volume had 500 slices. We also performed experiments on the LIDC-IDRI dataset [6]. For this dataset, we first applied data preprocessing by setting the entire background of the volumes to zero. We rescaled the images in the XY-plane to have size 256×256 , and, to compare with the TPDM method, only took the volumes with at least

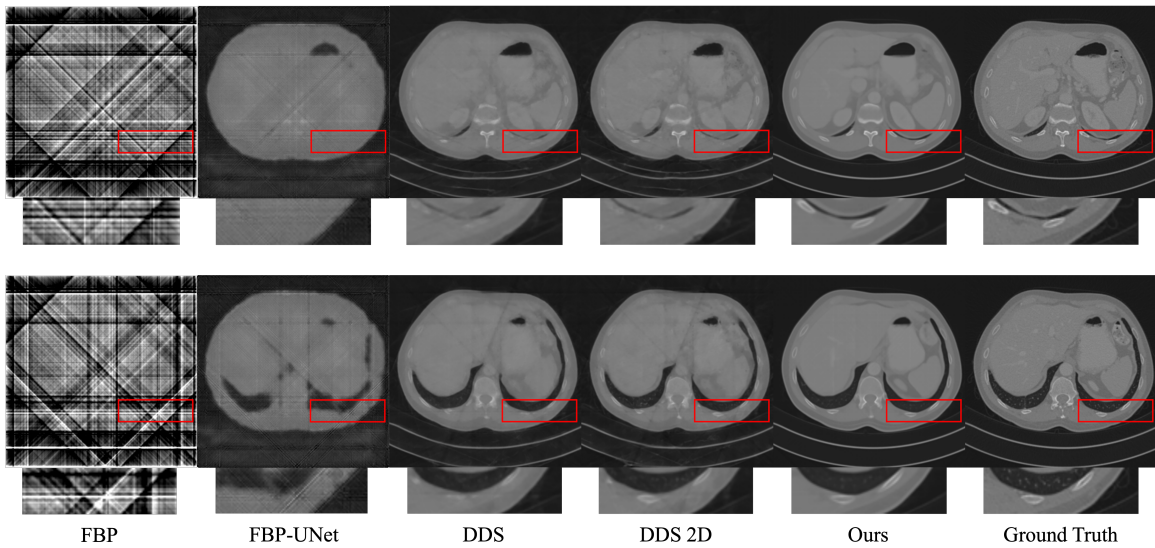


FIG 5.3 – Results of CT reconstruction with 4 views on AAPM dataset, axial view.

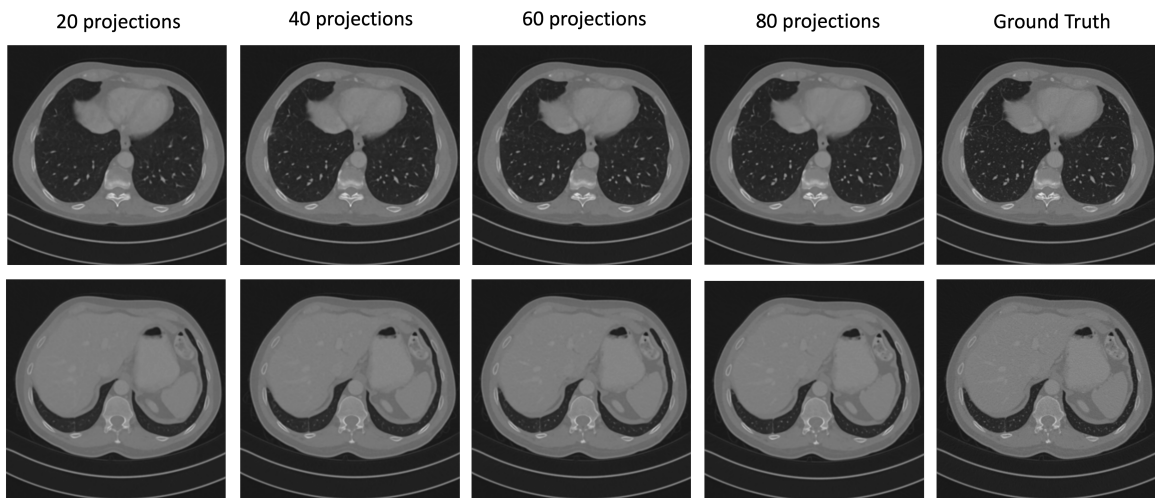


FIG 5.4 – Results of DiffusionBlend++ reconstruction with multiple views on AAPM dataset, axial view.

256 slices in the Z-direction, truncating the Z-direction to have exactly 256 slices. This resulted in 357 volumes which we used for training and one volume used for testing.

We performed experiments for sparse view CT (SVCT) and limited angle CT (LACT). The detector size was set to 512 pixels for all cases. For SVCT, we ran experiments on 4, 6, and 8 views. We also ran additional experiments on 20, 40, 60, 80, and 100 views and report the quantitative results in the Appendix. For LACT, we used the full set of views but only spaced around a 90 degree angle. In all cases, implementations of the forward and back projectors can be found in [42].

For a fair comparison between DiffusionBlend and DiffusionBlend++, we selected $j = 1$ for DiffusionBlend and each \mathcal{S}_i to contain 3 elements for DiffusionBlend++. In this manner, both methods involve learning a prior that involves products of joint distributions on 3 slices. To train the score function for DiffusionBlend, we started from scratch using the LIDC dataset. Since this dataset consisted of over 90000 slices, the network was able to properly learn this prior. We then fine tuned this network on the much smaller AAPM dataset. For DiffusionBlend++, the input and output images both had 3 channels from stacking the slices, so we fine-tuned the existing checkpoint from [73]. All networks were trained on PyTorch using the Adam optimizer with A40 GPUs. For reconstruction, we used 200 neural function evaluations (NFEs) for all the results. The appendix provides the full experiment hyperparameters. We observe that DiffusionBlend++ can reconstruct very high quality images that are free of artifacts as demonstrated in Fig.5.4 and Fig.5.3.

Comparison methods. We compared our proposed method with baseline methods for CT reconstruction and state of the art 3D diffusion model methods. We used the filtered back projection implementation found in [42]. For the other baseline, we used FBP-UNet [89] which is a supervised method that involves training a network for each specific task mapping the FBP reconstruction to the clean image. Since this is a 2D method, we learned a mapping between 2D slices and then stacked the 2D slices to get the final 3D volume. We also compared with classical CT reconstruction techniques such as SBT, SIRT, and CGLS [65] to benchmark our algorithm against traditional methods. Results for these methods are reported in the Appendix. For DiffusionMBIR [44], we fine-tuned the score function checkpoints on our data and used the same hyperparameters as the original work. We did the same for TPDM [109]; however, we ran TPDM only on the LIDC dataset because TPDM requires cubic volumes. Finally, we ran two variants of DDS [41]: one in which all the hyperparameters were left unchanged (DDS), and another in which no TV regularizer between slices was enforced (DDS 2D). Both of these methods were run with 200 NFEs. The appendix provides the experiment parameters.

Method	Sparse-View CT Reconstruction on AAPM						Sparse-View CT Reconstruction on LIDC					
	8 views		6 views		4 views		8 Views		6 Views		4 Views	
	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow
FBP	14.66	0.359	13.65	0.293	11.94	0.222	14.79	0.217	14.11	0.191	13.18	0.169
FBP-UNet	26.00	0.849	24.15	0.782	23.37	0.761	28.58	0.848	26.48	0.781	25.19	0.731
DiffusionMBIR	26.30	0.863	24.99	0.827	23.66	0.789	32.67	0.922	<u>31.18</u>	<u>0.901</u>	<u>29.02</u>	<u>0.863</u>
TPDM	-	-	-	-	-	-	27.51	0.816	25.60	0.776	21.99	0.695
DDS 2D	32.89	0.946	31.40	0.934	28.77	0.906	30.82	0.897	29.38	0.867	27.54	0.826
DDS	33.19	0.945	31.94	0.942	29.22	0.916	31.65	0.915	30.12	0.888	27.20	0.808
DiffusionBlend (Ours)	<u>34.29</u>	<u>0.955</u>	<u>33.26</u>	<u>0.949</u>	<u>31.84</u>	<u>0.944</u>	<u>33.34</u>	<u>0.933</u>	30.94	0.905	27.96	0.849
DiffusionBlend++ (Ours)	35.69	0.966	34.68	0.960	32.93	0.952	34.46	0.947	33.03	0.932	30.98	0.912

TBL 5.1 – Comprehensive comparison of quantitative results on Sparse-View CT Reconstruction on Axial View for AAPM and LIDC datasets. Best results are in bold.

Method	Sparse-View CT Reconstruction on AAPM						Sparse-View CT Reconstruction on LIDC					
	8 views		6 views		4 views		8 Views		6 Views		4 Views	
	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow
FBP	12.30	0.345	10.14	0.277	6.78	0.204	14.88	0.234	14.30	0.207	13.43	0.187
FBP-UNet	26.13	0.860	24.14	0.798	23.47	0.779	28.56	0.848	26.52	0.783	25.29	0.732
DiffusionMBIR	26.64	0.869	25.08	0.834	23.71	0.789	32.79	0.922	<u>31.30</u>	<u>0.900</u>	<u>28.98</u>	<u>0.862</u>
TPDM	-	-	-	-	-	-	27.66	0.819	25.57	0.784	21.87	0.708
DDS 2D	33.22	0.949	31.69	0.937	29.39	0.909	30.98	0.894	29.40	0.862	27.54	0.819
DDS	33.43	0.945	32.18	0.947	29.86	0.924	31.80	0.915	30.13	0.889	27.26	0.818
DiffusionBlend (Ours)	<u>35.09</u>	<u>0.958</u>	<u>33.97</u>	<u>0.952</u>	<u>32.38</u>	<u>0.943</u>	<u>33.73</u>	<u>0.934</u>	31.16	0.907	27.93	0.855
DiffusionBlend++ (Ours)	36.48	0.968	35.38	0.963	33.22	0.954	34.86	0.946	33.20	0.932	30.97	0.913

TBL 5.2 – Comprehensive comparison of quantitative results on Sparse-View CT Reconstruction on Sagittal View for AAPM and LIDC datasets. Best results are in bold.

Sparse-view CT. The results for different numbers of views and across different slices are shown in Tables 5.1, 5.2, and 5.3. DiffusionBlend++ exhibits much better performance over all the previous baseline methods (usually by a few dB) and outperforms DiffusionBlend. The second best method for each experiment is underlined and was, in most cases, DiffusionBlend. The exceptions are when the second best method is DiffusionMBIR, but this method was run with 2000 NFEs and took about 20 hours to run compared to 1-2 hours for both of our methods. The two DDS methods required similar runtime as our methods but in all cases exhibited inferior reconstruction results. Furthermore, DDS 2D generally performed worse than DDS. Thus, DDS failed to properly learn a 3D volume prior and still relied on the TV regularizer. Additionally, although TPDM should learn a 3D prior, the results were very poor compared to the other baselines. Our proposed

Method	Sparse-View CT Reconstruction on AAPM						Sparse-View CT Reconstruction on LIDC					
	8 views		6 views		4 views		8 Views		6 Views		4 Views	
	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow
FBP	14.64	0.325	13.18	0.268	11.16	0.236	14.78	0.206	14.10	0.181	13.10	0.165
FBP-UNet	27.34	0.878	25.12	0.827	24.10	0.810	28.87	0.858	26.59	0.793	25.37	0.744
DiffusionMBIR	29.86	0.908	28.12	0.875	25.68	0.843	33.29	0.922	<u>31.69</u>	<u>0.903</u>	<u>29.21</u>	<u>0.868</u>
TPDM	-	-	-	-	-	-	28.12	0.833	25.78	0.804	22.29	0.735
DDS 2D	33.64	0.950	32.33	0.939	30.25	0.916	31.60	0.898	29.99	0.871	28.03	0.830
DDS	33.97	0.934	32.95	0.930	30.89	0.932	32.51	0.920	30.83	0.898	27.61	0.828
DiffusionBlend (Ours)	<u>36.45</u>	<u>0.958</u>	<u>35.23</u>	<u>0.952</u>	<u>33.98</u>	<u>0.944</u>	<u>34.47</u>	<u>0.934</u>	31.48	0.908	28.24	0.859
DiffusionBlend++ (Ours)	37.87	0.968	36.66	0.963	34.27	0.955	35.66	0.947	33.97	0.935	31.38	0.913

TBL 5.3 – Comprehensive comparison of quantitative results on Sparse-View CT Reconstruction on Coronal View for AAPM and LIDC datasets. Best results are in bold.

Method	AAPM Dataset						LIDC Dataset					
	Axial		Sagittal		Coronal		Axial		Sagittal		Coronal	
	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow
FBP	16.36	0.643	16.36	0.524	15.62	0.531	18.79	0.672	19.84	0.675	20.01	0.676
FBP-UNet	27.38	0.910	27.81	0.918	28.44	0.930	29.42	0.885	29.50	0.884	29.54	0.887
DiffusionMBIR	25.98	0.872	27.14	0.877	27.74	0.903	<u>30.52</u>	0.906	30.57	0.906	30.68	0.907
TPDM	-	-	-	-	-	-	14.44	0.141	14.06	0.141	14.54	0.313
DDS 2D	28.05	0.916	27.99	0.916	28.82	0.922	27.92	0.843	27.89	0.835	27.96	0.842
DDS	28.20	0.918	28.17	0.926	29.03	0.934	28.12	0.865	28.06	0.869	28.13	0.879
DiffusionBlend (Ours)	<u>35.38</u>	<u>0.971</u>	<u>35.85</u>	<u>0.972</u>	37.62	<u>0.972</u>	30.43	<u>0.917</u>	<u>31.24</u>	<u>0.920</u>	<u>31.02</u>	<u>0.924</u>
DiffusionBlend++ (Ours)	35.86	0.975	36.03	0.976	<u>37.45</u>	0.976	34.33	0.957	34.48	0.957	34.64	0.956

TBL 5.4 – Comprehensive comparison of quantitative results on Limited-Angle CT Reconstruction on All Views for AAPM and LIDC datasets. Best results are in bold.

method learned a fully 3D prior and achieved the best results in the sagittal and coronal views.

Limited-angle CT. Table 5.4 shows all results for limited angle CT reconstruction for both the AAPM and LIDC datasets. Our DiffusionBlend++ method obtains superior performance over all the baseline methods and DiffusionBlend obtains the second best results. Similar to the SVCT experiments, DiffusionMBIR performed the best out of the baseline methods, but took approximately 40 hours to run. FBP-UNet performed reasonably well, but is a supervised method where the network must be retrained for each specific task. DDS is the most directly comparable to our method in runtime and methodology, but performed much worse quantitatively.

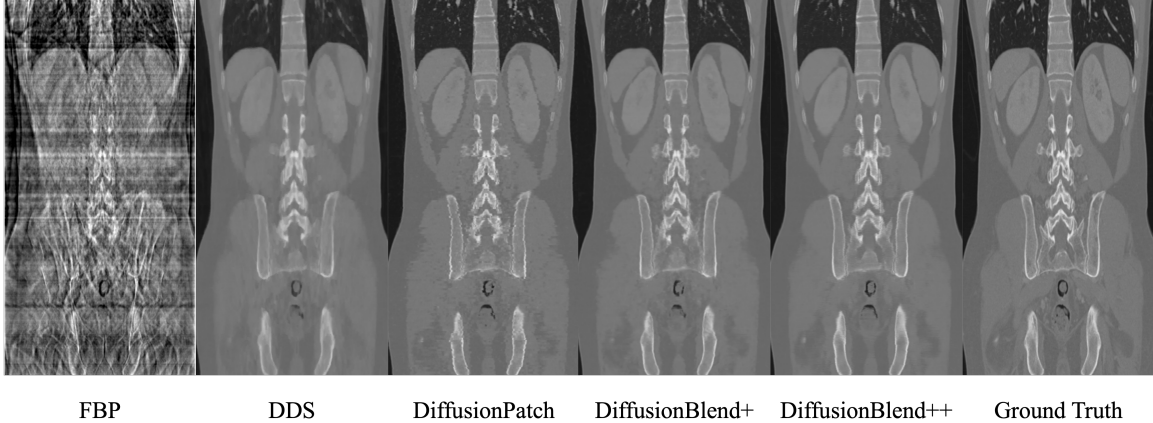


FIG 5.5 – Results of CT reconstruction with 8 views on AAPM dataset, coronal view. DiffusionPatch refers to Algorithm 1 with the same partition for every timestep, and DiffusionBlend+ refers to Algorithm 1 only with partitions of adjacency slices.

Algorithm	TV value	Difference with gt
DDS 2D	0.0104	0.0044
DDS	0.0031	-0.0034
DiffusionBlend++ (Ours)	0.0043	-0.0022
Ground Truth	0.0065	-

TBL 5.5 – TV values of different reconstruction algorithms on the AAPM test set

Inter-slice smoothness We demonstrate that DiffusionBlend++ learns the 3D prior internally, and achieves consistency and smoothness between 2D axial-plane slices without any external regularizations. In Table 5.5, we present the total variation (TV) value of the reconstructed images of different reconstruction algorithms on the test set of AAPM dataset, given by $\frac{1}{C \times W \times H} \|\mathbf{D}_z(x)\|_1$, where x is the image, \mathbf{D}_z is the total variation operator in z direction, and C , W , H are number of channels, width, and height. We find that both DiffusionBlend++ and DDS have TV less than the ground truth image, which implies that the reconstructed images are smooth in the z direction. However, we observe that DDS over-smooths the images as demonstrated in Fig. 5.5, which is represented by a much lower TV value than the ground truth. On the other hand, DiffusionBlend++ has smoothness level close to the ground truth without sacrificing sharpness of images.

Effectiveness of adjacency-slice blending and cross-slice blending We demonstrate that both the adjacency-slice blending and the cross-slice blending module are instrumental to a better reconstruction quality. Table 5.6 demonstrates the effectiveness of adding blending modules to the reverse sampling. Given the pretrained diffusion prior over slice

Adjacency	Cross	PSNR \uparrow	SSIM \uparrow
		34.85	0.954
✓		36.02	0.965
✓	✓	36.48	0.968

TBL 5.6 – Effectiveness of Blending Modules, Sagittal view performance on AAPM

patches, we observe that adding the adjacency-slice blending module improves the PSNR over a fixed partition by 1.17dB, and adding an additional cross-slice blending module further improves the PSNR by 1.63dB. Fig. 5.5 demonstrates that adding the cross-slice blending module removes artifacts and recovers sharper edges.

Ablation Studies We investigated the performance gain due to individual components. Details can be found in Appendix A.3.

5.4 Conclusion

In this work, we proposed two methods of using score-based diffusion models to learn priors of three dimensional volumes and used them to perform CT reconstruction. In both cases, we learn the distributions of multiple slices of a volume at once and blend the distributions together at inference time. Extensive experiments showed that our method substantially outperformed existing methods for 3D CT reconstruction both quantitatively and qualitatively in the sparse view and limited angle settings. In the future, more work could be done on other 3D inverse problems and acceleration through latent diffusion models. Image reconstruction methods like those proposed in this chapter have the potential to benefit society by reducing X-ray dose in CT scans.

CHAPTER 6

Local Patches Meet Global Context: Scalable 3D Diffusion Priors for Computed Tomography Reconstruction

6.1 Motivation

Diffusion models have emerged as a powerful class of generative model for image generation [168, 140, 144, 150], modeling the underlying image distribution $p(\mathbf{x})$ through learned denoising processes across progressively varied noise levels. The trained network approximates the score function of image distribution $\mathbf{s}(\mathbf{x}) = \nabla \log p(\mathbf{x})$. At inference time, starting from the initialization of Gaussian noise image, the noise signal is iteratively removed using the trained denoiser network, ultimately resulting in clean image samples drawn from the learned distribution prior $p(\mathbf{x})$.

Inverse problems can be formulated as $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{n}$, where \mathbf{y} is a known measurement, \mathcal{A} is the forward operator, \mathbf{x} is the unknown image to be recovered, and \mathbf{n} is random noise. Specifically, computed tomography (CT) reconstruction can also be formulated as an inverse problem that aims to reconstruct a 3D volume \mathbf{x} from the measurement \mathbf{y} , containing projection views of \mathbf{x} acquired from different angles [59]. Furthermore, to minimize patient radiation exposure, it is desirable to reduce the number of acquired projection views, typically on the order of 100-1000 views, to substantially fewer [162, 42], which is called *sparse-view CT reconstruction*. In this case, the inverse problem becomes ill-posed, as there are many potential images \mathbf{x} that could correspond to a measurement \mathbf{y} . Hence, it is necessary to impose prior knowledge on \mathbf{x} to obtain a reasonable reconstruction, such as [101, 104] in previous works.

This chapter based on [199]. Equal contribution with Taewon Yang.

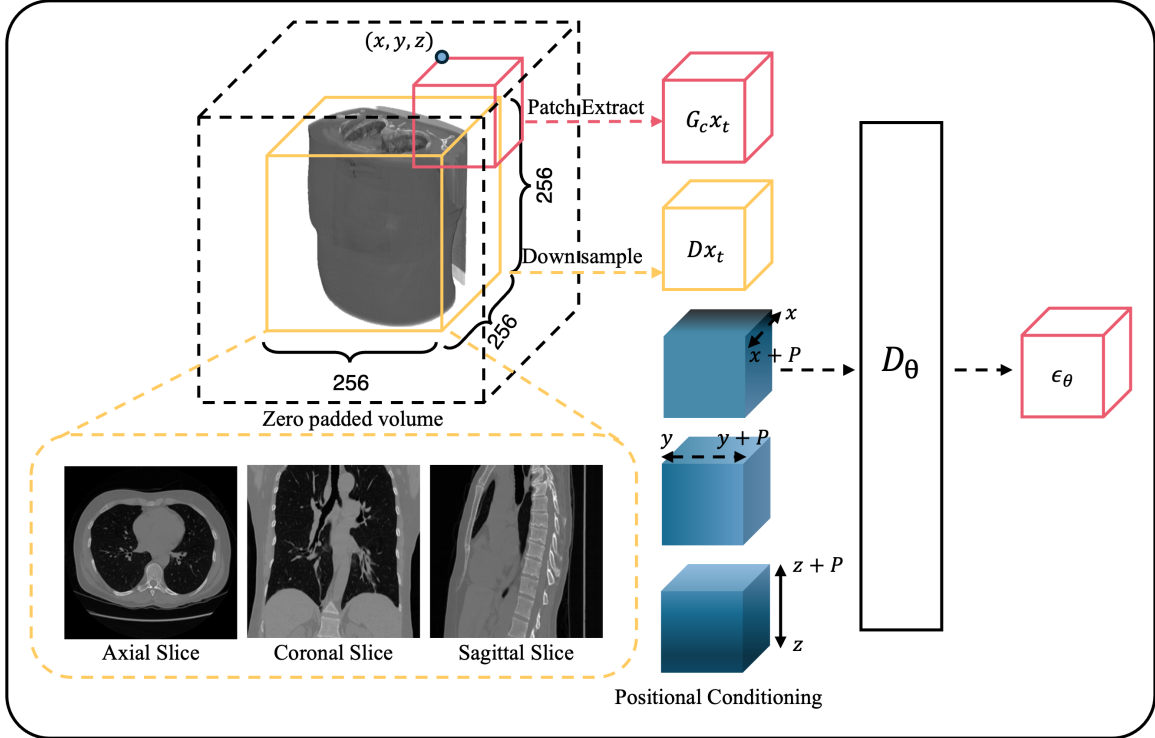


FIG 6.1 – Training the 3D patch diffusion model. Noisy 3D patch $G_c x_t$ (local details), downsampled volume $D x_t$ (global context), and positional encoding are concatenated as inputs to the denoiser network D_θ in each training iteration, to predict the noise ϵ_θ in the input noisy 3D patch. Positional encoding consists of voxel-based coordinates in x, y, z that are normalized to $[-1, 1]$.

Given that diffusion models can learn powerful generative priors over image data, it is natural to leverage these priors to address sparse-view CT reconstruction problem. Previous works [39, 95, 181, 115, 45] have proposed various approaches leveraging diffusion models to solve different inverse problems, including 2D CT reconstruction. These approaches generally consist of, at inference time, interweaving the original diffusion model steps with data fidelity enforcement steps that push the image to be consistent with the measurements.

However, challenges emerge when directly applying these existing approaches to real-world sparse-view CT tasks, where the images are inherently three-dimensional volumes. Standard clinical CT images have the full size of $512 \times 512 \times 512$, making it almost infeasible to fit the entire image into the network training due to memory constraints. Furthermore, training diffusion models typically requires large-scale datasets [203], often comprising hundreds of thousands of images such as ImageNet [153] or LSUN [200]. For CT reconstruction, there is a lack of high-quality volumetric CT data in such scale. Finally,

training time of diffusion models on 3D images would be prohibitively expensive. Thus, the high computational costs in memory, time, and data pose significant challenges to directly train diffusion models on 3D CT images.

To address these challenges, prior works have proposed leveraging 2D image priors, either by combining them with hand-crafted external regularization terms [42] or by training an additional 2D image prior along a perpendicular plane [109]. However, these approaches tend to produce blurry artifacts due to sensitivity to the hand-crafted regularization term, or require long inference times due to the use of two perpendicular 2D CT priors. Recently, [162] introduced a new solution to learn the joint distribution over a small set of 2D slices simultaneously, but the underlying network remains a 2D architecture with multiple channels, needing significant reconstruction time comparable to that of [41]. However, these methods do not fully realize and exploit the generative capacity of diffusion models for high-resolution 3D volumetric images. Recent work by Teng et al. [174] trained a 3D convolutional UNet with $256 \times 256 \times 30$ slabs and used DPS to reconstruct $256 \times 256 \times 120$ CT volumes. No 3D image generation was shown.

To overcome these limitations, we propose a novel approach that learns scalable 3D diffusion priors by representing 3D volumes as **position-encoded local patches**. To ensure high-quality 3D image generation, our study emphasizes the importance of incorporating global context. Specifically, by jointly modeling the distribution of local patches alongside a downsampled version of the entire volume, our method effectively unifies global contextual information with local image statistics, to obtain improved **global-aware patch priors**. Beyond image generation, the resulting fully 3D prior demonstrates superior performance in large-scale inverse problem solving, such as high-resolution 3D CT reconstruction. Overall, our contributions are three-fold.

- We propose a novel 3D patch-based diffusion model that learns a fully 3D diffusion prior, enabling the scalable generation of high-resolution 3D images, and offering an unprecedentedly efficient and accurate solution to address high-resolution 3D inverse problems.
- Our key technical contribution is a framework that models the joint distribution of position-aware 3D local patches and a downsampled 3D volume, enabling scalable learning of a global-contextual 3D patch-based prior by effectively integrating local and global information to achieve high-quality 3D image generation.
- Extensive experiments on different datasets and different-sized images show that our proposed method achieves SOTA performance for 3D sparse-view CT reconstruction in various settings.

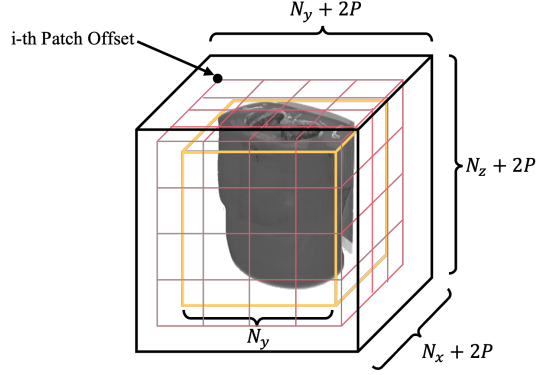


FIG 6.2 – Schematic illustration for zero padding and partitioning image into 3D patches. Each index i represents one of P^3 possible ways to choose a patch offset tuple.

6.2 Methods

6.2.1 Global-Aware Patch Prior

We begin by following the formulation of learning patch-based priors as described in [82], but generalized for 3D images. In patch-based diffusion models, one needs first to zero pad each $N_x \times N_y \times N_z$ image by an amount P on each side. Slightly recycling notation, we let \mathbf{x} denote the resulting padded image. We partition \mathbf{x} into many square patches, with one bordering region of zeros by first choosing the i th patch offset tuple $(o_1, o_2, o_3) \in \{0, \dots, P-1\}^3$ according to Figure 6.2. The number of patches needed in the x direction to perfectly tile the image is $k_1 = N_x/P$, and define $k_2 = N_y/P$ and $k_3 = N_z/P$ as the number of patches to perfectly tile the image in the y and z directions, respectively. Thus, the total number of patches needed to cover the image with an offset is $l = (k_1 + 1)(k_2 + 1)(k_3 + 1)$. Hence, our model for the data distribution has the form

$$p(\mathbf{x}) = \frac{1}{Z} \left(\prod_{i=1}^{P^3} p_{i,B}(\mathbf{x}_{i,B}) \prod_{r=1}^l p_{i,r}(\mathbf{x}_{i,r}) \right)^{1/P^3}, \quad (6.1)$$

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 $P \times P \times P$ patch when using the partitioning scheme corresponding to the i th patch offset, $p_{i,r}$ is the probability distribution of that region, and Z is a normalizing factor. This model takes a product over all possible patch offsets, eliminating boundary artifacts that would occur if a fixed patch offset was used.

Crucially, however, the model (6.1) does not account for global structure within the whole image or correlation between two distant patches. This can lead to reasonable local structure within each generated local patches, but a lack of coherent and realistic longer range details. At the same time, memory constraints in 3D preclude inputting the entire image into the network that approximates any part of the score function. Motivated by these considerations, we develop an improved **patch-based prior with global context** as follows:

$$p(\mathbf{x}) = \left(\prod_{i=1}^{P^3} p_{i,B}(\mathbf{x}_{i,B}) \prod_{r=1}^l p_{i,r}(\mathbf{x}_{i,r}, \mathbf{D}\mathbf{x}) \right)^{1/P^3} / Z, \quad (6.2)$$

so that now $p_{i,r}$ represents a joint distribution between the local patch $\mathbf{x}_{i,r}$ and the down-sampled image $\mathbf{D}\mathbf{x}$, with \mathbf{D} representing an operator that downsamples the whole image to a smaller image size. (Note that in practice, \mathbf{D} may require pixel interpolation if the whole image dimension is not a multiple of the downsampled image.) Then

$$\log p(\mathbf{x}) = \frac{1}{P^3} \left(\sum_{i=1}^{P^3} (\log p_{i,B}(\mathbf{x}_{i,B}) + \sum_c \log p_c(\mathbf{G}_c \mathbf{x}, \mathbf{D}\mathbf{x})) \right), \quad (6.3)$$

where for convenience we reindex the patches with c and \mathbf{G}_c denotes the same patch grabbing operator as aforementioned. To take the gradient of this model, first define $\mathbf{u} = \mathbf{G}_c \mathbf{x}$ and $\mathbf{v} = \mathbf{D}\mathbf{x}$. Then

$$\nabla \log p(\mathbf{x}) = \frac{1}{P^3} \left(\sum_{i=1}^{P^3} \sum_c (\mathbf{G}'_c \nabla_{\mathbf{u}} \log p_c(\mathbf{G}_c \mathbf{x}, \mathbf{D}\mathbf{x}) + \mathbf{D}' \nabla_{\mathbf{v}} \log p_c(\mathbf{G}_c \mathbf{x}, \mathbf{D}\mathbf{x})) \right), \quad (6.4)$$

where we dropped the gradient of the bordering term since it is known to be zero, and we applied the chain rule. Hence, we can approximate the score functions $\nabla_{\mathbf{u}} \log p_c(\mathbf{G}_c \mathbf{x}, \mathbf{D}\mathbf{x})$ and $\nabla_{\mathbf{v}} \log p_c(\mathbf{G}_c \mathbf{x}, \mathbf{D}\mathbf{x})$ with the neural networks \mathbf{s}_u and \mathbf{s}_v respectively. The denoising score matching loss is then

$$L(\boldsymbol{\theta}) = \frac{1}{P^3} \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma_t^2 \mathbf{I})} \left\| \sum_{i=1}^{P^3} \sum_c (\mathbf{G}'_c \mathbf{s}_u(\mathbf{u}, \mathbf{v}; \boldsymbol{\theta}) + \mathbf{D}' \mathbf{s}_v(\mathbf{u}, \mathbf{v}; \boldsymbol{\theta})) - P^3 \boldsymbol{\epsilon} / \sigma_t^2 \right\|_2^2. \quad (6.5)$$

Note that $\hat{\mathbf{x}} = \mathbf{x} + \boldsymbol{\epsilon}$ and now $\mathbf{u} = \mathbf{G}_c \hat{\mathbf{x}}$, $\mathbf{v} = \mathbf{D}\hat{\mathbf{x}}$, and \mathbf{x} is sampled from the clean image distribution but we dropped these terms from the expectation for simplicity.

Unfortunately, due to the double summation inside of the norm, it is computationally infeasible to train the network using this loss function. Hence, we first use the inequality

$$\|X_1 + \dots + X_n\|_2^2 \leq n \cdot (\|X_1\|_2^2 + \dots + \|X_n\|_2^2), \quad (6.6)$$

to upper bound the loss (6.5) without the expectation to get

$$L \leq C \sum_{i=1}^{P^3} \left\| \sum_c \mathbf{G}'_c \nabla_{\mathbf{u}} \log p_c(\mathbf{u}, \mathbf{v}) + \mathbf{D}' \nabla_{\mathbf{v}} \log p_c(\mathbf{u}, \mathbf{v}) - \epsilon / \sigma_t^2 \right\|_2^2, \quad (6.7)$$

for a constant $C > 0$. Therefore, during training we can choose a random i and then minimize the inner sum. For fixed i , we usually choose each of the patches c in a non-overlapping way. Unfortunately, the \mathbf{D}' operator upscales the $\nabla_{\mathbf{v}}$ patch-sized term into a whole image sized term, so the sum over c will contain contributions from each patch for every pixel. However, the part that is outside the patch c is difficult to learn as the only contribution is from the upscaled part involving \mathbf{D}' . Thus we instead focus on learning the patch-specific part inside of the loss. Rewriting the norm part of the loss, this is

$$L_c \approx \left\| \sum_c \mathbf{G}_c \mathbf{G}'_c \nabla_{\mathbf{u}} \log p_c(\mathbf{u}, \mathbf{v}) + \mathbf{G}_c \mathbf{D}' \nabla_{\mathbf{v}} \log p_c(\mathbf{u}, \mathbf{v}) - \mathbf{G}_c \epsilon / \sigma_t^2 \right\|_2^2. \quad (6.8)$$

When factoring out the \mathbf{G}_c term, we observe that each term in this sum independently contributes to one of the non-overlapping patches of the entire image. Hence, when training the network, we do not have to perform denoising score matching across the entire image: we need only add noise to the entire image, provide one of the patches and the downsampled image to the network, and have the network learn to denoise just the patch $\mathbf{G}_c \epsilon$. In practice, to represent the joint score functions of $\mathbf{G}_c \mathbf{x}$ and $\mathbf{D} \mathbf{x}$, we input $\mathbf{G}_c \mathbf{x}$ and $\mathbf{D} \mathbf{x}$ to the network via concatenation along the channel dimension (where \mathbf{D} is chosen appropriately such that $\mathbf{G}_c \mathbf{x}$ and $\mathbf{D} \mathbf{x}$ are the same size), and represent $\nabla_{\mathbf{u}} \log p_c(\mathbf{G}_c \mathbf{x}, \mathbf{D} \mathbf{x}) + \nabla_{\mathbf{v}} \log p_c(\mathbf{G}_c \mathbf{x}, \mathbf{D} \mathbf{x})$ using single channels of the network output.

To facilitate the network to learn different prior distributions of 3D local patches at different locations, it is desirable to incorporate the positional information of the 3D patches. We do this by defining the x positional array as a 3D array representing the x coordinates of each pixel in the image, and normalizing to be scaled between -1 and 1. Similarly, y and z positional arrays are defined for y and z coordinates of the pixels, respectively. To enable the network to learn patch distributions that vary based on location within the image, we extract patches from these positional arrays corresponding to the image patches.

As shown in Figure 6.1, these positional patches are also concatenated along the channel dimension with $G_c \mathbf{x}$ and $D\mathbf{x}$.

In summary, we use a 3D UNet with five input channels: the noisy patch, the downsampled (noisy) image, and the three positional arrays. The network has one output channel, used to approximate $\nabla_{\mathbf{u}} \log p_c(\mathbf{u}, \mathbf{v})$ and $\nabla_{\mathbf{v}} \log p_c(\mathbf{u}, \mathbf{v})$, and is trained according to the loss (6.8). Most importantly, the spatial dimensions of the input are only in the 3D patch size of $P \times P \times P$ as shown in Figure 6.1 Therefore, we never need to input the whole 3D image into the network. The appendix shows the pseudocode for the training algorithm.

6.2.2 Sampling and Reconstruction Algorithm

Diffusion models are inevitably slow, as they typically require hundreds or thousands of steps to sample one image. Prior works have explored accelerating the diffusion sampling stage [165, 93, 128], For example, to expedite the sampling process, DDIM [165] proposed the following update rule

$$\mathbf{x}_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \hat{\mathbf{x}}_t + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \hat{\epsilon}_t + \sigma_t \epsilon_t \quad (6.9)$$

where $\epsilon_t \sim \mathcal{N}(0, I)$, $\bar{\alpha}_{t-1}$ is the noise schedule at timestep $t - 1$, $\hat{\epsilon}_t$ is the predicted noise and $\hat{\mathbf{x}}_t$ is the denoised estimate

$$\hat{\mathbf{x}}_t = \frac{1}{\sqrt{\bar{\alpha}_t}} (\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \hat{\epsilon}_t). \quad (6.10)$$

In [165], the stochasticity is controlled by σ_t that is defined as

$$\sigma_t = \eta \sqrt{(1 - \bar{\alpha}_{t-1}) / (1 - \bar{\alpha}_t)} \sqrt{1 - \bar{\alpha}_t / \bar{\alpha}_{t-1}} \quad (6.11)$$

where $\eta \in [0, 1]$ so that $\eta = 0$ leads to fully deterministic sampling. Alternatively, [41] considers

$$\sigma_t = \eta \sqrt{1 - \bar{\alpha}_{t-1}}, \quad (6.12)$$

for which (6.9) becomes

$$\mathbf{x}_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \hat{\mathbf{x}}_t + \sqrt{1 - \bar{\alpha}_{t-1}} (\sqrt{1 - \eta^2} \cdot \hat{\epsilon}_t + \eta \epsilon_t). \quad (6.13)$$

Assuming the estimated $\hat{\epsilon}_t$ and the stochastic $\epsilon \sim \mathcal{N}(0, I)$ are independent Gaussians, $\tilde{\epsilon} = \sqrt{1 - \eta^2} \cdot \hat{\epsilon}_t + \eta \epsilon_t$ is also a Gaussian. Hence, $\tilde{\epsilon}$ can equivalently be represented as

a sample from $\tilde{\epsilon} \sim \mathcal{N}(0, \mathbf{I})$. This preservation of noise variance ensures that the model samples \mathbf{x}_{t-1} from the marginal distribution $q(\mathbf{x}_{t-1}|\mathbf{x}_0)$.

For the underlying sampling algorithm, we choose DDIM [165] as a fast and reliable sampler that many other works [181, 95, 47] have also used as the backbone for solving inverse problems. During the sampling step, each update uses (6.13).

The score function model (6.4) allows for estimation of the score of the whole image from the trained network. However, it requires averaging over P^3 terms in the outer loop, which would be prohibitively expensive. To address this issue, [82] proposed a method that, instead of taking the average, chooses one of the terms at random and uses that as the approximation of the score.

However, when the number of sampling steps is limited, this approach becomes inapplicable, introducing boundary artifacts between patches. Also, the sampling process on one initial point does not consider the joint distribution of pixels across different patches at the same time step. Here, we propose a novel recurrent noising strategy to further mitigate the boundary artifacts when reducing the sampling steps for improving sampling efficiency. Specifically, at each iteration of DDIM, we denoise and renoise the image for K times, where in each time a different randomly chosen patch location, and finally average the results over K times to obtain the sample. The sampling equation is

$$\mathbf{x}_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \hat{\mathbf{x}}_t^{avg} + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \cdot \hat{\epsilon}_t + \sigma_t \epsilon_t \quad (6.14)$$

where σ_t is equal to (6.12), $\hat{\mathbf{x}}_t^{avg}$ is the averages of $\{\hat{\mathbf{x}}_{i^w, t}\}_{w=1}^K$ and $\hat{\epsilon}_t$ is the scaled sum of $\{\epsilon_{i^w, t}\}_{w=1}^K$, computed over K randomly chosen patch locations i^w . Here, noise variance is preserved similar to (6.13), allowing this method to sample from $q(\mathbf{x}_{t-1}|\mathbf{x}_0)$.

Previous works [41, 162] showed that conjugate gradient method allows the model to apply data consistency with measurement much faster. At each timestep t , using the denoised estimate (6.10) we apply conjugate gradient method on the following forward operator \mathbf{A} :

$$\hat{\mathbf{x}}_t' = \text{CG}(\mathbf{A}^* \mathbf{A}, \mathbf{A}^* \mathbf{y}, \hat{\mathbf{x}}_t, M) \quad (6.15)$$

where the CG operator runs M CG steps for the equation $\mathbf{A}^* \mathbf{y} = \mathbf{A}^* \mathbf{A} \mathbf{x}$. We apply this data consistency term after K iterations and then apply DDIM sampling to solve the inverse problem. To summarize, our proposed sampling algorithm is shown in 4.

Algorithm 4 3D Patch Diffusion Sampling with Recurrent Noising

Require: Recurrent noising number K , \mathbf{A} , \mathbf{y} , P

```
1: Initialize  $\mathbf{x}_T \sim \mathcal{N}(0, \mathbf{I})$ 
2: for  $t = N : 1$  do
3:   for  $w = 1 : K$  do
4:     Randomly select integer  $i^w \in [1, P^3]$ 
5:     For  $1 \leq r \leq l$ , extract 3D patches  $\mathbf{x}_{i^w, r, t}$ 
6:     for  $r = 1 : l$  do
7:        $\hat{\mathbf{x}}_{i^w, r, t}, \boldsymbol{\epsilon}_{i^w, r, t} \leftarrow D_{\theta}(\mathbf{x}_{i^w, r, t}, \sigma_t)$ 
8:     end for
9:     Concatenate all  $\hat{\mathbf{x}}_{i^w, r, t}$  to entire volume  $\hat{\mathbf{x}}_{i^w, t}$ 
10:    Concatenate all  $\boldsymbol{\epsilon}_{i^w, r, t}$  to entire volume  $\boldsymbol{\epsilon}_{i^w, t}$ 
11:     $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I})$ 
12:     $\mathbf{x}_t \leftarrow \sqrt{\bar{\alpha}_t} \hat{\mathbf{x}}_{i^w, t} + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}$ 
13:  end for
14:   $\hat{\mathbf{x}}_t^{\text{avg}} = \sum_{w=1}^K \hat{\mathbf{x}}_{i^w, t} / K$ ,  $\hat{\boldsymbol{\epsilon}}_t = \sum_{w=1}^K \boldsymbol{\epsilon}_{i^w, t} / \sqrt{K}$ 
15:  Set  $\hat{\mathbf{x}}_t' = \text{CG}(\mathbf{A}^* \mathbf{A}, \mathbf{A}^* \mathbf{y}, \hat{\mathbf{x}}_t^{\text{avg}})$ 
16:  Sample  $x_{t-1}$  via DDIM using (6.9)
17: end for
```

6.3 Results

6.3.1 Experimental Setup

Dataset preparation. We conducted experiments using the LIDC-IDRI dataset [6]. The original dataset has a spatial resolution of 512×512 in the XY-plane, with varying sizes along the Z-axis. In our experiments, we selected the first 256 slices along the Z-axis and filtered out samples with inconsistent shapes. To ensure comparability with contemporary comparison methods, we initially rescaled each XY-plane to 256×256 . This preprocessing resulted in 90 volumes for training and 1 volume for testing. We also performed experiments on the AAPM dataset [135] that contains only 10 CT volumes, to evaluate our method’s performance with a limited amount of data. This dataset was preprocessed in the same manner as the LIDC-IDRI dataset. To demonstrate that our approach can be extended to larger volumes, we additionally prepared a version of the LIDC-IDRI dataset with clinically relevant dimensions of $512 \times 512 \times 256$.

Baseline setting. For the LIDC-IDRI dataset, the model was trained for approximately 8 days with a batch size of 64 and a patch size of $32 \times 32 \times 32$. For the AAPM dataset, which contains only 9 volumes, we trained the model from scratch for 5 days using the

same batch and patch sizes. For the larger LIDC-IDRI dataset, the model was trained for 28 days with a batch size of 16 and a patch size of $64 \times 64 \times 32$. All training was done on a single NVIDIA A100 GPU. We applied our proposed sampling method using $K = 2$ and $\eta = 0.8$ for both $256 \times 256 \times 256$ and $512 \times 512 \times 256$ volumes. The total number of model parameters is 68.59 million.

6.3.2 Unconditional 3D Image Generation

Our proposed method is able to learn a fully 3D prior on $256 \times 256 \times 256$ sized CT volumes, allowing it to generate high quality volumes. We used (6.13) with $K = 1$ and $\eta = 0.4$ and 200 steps for unconditional sampling. Figure 6.3 shows unconditional volume samples from LIDC-IDRI prior, visualized across axial, coronal, and sagittal slices. The top row shows slices from an unconditionally sampled volume, while the bottom row shows the nearest-neighbor volume from the training dataset. These results indicate that the proposed prior does not simply memorize from the training dataset and is capable of generating realistic high-resolution 3D CT volumes with fine-grained details of anatomic structure. Further generation results on different prior are in the appendix.

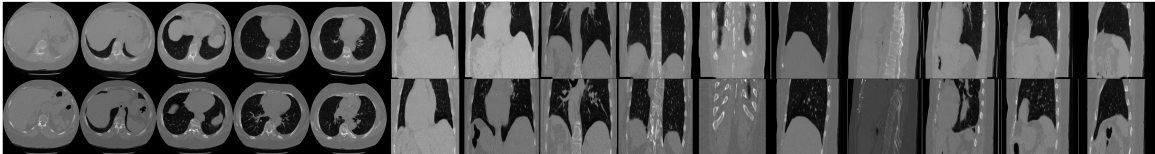


FIG 6.3 – Unconditional 3D image generation results using the LIDC-IDRI prior. The top row shows axial, coronal, and sagittal slices from a generated volume, and the bottom row shows the corresponding slices from its nearest-neighbor volume in the training dataset. The slice indices for the axial, coronal, and sagittal views are $[30, 80, 130, 180, 230]$, $[70, 100, 130, 160, 190]$, and $[60, 130, 160, 190, 210]$, respectively

6.3.3 Solving Inverse Problems

We conduct experiments of sparse-view CT (SVCT) reconstruction from 8, 20, and 60 views on one of the test volumes. The results in Table 6.1 and Figure 6.4 demonstrate that our method is able to reconstruct high-quality images. Furthermore, we compared our proposed method with various classical methods and diffusion model methods for 3D CT reconstruction. We used the filtered back projection method with a ramp filter, whose implementation is found in [42]. We also used the total variation (TV) regularizer and

TBL 6.1 – PSNR (dB) comparison of various methods for CT reconstruction on different datasets and resolutions with best results in bold.

Dataset / Method	LIDC ($256 \times 256 \times 256$)			LIDC ($512 \times 512 \times 256$)			AAPM ($256 \times 256 \times 256$)			
	Views	8	20	60	8	20	60	8	20	60
FBP		15.08	20.21	29.37	14.46	20.05	28.02	13.81	18.39	27.71
ADMM-TV		21.15	24.12	29.43	22.47	25.39	29.38	23.71	26.75	31.28
FBP-UNet [89]		22.60	27.77	32.78	26.09	30.32	35.34	25.93	30.06	36.93
DDS [41]		23.23	29.92	35.79	24.37	27.33	28.61	30.18	34.94	38.87
DiffusionBlend [162]		30.43	35.89	40.87	31.69	35.94	39.32	–	–	–
Blend + FT [162]		32.94	37.05	42.72	33.01	36.44	39.31	32.25	37.21	42.21
Proposed		33.06	38.56	43.70	33.17	37.33	40.16	32.96	38.47	42.66

solved the optimization problem using ADMM (ADMM-TV) with the implementation in [74]. For the other baseline, we implemented FBP-UNet [89] which is a supervised method that involves training a UNet that maps FBP reconstruction to the clean image. Since FBP-UNet is a 2D method, we learned a mapping between 2D slices and then stacked the 2D slices to get the final 3D volume.

For diffusion model methods, we ran DDS [41] and DiffusionBlend [162]. To obtain the best possible results for DiffusionBlend, we used the trick in the original work that involved taking a network that was pretrained on ImageNet and then finetuning it on the relevant CT dataset. Thus, we specifically ran DiffusionBlend++ with the prior trained on groups of 3 slices. To illustrate the effect of using the pretrained network, we show experiments both with the finetuning method (Blend+FT) and with networks that were trained from scratch on the CT datasets. Since our proposed method does not rely on a pretrained network, the most fair comparison is with [162] without a pretrained network; nevertheless we showed both results to illustrate how our proposed method can obtain superior results. Both of these methods were run with 200 sampling steps. The appendix provides the experiment parameters.

6.3.4 Ablation studies

We conducted ablation studies to analyze the factors that may influence the performance of our 3D patch diffusion model.

TBL 6.2 – Per volume runtimes of different methods for 8 view CT recon.

Method	Runtime (minutes) ↓
FBP	0.1
ADMM-TV	6
FBP-UNet	1
DDS	64
DiffusionBlend	55
Proposed	20

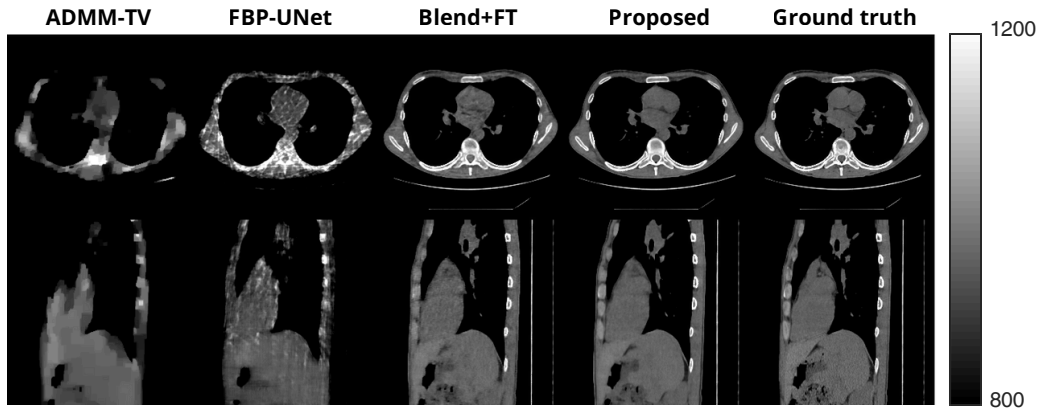


FIG 6.4 – Results of our proposed method and comparison methods for 20 view CT recon on LIDC $256 \times 256 \times 256$ dataset. Images are shown in modified Hounsfield units. The top row shows the axial slice and the bottom row shows the sagittal slice from the reconstructed volume.

Impact of downsample channel. Downsample channel plays a critical role in our method, allowing the model to learn the global structure. Without the downsample channel the model will only learn the local relationship between patches, failing to generate realistic CT images. To validate this, we compared the generation quality of the 3D patch-based diffusion model with and without the downsample channel using the Fréchet Inception Distance (FID), as shown in Fig. 6.3.

Impact of patch size. Consider an extreme case where the patch size is 1×2 . According to our modeled prior distribution in (6.1), using such small patches is theoretically valid but unlikely to fully capture image statistics. To better understand this behavior, we examined how patch size affects the model’s performance on inverse problem tasks. We trained

TBL 6.3 – FID comparison with and without downsample channel

Method	w. downsample	w.o. downsample
FID ↓	40.80	112.12

models with varying patch sizes, which is feasible due to the input-size-agnostic property of U-Net architectures. Patch sizes were chosen by dividing the full volume size by a scaling factor. For example, a scaling factor of 8 on a $256 \times 256 \times 256$ volume corresponds to patches of size $32 \times 32 \times 32$. Since the data are volumetric, doubling the patch size in each dimension reduced the batch size by a factor of 8. As a result, smaller patches train faster in wall-clock time: the factor-16 setting converges within 3 days, whereas the factor-8 setting requires approximately 8 days on $256 \times 256 \times 256$ volumes.

Experimental results indicate that as the patch size decreases, the model failed to capture meaningful semantic structure, leading to reduced PSNR on the sparse-view CT (SVCT) task. As shown in Table 6.4, model performance degrades as the scaling factor increases. This observation suggests that when patches become too small relative to the full image, they behave more like individual pixels. Consequently, the model can no longer learn a coherent image-level prior, causing (6.1) to describe a less effective model.

TBL 6.4 – PSNR in dB of different patch and volume size

Factor	$256 \times 256 \times 256$	$512 \times 512 \times 256$
8	33.06	33.17
16	27.87	31.71

Impact of recurrent noising during sampling process. This section explores how recurrent noising affects the reconstruction quality. To understand how the number of iterations K affects the reconstructed images, we tested on the 8-view SVCT task with 200 sampling steps. Table 6.5 shows that recurrent noising during sampling ($K > 1$) yielded higher PSNR than reconstruction without any recurrent noising ($K = 1$). The PSNR peaked at $K = 2$ and gradually decreased as K increased. We conjecture that this

TBL 6.5 – PSNR in dB of different K , with best results in bold

K	8 view	20 view	60 view
1	32.53	38.18	43.46
2	33.06	38.56	43.70
3	32.97	38.43	43.52
4	32.74	38.15	43.44

is because the estimated $\hat{\epsilon}_t$ is not strictly sampled from a Gaussian distribution, causing error accumulation as K increases.

6.4 Conclusion

In this chapter, we propose a novel 3D patch-based diffusion model that is conditioned with downsampled volume and coordinates to enable learning the 3D CT image prior in an efficient way. By learning from the local patches coupled with global context, our model is able to generate high-resolution 3D images. Experiments across different volume sizes and CT datasets showed that our method achieve SOTA results on SVCT reconstruction while reducing inference time by more than $2\times$. One of the limitations is that our approach assumes the data is well aligned to a consistent coordinate, which may reduce applicability to less structured image domains. In the future, we plan to investigate alternative methods that allow the model to learn the positional information of the patch, thereby improving robustness in less structured data.

CHAPTER 7

Test-Time Adaptation Improves Inverse Problem Solving with Patch-Based Diffusion Models

7.1 Motivation

In imaging, inverse problems are important and consist of reconstructing an image \mathbf{x} from a measurement $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \epsilon$. 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 competitive results can be obtained when using them to solve inverse problems [43, 39, 162, 181, 96, 115].

However, these diffusion inverse solvers (DIS) require well-trained diffusion models, which in turn require large amounts of clean training data [168, 73]. It may be infeasible to collect large training data sets in many applications such as medical imaging [43, 167, 86], black hole imaging [61, 60], and phase retrieval [115, 188]. For example, in practical applications such as dynamic CT reconstruction [147] and single photon emission CT [118], obtaining high quality measurements, which can lead to reconstructions closely approximating the ground truth, can be slow or potentially harmful to patients, so only very small datasets of clean images are available. Moreover, for very challenging inverse problems such as black hole imaging [61] and Fresnel phase retrieval [71], no ground truth images are known, so one must obtain a reconstruction from only a single measurement \mathbf{y} .

This chapter based on [81].

In these practical applications with extremely limited or even nonexistent available data (measurement-only), it can be difficult or even impossible to train a diffusion model to approximate the underlying distribution well. Therefore, it is desirable to develop a method that can be applied regardless of dataset limitations. In this paper, we investigate applying a self-supervised loss at reconstruction time that allows the network to be adapted to the test data. Specifically, we investigate three different settings with various data availability ranging from nonexistent, limited, and abundant training data: (1) the *dataless* setting in which no training data is available and we are only given measurement \mathbf{y} , (2) 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, (3) and the *in-distribution* setting in which we have sufficient training data of the same distribution as the test dataset. Our goal is to develop methods that can be used to adapt the trained model to testing data in these settings with varying training data availability.

Recently, some previous works have aimed to address these problems by demonstrating that diffusion models have a stronger generalization ability than other deep learning methods [86], and that slight distribution mismatches between the training data and test data may not significantly degrade the reconstructed image quality. However, in cases of particularly sparse 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 [61, 10]. To address these challenges in the dataless setting, recent works first train a network on a large sample of synthetic data that may differ greatly from the test data. Then, at reconstruction time, the trained diffusion network is updated based on the measurement \mathbf{y} [10, 47]. This adaptation aims to shift the underlying prior learned by the network towards the appropriate prior corresponding to the test data distribution. However, with the large number of network parameters, a parameter-sensitive network adaptation approach is required at test time to avoid overfitting to the measurement. Furthermore, these methods have not been tested in the small dataset setting, and in the in-distribution setting, excessive network refining yields inconsistent improvement and degradation in the reconstructed image [10, 47].

Patch-based diffusion models have shown success both for image generation [183, 55] and for inverse problem solving [82]. In particular, the method of [82] involves training networks that take in only patches of images at training and reconstruction time, learning priors of the entire images from only image patches based on positional encoding. In cases of limited training data, [82] 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 dataless setting and the small dataset setting, given a severe lack of data. Inspired by this, we propose to use patch-based diffusion models to tackle the challenges arising from data limitations and mismatched distributions in a unified way. We first develop a test-time adaptation method to take a network trained on patches from a mismatched distribution and adapt it on the fly at reconstruction time. We then show that this method can improve image quality consistently *in all three settings* with varying data availability ranging from abundant, to limited and nonexistent.

In summary, our contributions are as follows:

- We integrate the patch-based diffusion model with the deep image prior (DIP) framework to take a network trained on a mismatched distribution and adapt it at reconstruction time towards the test-time distribution using a self-supervised loss.
- Experimentally, we find this approach of test-time diffusion adaptation leads to improved image quality both in terms of quantitative and qualitative metrics, outperforming using whole-image diffusion models under the same setting. In particular, even when the diffusion model was trained in-distribution, test-time adaptation offers an additional improvement.
- Theoretically, we show that in certain cases, using the self-supervised loss with a patch-based prior leads to a beneficial form of data augmentation compared to whole-image diffusion prior, providing further justification for the superior generalizability of adapting the patch-based prior.

7.2 Methods

This section extends the patch-based diffusion model framework of [82] so that it can be applied with the self-supervised loss for test-time adaptation.

7.2.1 Patch-based prior

We first zero pad the $N \times N$ image by an amount P on each side and model the resulting padded image \mathbf{x} . When choosing the i th patch offset tuple $(o_1, o_2) \in \{0, \dots, P - 1\}^2$ in Figure 7.1, we partition \mathbf{x} 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

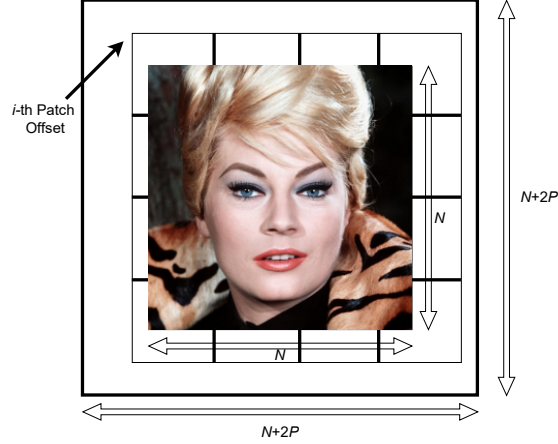


FIG 7.1 – Schematic for zero padding and partitioning image into patches. Each index i represents one of P^2 possible ways to choose a patch offset tuple.

image, our model for the data distribution has the form

$$p(\mathbf{x}) = \frac{1}{Z} \prod_{i=1}^{P^2} p_{i,B}(\mathbf{x}_{i,B}) \prod_{r=1}^{(k+1)^2} p_{i,r}(\mathbf{x}_{i,r}), \quad (7.1)$$

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 uses many 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} (or a patch of that image) 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. 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

$$\arg \min_{\theta} \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \sigma_t^2 I)} \|D_\theta(\mathbf{x} + \epsilon, \sigma_t) - \mathbf{x}\|_2^2, \quad (7.2)$$

where here $\mathbf{x} \sim p(\mathbf{x})$ 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. See [82] and Appendix A.4.3 for more details about the patch-based training process.

7.2.2 Diffusion self-supervision

We first review how to apply the [self-supervised loss](#) in conjunction with diffusion models. For each specific measurement \mathbf{y} , the original DIP framework optimizes the network parameters θ via the self-supervised loss (2.5) 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. That conditional expectation should be used in [test-time adaptation](#) of diffusion models.

We begin with the unconditional expectation by leveraging the patch-based prior. Following (7.1), we apply Tweedie’s formula to express the denoiser of \mathbf{x} solely in terms of 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 offset index i for each denoiser evaluation:

$$D_\theta(\mathbf{x}) \approx D_{i,B}(\mathbf{x}_{i,B}) + \sum_{r=1}^{(k+1)^2} D_{i,r}(\mathbf{x}_{i,r}). \quad (7.3)$$

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 (7.3) provides an *unconditional* estimate of the clean image; to obtain an approximate 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 [41].

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 such that $D_\theta(\mathbf{x}_t|\mathbf{y})$ becomes more consistent with the measurement:

$$\theta \leftarrow \arg \min_{\theta} \|\mathbf{y} - \mathbf{A} D_\theta(\mathbf{x}_t|\mathbf{y})\|_2^2. \quad (7.4)$$

Previously, additional LoRA parameters [78] were used as an injection to the network to leave the original parameters unchanged during this process [10, 47]. However, the

Algorithm 5 Test-time Adapted Diffusion 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 (7.3) 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 $L(\boldsymbol{\theta}) = \|\mathbf{y} - \mathbf{A}D_\theta(\mathbf{x}_t|\mathbf{y})\|_2^2$

 Update $\boldsymbol{\theta}$ by backpropagating $L(\boldsymbol{\theta})$

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 (7.3) with a random index i

 Run C iterations of CG for (2.7) 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

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.4.1 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 [39] or [181] 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.

After this step, we apply the refined network to compute a new estimate of the score of \mathbf{x}_t and then use it to update \mathbf{x}_t . Similar to the network refining step, we use the stochastic version of the denoiser given by (7.3) rather than the full version. Ref. [82] showed that for patch-based priors, Langevin dynamics [166] worked particularly well as a sampling algorithm, so we used it here in conjunction with CG steps to enforce data fidelity. Algorithm 5 summarizes the entire test-time adapted method for solving inverse problems.

7.2.3 Patch-based training

In the small dataset and in-distribution settings, we are provided clean training images that can be used to train the patch-based network. In the small dataset setting, the dataset

is too small to directly train a diffusion model from scratch, so we initialize the training process with a checkpoint trained on a different large dataset and then fine-tune this checkpoint on the small dataset. Works such as [138] and [207] found that this fine-tuning process allows diffusion models to be trained with a much smaller dataset than would otherwise be required. On the other hand, for the in-distribution setting, we initialize the weights of the network randomly and train on the large in-distribution dataset.

Ref. [183] found that training with varying patch sizes improved image generation performance compared to fixing the patch size to that used during inference. Here, we also applied a varying patch size scheme during fine-tuning as a method of data augmentation. We used the UNet architecture in [73] that can accept images of different sizes. Hence, the loss becomes

$$\arg \min_{\theta} \mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x} \sim p_d(\mathbf{x})} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \sigma_t^2 I)} \|D_{\theta}(\mathbf{x} + \epsilon, \sigma_t) - \mathbf{x}\|_2^2, \quad (7.5)$$

where $\mathbf{x} \sim p_d(\mathbf{x})$ represents drawing a patch of random size and location from an image belonging to the fine-tuning dataset. Appendix A.4.3 provides full details of the training process.

7.2.4 Theoretical Analysis

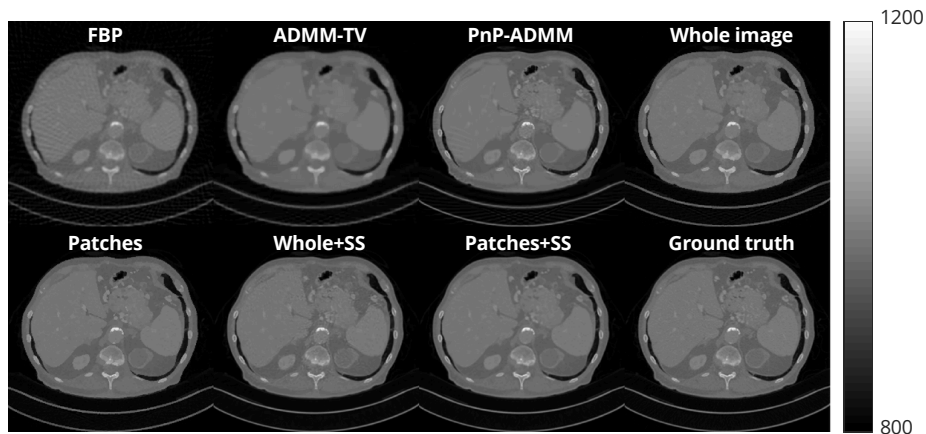


FIG 7.2 – Dataless setting: Results of 60 view CT reconstruction using self supervised (SS) loss. The display uses modified HU units to show more contrast between organs.

This section provides a rough sketch of why Algorithm 5 should work better for the patch-based model compared to a whole-image model. Firstly, rewrite (7.3) as follows:

$$D(\mathbf{x}) = \sum_c \mathbf{G}'_c D_\theta(\mathbf{G}_c \mathbf{x}, c), \quad (7.6)$$

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

Now we analyze (7.4) using this framework. The loss function for self-supervision is given by:

$$L(\boldsymbol{\theta}) = \left\| \mathbf{y} - \mathbf{A} \sum_c \mathbf{G}'_c D_\theta(\mathbf{G}_c \mathbf{x}, c | \mathbf{y}) \right\|_2^2. \quad (7.7)$$

For inverse problems such as superresolution and CT reconstruction, \mathbf{A} is a wide matrix that has full row rank. Hence, even when the measurement \mathbf{y} is noisy, there exists some \mathbf{x}_0 with $\mathbf{y} = \mathbf{A}\mathbf{x}_0$. Then we have

$$L(\boldsymbol{\theta}) = \left\| \mathbf{A}(\mathbf{x}_0 - \sum_c \mathbf{G}'_c D_\theta(\mathbf{G}_c \mathbf{x}, c | \mathbf{y})) \right\|_2^2 \quad (7.8)$$

$$L(\boldsymbol{\theta}) \leq \|\mathbf{A}\|_2^2 \left\| \mathbf{x}_0 - \sum_c \mathbf{G}'_c D_\theta(\mathbf{G}_c \mathbf{x}, c | \mathbf{y}) \right\|_2^2 \quad (7.9)$$

$$= \|\mathbf{A}\|_2^2 \left\| \sum_c \mathbf{G}'_c \mathbf{G}_c \mathbf{x}_0 - \sum_c \mathbf{G}'_c D_\theta(\mathbf{G}_c \mathbf{x}, c | \mathbf{y}) \right\|_2^2, \quad (7.10)$$

where in the last step, we have used the fact that at each diffusion iteration, the patches are nonoverlapping. Thus,

$$L(\boldsymbol{\theta}) \leq \|\mathbf{A}\|_2^2 \left\| \sum_c \mathbf{G}'_c (\mathbf{G}_c \mathbf{x}_0 - D_\theta(\mathbf{G}_c \mathbf{x}, c | \mathbf{y})) \right\|_2^2. \quad (7.11)$$

Now we have a sum of the form $\|\mathbf{z}_1 + \dots + \mathbf{z}_n\|_2^2$, which [162] showed to be upper bounded by $K(\|\mathbf{z}_1\|_2 + \dots + \|\mathbf{z}_n\|_2^2)$ for a fixed constant K . Applying this inequality and absorbing

$\|\mathbf{A}\|_2^2$ into the constant, we have

$$L(\boldsymbol{\theta}) \leq K \sum_c \|\mathbf{G}'_c(\mathbf{G}_c \mathbf{x}_0 - D_\theta(\mathbf{G}_c \mathbf{x}, c|\mathbf{y}))\|_2^2 \quad (7.12)$$

$$= K \sum_c \|\mathbf{G}_c \mathbf{x}_0 - D_\theta(\mathbf{G}_c \mathbf{x}, c|\mathbf{y})\|_2^2 \quad (7.13)$$

A similar derivation for the whole image model shows that the loss in that case is bounded by

$$L_w(\boldsymbol{\theta}) \leq K \|\mathbf{x}_0 - D_\theta(\mathbf{x}|\mathbf{y})\|_2^2. \quad (7.14)$$

Table A.18 shows that performance is improved by using more backpropagation iterations for the loss function; hence, although in practice we only perform a fixed number of iterations for speed, optimally we should aim to reduce the loss $L(\boldsymbol{\theta})$ to zero. Observe that (7.13) has the same form as the loss that would be used for refining the network with a whole image model (7.14). However, now instead of a loss of a single image, we now have individual losses of many patches of an image. For example, the experiments of Table 7.1 used 25 patches to tile each image for each diffusion iteration, so we had 25 losses. This method of data augmentation helps explain why the patch-based model obtains better performance than the whole-image model when performing test-time adaptation. We additionally note that although the positional encoding input into the network is different for each patch, the network does not separately learn a distribution for each position, as the weights are shared across these different positions. This is analogous to the analysis of [54], where a single diffusion model was trained on the 1000 classes of ImageNet with the class label of the image being included as an additional input to the network. Since each class only had around 1000 images, it would have been very difficult to train a diffusion model on only one of the classes, but by training across all the classes at once, a much better network can be trained.

7.3 Experiments

This section reports empirical results on sparse-view CT reconstruction, image deblurring, and image super-resolution. For each computational imaging application, we illustrate the benefits of using the self supervised (SS) loss in all three of the settings described in the introduction: *dataless* (zero-shot SS), *small dataset* (fine tuning followed by SS), and *in-distribution* (where one might expect SS to be counter-productive).



FIG 7.3 – Dataless setting: Results of deblurring using test-time adaptation via the self supervised (SS) loss and several comparison methods.

Experimental setup. For the CT experiments, we used the AAPM 2016 CT challenge data from [135]. We applied the same data processing methods as in [82] with the exception that we used all the XY (transaxial) slices of size 256×256 from the 9 training volumes to train the in distribution networks, yielding a total of 5936 slices. For the deblurring and superresolution experiments, we used the CelebA-HQ dataset [126] with each image having size 256×256 . The test data was a randomly selected subset of 10 of the images not used for training. In all cases, we report the average metrics across the test images: peak SNR (PSNR) in dB, and structural similarity metric (SSIM) [185].

In the dataless and small dataset settings, since there is insufficient data from the test distribution to train the diffusion model well, we first trained the network using generated synthetic data. This data consisted of phantom images consisting of randomly placed ellipses of different shapes and sizes. See Fig. A.37 for examples. These phantoms can be generated on the fly in large quantities. We used networks trained on grayscale phantoms for the CT experiments and networks trained on RGB phantoms for the deblurring and superresolution experiments. Appendix A.4.2 contains precise specifications of the phantoms. Then in the small dataset setting, we fine-tuned the network (originally trained on the phantom images) on the small dataset. The small dataset consisted of 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.4.1. In contrast, for the in-distribution setting, there was sufficient data to directly train a diffusion model on the in-distribution data, so we initialized the network weights randomly and trained only on the in-distribution dataset.

We trained the patch-based networks with 64×64 patches and used a zero padding value of 64, so that 5 patches in both directions were used to cover the target image. We used the network architecture in [93] for both the patch-based networks and whole-image networks. All networks were trained on PyTorch using the Adam optimizer with 2 A40 GPUs.

Diffusion test-time adaptation. To evaluate the effectiveness of test-time adaptation, we applied Algorithm 5 to solve each of the inverse problems. For the forward and backward projectors in CT reconstruction, we used the implementation provided by the ODL [173]. We performed two sparse-view CT (SVCT) experiments: one using 20 projection views, and one using 60 projection views. Both of these used a parallel-beam forward projector where the detector size was 512 pixels. For the deblurring experiments, we used a uniform blur kernel of size 9×9 and added white Gaussian noise with $\sigma = 0.01$ where the clean image was scaled between 0 and 1. For the superresolution experiments, 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 directly used the diffusion model without test-time adaptation. In particular, we used the network trained only on phantoms for the dataless setting, the fine-tuned network in the small dataset setting, and the in-distribution network in the last setting. We used the same sampling algorithm (Langevin dynamics) and inverse problem solving method (conjugate gradient descent) as Algorithm 5, but removed the test-time adaptation step. 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 (4.3) 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 [196] and PnP-RED [83]. 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 \mathbf{y} . For superresolution, the baseline was obtained by upsampling the low resolution image \mathbf{y} using nearest-neighbor interpolation. The implementation of ADMM-TV is in [74]. Finally, since we assume we do not have access to a large sample of clean training data, we used the off-the-shelf denoiser BM3D [49]. Appendix A.4.3 contains the values of all the parameters of the algorithms.

Tables 7.1, 7.3, and 7.4 show the main results for the dataless, small dataset, and in-distribution settings, respectively. In the latter two settings, we chose to redisplay only

TBL 7.1 – Comparison of quantitative results on four different inverse problems in the dataless setting with the self-supervised loss (SS). 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	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	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 [196]	30.20	0.838	36.75	0.932	28.98	0.815	27.29	0.796
PnP-RED [83]	27.12	0.682	32.68	0.876	28.37	0.793	27.73	0.809
Whole image	28.11	0.800	33.10	0.911	25.85	0.742	25.65	0.742
Patches [82]	27.44	0.719	33.97	0.934	26.77	0.782	26.12	0.759
Whole+SS [10]	33.19	0.861	40.47	0.957	29.50	0.831	27.07	0.701
Patches+SS (Ours)	33.77	0.874	41.45	0.969	30.34	0.860	28.10	0.827

the results of the best baseline out of the baseline, ADMM-TV, PnP-ADMM, and PnP-RED shown in Table 7.1 and labeled that row "Best baseline." Particularly for the dataless setting, applying the test-time adaptation method yields much higher quantitative results when averaged across the test dataset than simply using the pretrained diffusion model in all the inverse problems. However, we observe that even when the pretrained diffusion model was trained on the large in-distribution dataset, including the test-time adaptation step still resulted in further improvement in image quality. Thus, in contrast with many other self-supervised methods such as DIP, our method can avoid overfitting to the measurement and even benefit from test-time refinement. Appendix A.4.1 further analyzes overfitting and shows that by increasing the number of network refining iterations done per diffusion iteration, the image quality does not drop, indicating that overfitting is avoided. We further note that in all three settings, when using the self-supervised loss, the patch-based prior outperformed the whole image prior, which is consistent with our theoretical analysis of Algorithm 5. Lastly, Fig. 7.2 shows that some artifacts appear in the whole-image SS method that are absent in our patch SS method.

We also ran ablation studies to examine the effect of various parameters on the proposed method. [10] and [47] 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.4.1 shows the results of these experiments.

The initial motivation of using patch-based diffusion models was partially to solve high resolution imaging problems [82]. To show that our method scales to larger images, we ran experiments on 60 view CT reconstruction and deblurring with 512×512 images in the dataless setting. For the CT experiments, we still used the AAPM dataset [135] processed in the same way as for Table 7.1, but kept the slices in their original size of 512×512 . For deblurring, we used the FFHQ dataset [94] 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 7.1 as initializations for these out of distribution experiments. Since the patch size at reconstruction time was kept to be 64×64 as before, we used 9 patches in both directions (for 81 total) to tile each image. Table 7.2 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.

TBL 7.2 – Results of inverse problem solving in dataless setting for 512×512 images.

Method	CT, 60 views		Deblur	
	PSNR↑	SSIM↑	PSNR↑	SSIM↑
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

Small dataset fine-tuning. We further examined the effects of fine-tuning the networks using the small dataset with the patch-based model and the whole image model. Figures 7.6 and 7.7 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

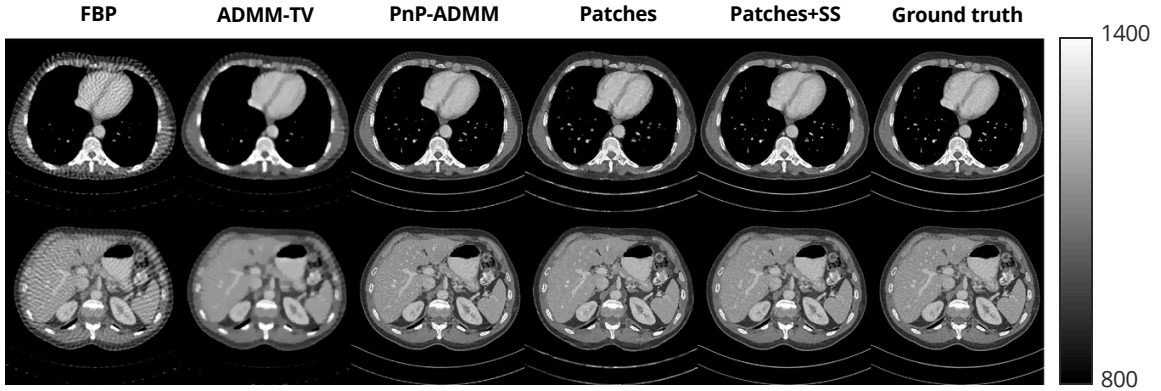


FIG 7.4 – Results of 60 view reconstruction on 512×512 images in dataless setting displayed in Hounsfield units.

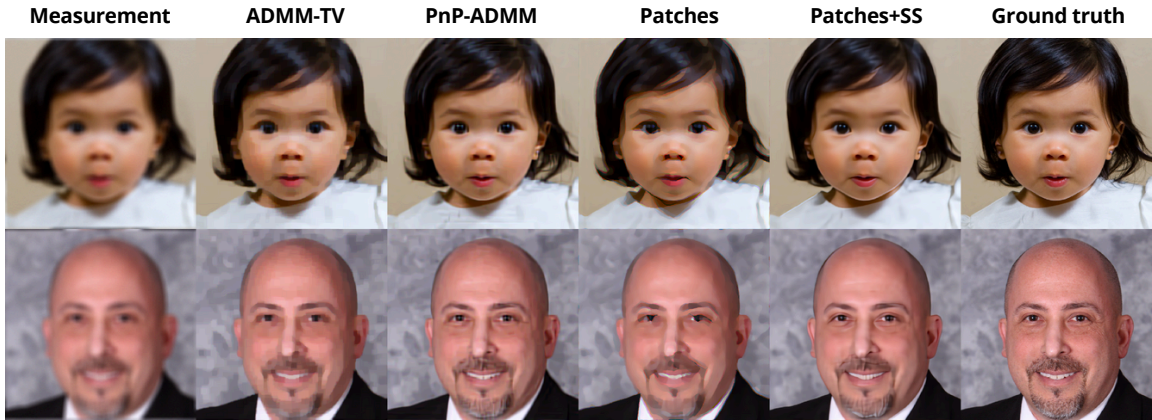


FIG 7.5 – Results of deblurring on 512×512 images in dataless setting.

stable performance throughout the entire training process. This illustrates that whole-image diffusion models exhibits severe overfitting problems when only a small amount of training data is unavailable, similar to the original DIP method. Furthermore, patch-based diffusion models assist greatly with this problem and the results are evident for solving inverse problems.

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 7.9 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

TBL 7.3 – Comparison of results for using self-supervised (SS) diffusion models in small dataset setting. The diffusion model is first trained on ellipse phantoms and then fine-tuned with the small dataset. Best results are in bold.

Method	CT, 20 Views		CT, 60 Views		Deblurring		Superresolution	
	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	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	33.44	0.875	41.21	0.965	29.25	0.840	28.10	0.827
Whole image+SS	34.57	0.881	41.34	0.962	30.16	0.852	27.72	0.814
Patches+SS	36.43	0.914	42.42	0.971	30.56	0.867	28.60	0.834

TBL 7.4 – Comparison of results for using diffusion models trained on thousands of in-distribution images to solve inverse problems. Best results are in bold. Self-supervised refinement of network weights is beneficial even in this setting.

Method	CT, 20 Views		CT, 60 Views		Deblurring		Superresolution	
	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow
Best baseline	30.20	0.838	36.75	0.932	28.98	0.815	27.73	0.809
Whole image	33.99	0.886	41.67	0.969	29.87	0.851	28.33	0.801
Patches	34.02	0.889	41.70	0.967	30.12	0.865	28.49	0.835
Whole image+SS	35.38	0.897	41.68	0.966	30.31	0.854	27.96	0.797
Patches+SS	36.82	0.923	42.33	0.970	30.78	0.875	28.72	0.842

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.

Finally, 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 of CT scans [6]. 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

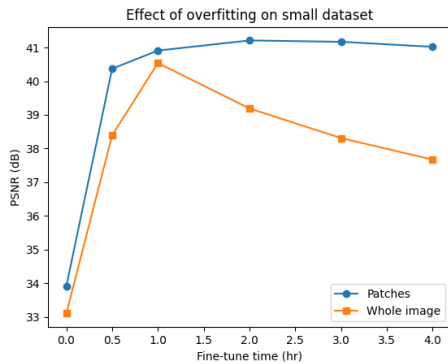


FIG 7.6 – Comparison of PSNR between patch-based model and whole-image model for overfitting in small dataset setting.

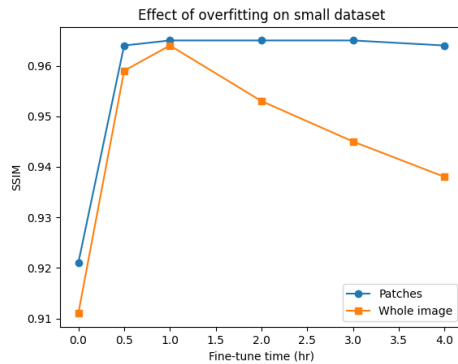


FIG 7.7 – Comparison of SSIM between patch-based model and whole-image model for overfitting in small dataset setting.



FIG 7.8 – Small dataset setting: Results of inverse problem solving. Top row is 60 view CT recon, middle row is deblurring, and bottom row is superresolution.

dataset was the same as the one used in Table 7.1. Table 7.5 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 test-time adaptation.

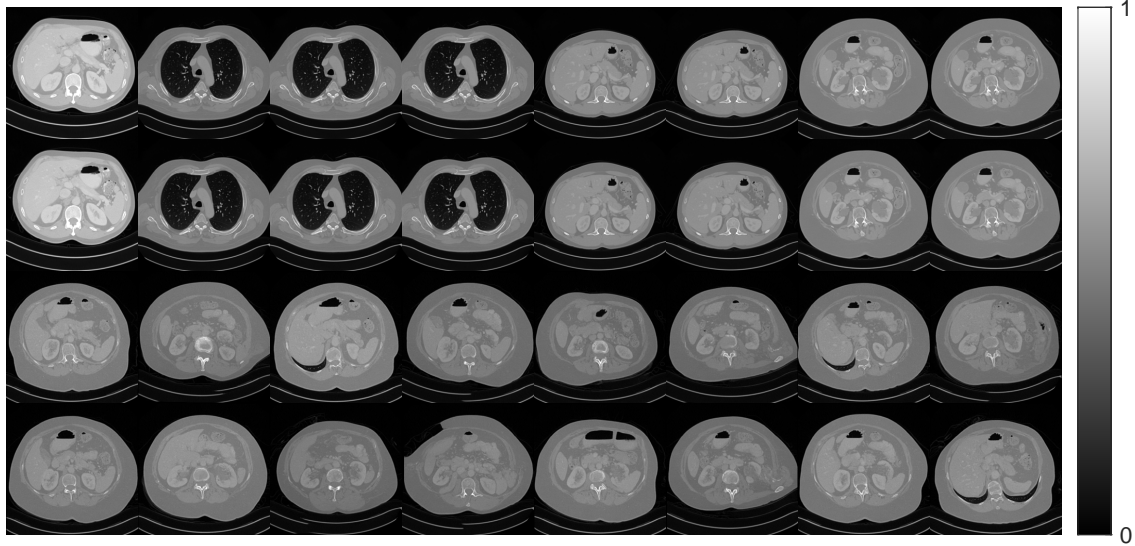


FIG 7.9 – Unconditional generation of CT images from networks fine-tuned in the small dataset setting. Top two rows were generated with the whole image model; bottom two rows were generated with the patch-based model.

7.4 Conclusion

This paper presented a method of using patch-based diffusion models with test-time adaptation to solve inverse problems when the data distribution might be mismatched from the trained network. In particular, we conducted experiments in setting when no data is available, when a small dataset of training images is available, and when a large in-distribution dataset is available. In all the settings, applying the self-supervised loss improved image quality, even for a well-trained network. Furthermore, the patch-based method outperformed whole-image methods in a variety of inverse problems and we provided theoretical justifications to explain this improvement. 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 [198]. Limitations of the work include a slow runtime for the test-time adapted algorithm and a lack of theoretical guarantees for dataset size requirements. Providing uncertainty quantification is also an open problem for such self-supervised methods.

TBL 7.5 – CT reconstruction results where the initial checkpoint was trained on LIDC dataset and refined on the fly with the AAPM measurement.

Dataset size	CT, 20 views		CT, 60 views	
	PSNR↑	SSIM ↑	PSNR↑	SSIM ↑
Whole image	33.04	0.874	40.43	0.949
Patches	33.88	0.886	40.96	0.955
Whole image+SS	35.01	0.894	41.95	0.967
Patches+SS (Ours)	36.34	0.918	42.32	0.972

CHAPTER 8

SPAR: Refine a Single Pretrained Diffusion Model to Solve Inverse Problems in Many Modalities

8.1 Motivation

Inverse problems arise in real-world settings where an object cannot be observed directly and indirect measurements are obtained. Such inverse problems commonly occur in many imaging applications, for example, deblurring and super-resolution with natural images; as well as in medical imaging, such as computed tomography (CT) and magnetic resonance imaging (MRI) reconstruction. Many real-world inverse problems are highly ill-posed and difficult to solve because the solution may not be unique and the measurement noise can further complicate the reconstruction process. Furthermore, in practical applications, imaging modalities and resolutions can vary significantly, complicating the reconstruction process. Generally, the measurement acquisition process is represented as $\mathbf{y} = \mathcal{A}(\mathbf{x}) + \epsilon$, where \mathbf{y} is the unknown measurement, \mathbf{x} is the unknown clean image to be recovered, \mathcal{A} is an operator representing the measurement system, and ϵ is unknown random noise.

A common approach to solving inverse problems is to enforce prior knowledge about \mathbf{x} to obtain a high-quality reconstruction, such as sparsity in compressed sensing [148]. Recently, diffusion models have become state-of-the-art methods for generating high-quality image data from random noise samples [73, 166] and, as such, can be integrated as a learned prior into many inverse problem solvers. Since diffusion models were proposed to be

This chapter is based on a paper currently under review [79]. Code at <https://github.com/jasonhu4/SPAR>.

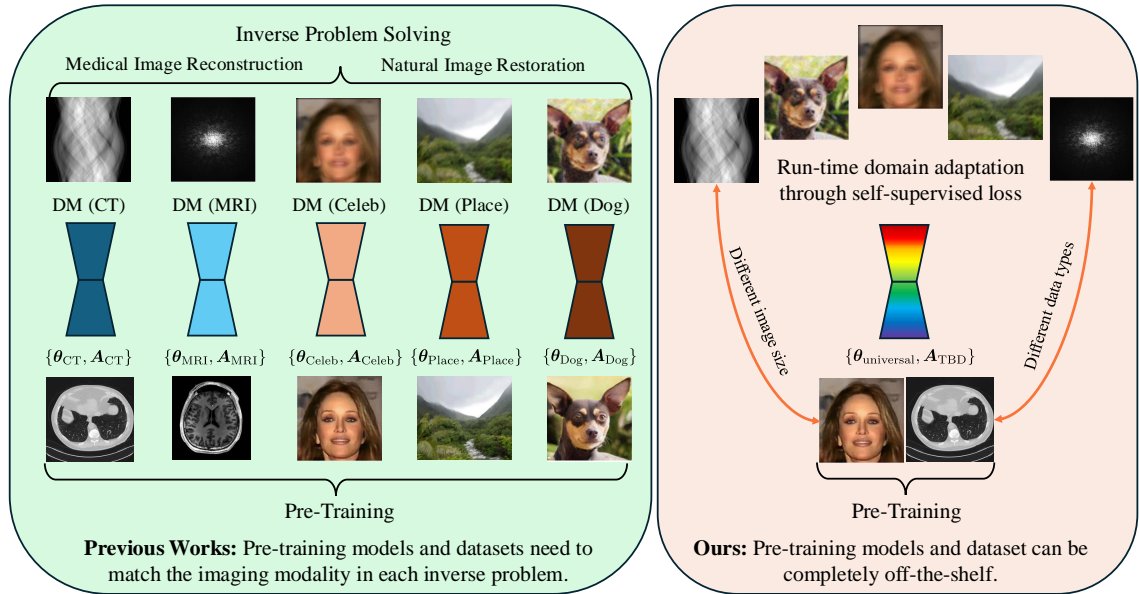


FIG 8.1 – Comparison of traditional domain-specific diffusion models (DM) versus the proposed SPAR method. Left: previous works require separate pre-training for each domain/modality (CT, MRI, CelebA, Places, Dogs). Right: The proposed method enables domain adaptation at testing time through self-supervised learning, avoiding the need for domain-specific pre-training.

trained unconditionally on natural images [73, 168], the key challenge in many diffusion inverse solvers (DIS) [181, 95, 39, 163, 115] is enforcing consistency with observed measurements while simultaneously leveraging the image prior learned by the unconditional diffusion model. Owing to the strong learned prior, diffusion models provide excellent performance for various inverse problems, especially when the measurements are very compressed or noisy [39, 42, 41, 115, 51, 5, 143].

Nevertheless, separate diffusion models are usually trained for different imaging applications, including different modalities or for various anatomic sites *e.g.*, contrast CT, low-dose CT, brain MRI, knee MRI, natural RGB images, etc. Training diffusion models from scratch is very computationally expensive and typically requires large-scale high-quality training datasets and substantial GPU resources [73, 168], which may be infeasible in many practical scientific or medical application domains having only scarce data. Therefore, on-the-fly network fine-tuning methods have been developed to tackle this challenge. Specifically, [47] and [10] first train a diffusion model on a dataset with constrained distribution, then perform network refining on the fly to solve an inverse problem where the clean image corresponds to a different data distribution but the same modality (*e.g.*,

training with brain MRI and then reconstructing knee MRI). This network refining relies on a self-supervised loss that enforces the clean image (and by extension, the output of the denoiser learned by the diffusion network) to be consistent with the measurement. By enforcing measurement consistency as a constraint, the pretrained diffusion model weights are updated as part of the reconstruction process, adapting the diffusion prior to out-of-distribution (OOD) test samples that were unseen during the original diffusion training phase. A patch-based prior [81] enhances these methods, improving robustness to overfitting to the measurement y . However, these methods still exhibit instability at the tail ends of the noise scale [47], which requires technical tricks such as truncating the diffusion process to resolve. Furthermore, although existing methods have explored how to generalize across different distributions within the same imaging modalities (such as different CT datasets), they still require training multiple networks for applications of different modalities (e.g., CT and MRI), as the image contents, pixel sizes and data types (e.g., real and complex valued) differ across these modalities.

Our proposed ***Single Pretrain Always Refine (SPAR)*** method tackles these challenges as follows. First, we developed a *unified model architecture* based on a channel-wise conditioning mechanism that accommodates images of different modalities, including grayscale CT images, complex-valued MRI images, and RGB colored images. This versatile framework enables the diffusion network to be pretrained on a mixture of images from diverse modalities, thereby ensuring that the learned distribution captures a broad range of image types that can be reconstructed with test-time refinement. To further enhance the model flexibility, we adopt a *patch-based diffusion prior* that handles images of different sizes. Furthermore, motivated by the observation that diffusion sampling behavior varies across timesteps (i.e., noise levels), we propose a *novel majorizer-based self-supervised loss* function for a more effective test-time adaptation, especially at lower noise levels, while mitigating the risk of overfitting. Overall, our proposed method provides a unified framework that enables a single versatile pretrained diffusion model to address a wide range of inverse problems across different imaging modalities through test-time domain adaptation, as illustrated in Fig. 8.1. Our key contributions are as follows:

- We introduce ***SPAR***, a novel test-time adaptation framework for solving diverse inverse problems using a single, pretrained patch-based diffusion model. A core innovation is a channel-wise conditioning mechanism that allows the pretrained diffusion network to process different image modalities (e.g., grayscale, RGB, complex-valued) within a unified architecture.

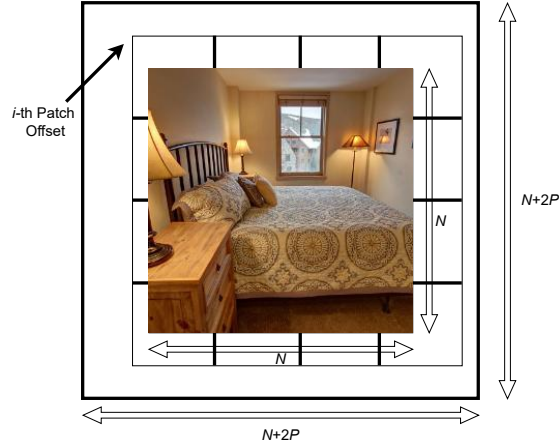


FIG 8.2 – Schematic illustration for zero padding and partitioning image into patches. Each index i represents one of P^2 possible ways to choose a patch offset tuple.

- To address the distributional mismatch between training and test data, we propose a novel self-supervised loss function to continuously refine the model during reconstruction. We provide theoretical justification to show that, at lower noise levels, this loss function is less prone to errors that arise from imperfect denoising score matching at training time.
- With the proposed modality- and resolution-agnostic framework, we experimentally validate that SPAR achieves improved image quality across a wide spectrum of inverse problems across imaging modalities and resolutions, even with only a single pretrained model.

8.2 Methods

8.2.1 Patch-based Prior

Previous works show that patch-based diffusion models achieve consistent performance in a wide variety of inverse problems where the amount of training data varies significantly [82, 81]. The patch-based diffusion models exhibit generalization capabilities especially under data-constrained conditions, which motivates our adoption of this framework as the foundation for our method.

In patch-based diffusion models, one firsts zero pads each $N \times N$ training image by P pixels on each side. Slightly recycling notation, we let \mathbf{x} denote the resulting padded image. We partition \mathbf{x} into many square patches, with one bordering region of zeros, by first choosing the i th patch offset tuple $(o_1, o_2) \in \{0, \dots, P-1\}^2$ according to Figure 8.2.

Since $k = N/P$ patches are needed in one direction to perfectly tile the image, our model for the data distribution has the form

$$p(\mathbf{x}) = \frac{1}{Z} \prod_{i=1}^{P^2} p_{i,B}(\mathbf{x}_{i,B}) \prod_{r=1}^{(k+1)^2} p_{i,r}(\mathbf{x}_{i,r}), \quad (8.1)$$

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 i th patch offset, $p_{i,r}$ is the probability distribution of that region, and Z is a normalizing factor. This model takes a product over all possible patch offsets, eliminating boundary artifacts that would occur if a fixed patch offset was used. However, when computing the denoiser of \mathbf{x} using Tweedie’s formula [56] in practice, we use a stochastic version of (8.1) corresponding to selecting a random index i for each denoiser evaluation:

$$D_\theta(\mathbf{x}) \approx D_{i,B}(\mathbf{x}_{i,B}) + \sum_{r=1}^{(k+1)^2} D_{i,r}(\mathbf{x}_{i,r}). \quad (8.2)$$

8.2.2 Model Pretraining

For pretraining, we use a neural network $D_\theta(\mathbf{x}, \sigma_t)$ that takes as input a noisy image \mathbf{x} (or a patch of that image) along with the noise level σ_t . To incorporate positional information, we define the x positional array as a 2D array representing the x coordinates of each pixel in the image, scaled between -1 and 1. Similarly, the y positional array is defined for the y coordinates of the pixels. To enable the network to learn patch distributions that vary based on location within the image, we extract patches from these positional arrays corresponding to the image patches. These positional patches are concatenated along the channel dimension of the noisy image patch, forming a unified input array for the network. (The network output is just the image patch, not the position channels.) As our approach is based on patch-based priors, we apply denoising score matching on image patches rather than the whole image. The training loss is:

$$\arg \min_{\theta} \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \sigma_t^2 I)} \|D_\theta(\mathbf{x} + \epsilon, \sigma_t) - \mathbf{x}\|_2^2, \quad (8.3)$$

where $\mathbf{x} \sim p(\mathbf{x})$ represents a patch drawn from an image in the training dataset, σ_t is a predetermined noise schedule, and \mathcal{U} denotes the uniform distribution. Training on patches of varying sizes can reduce training time and improve the network’s generation

capability [183]. Thus, we use a UNet architecture as the backbone [73] to take in patches of different sizes, which is trained with varying patch sizes [183, 82].

Channel-wise conditioning mechanism. Real-world images span a wide range of modalities and types, so our aim is to train a diffusion model that can be adapted to the different possible images that may be required at reconstruction time. However, grayscale images only have 1 channel, whereas natural images have 3 channels. We handle both types of images with a single network that inputs and outputs 3-channel images by concatenating the image with itself. Specifically, we duplicate each grayscale image three times and concatenate it along the channel dimension so it becomes a 3-channel image that can fit as inputs to the network. We add the *same* random Gaussian noise to each channel; Appendix A.5.2 justifies this process over the alternative approach of using different noises across channels. Thus, the network should learn to denoise all 3 identical channels of this image.

Finally, we include an additional class label [93] as a network input to inform the network whether the image input was a grayscale or a colored image. References [140] and [93], among others, have shown that including class labels to delineate between different image distributions can improve training and inference results. Appendix A.5.2 provides a theoretical derivation that shows this method of performing denoising score matching for grayscale images is equivalent to training diffusion models for grayscale images in the traditional way (with single-channel networks) under certain regularity assumptions. Appendix A.5.4 provides the pseudocode for the pretraining stage of SPAR.

8.2.3 Self-Supervised Model Finetuning

To refine a diffusion model for OOD testing samples at reconstruction time, we use the measurement-consistency constraints to update the model weights (in a similar manner as [47]) by

$$\boldsymbol{\theta} \leftarrow \arg \min_{\boldsymbol{\theta}} L(\boldsymbol{\theta}; \mathbf{x}_t), \quad L(\boldsymbol{\theta}; \mathbf{x}_t) \triangleq \Psi(D_{\boldsymbol{\theta}}(\mathbf{x}_t | \mathbf{y})) \quad (8.4)$$

$$\Psi(\mathbf{x}) \triangleq \frac{1}{2} \|\mathbf{y} - \mathbf{A} \mathbf{x}\|_2^2. \quad (8.5)$$

We replace (8.5) with a quadratic majorizer of the form:

$$\phi(\mathbf{x}; \tilde{\mathbf{x}}) = \Psi(\tilde{\mathbf{x}}) + \langle \nabla \Psi(\tilde{\mathbf{x}}), \mathbf{x} - \tilde{\mathbf{x}} \rangle + \frac{\mathcal{L}}{2} \|\mathbf{x} - \tilde{\mathbf{x}}\|_2^2, \quad (8.6)$$

where \mathcal{L} denotes a Lipschitz constant of $\nabla\Psi$ and $\tilde{\mathbf{x}}$ can be any point. Appendix 8.2.4 shows that $\phi(\mathbf{x}; \tilde{\mathbf{x}}) \geq \Psi(\mathbf{x})$. When applying this majorizer with the self-supervised loss $L(\boldsymbol{\theta}, \tilde{\mathbf{x}})$, the majorizer of the network refining loss is

$$\check{L}(\boldsymbol{\theta}; \tilde{\mathbf{x}}) \triangleq \langle \mathbf{A}'(\mathbf{A}\tilde{\mathbf{x}} - \mathbf{y}), \mathbf{x}_\theta \rangle + \frac{\mathcal{L}}{2} \|\mathbf{x}_\theta - \tilde{\mathbf{x}}\|_2^2, \quad (8.7)$$

where $\mathbf{x}_\theta = D_\theta(\mathbf{x}_t|\mathbf{y})$, \mathbf{A}' denotes the (Hermitian) transpose of \mathbf{A} , and we set $\tilde{\mathbf{x}}$ to be the denoised image from the previous diffusion iteration with a detached gradient. We drop constants that do not depend on \mathbf{x}_θ . In practice, we treat \mathcal{L} as a tunable parameter.

The benefit of using this surrogate loss lies in the observation of [98] that diffusion models tend to have problems learning the score function at low levels of added noise. Section 8.2.4 shows that for low noise levels, (8.7) is less sensitive to inaccuracies in the learned score function than (8.4). However, for high noise levels, (7.4) is more accurate because the majorizer is an upper bound for the true loss function, with $\phi(\mathbf{x}; \tilde{\mathbf{x}}) = \Psi(\mathbf{x})$ only when $\mathbf{x} = \tilde{\mathbf{x}}$. At high noise levels, $\tilde{\mathbf{x}}$ is likely to change significantly from one iteration to the next, because each diffusion iteration could substantially alter the image with significant added noise. Hence, the gap between the true loss (7.4) and the majorizer loss will be large. Motivated by such observations, we propose a novel noise-varying self-supervised loss function for test-time diffusion model adaptation, which is designed to use (7.4) for the higher noise levels while using (8.7) for lower noise levels. Section 8.2.4 provides a more rigorous theoretical analysis to justify the benefit of this new proposed self-supervised loss. In the implementation, we used a parameter search to find the optimal threshold σ' in Algorithm 6; Appendix A.5.3 compares the results of using different values of σ' .

Finally, though the original network is trained on 3-channel images, we also wish to also perform reconstruction on 1-channel images (grayscale images) and 2-channel images (MRI images with the real and imaginary parts separated and concatenated along the channel dimension). Therefore, when using (8.2) to denoise these images, we need to first expand the image to 3 channels, e.g., by repeating one of the channels. We then process the resulting 3-channel output by using appropriate channel averaging to get a denoised 1- or 2-channel image. Mathematically, this is defined as follows:

- When \mathbf{x} has 1 channel, we define the operator $\mathbf{H} : \mathbb{R}^{1 \times M \times M} \rightarrow \mathbb{R}^{3 \times M \times M}$ that concatenates \mathbf{x} with itself along the channel dimension 3 times. In essence: $\mathbf{H} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \otimes \mathbf{I}_{M^2}$.

- When \mathbf{x} has 2 channels, we define the operator $\mathbf{H} : \mathbb{R}^{2 \times M \times M} \rightarrow \mathbb{R}^{3 \times M \times M}$ that concatenates an additional copy of the first channel of \mathbf{x} with itself before the second channel. In essence: $\mathbf{H} = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \otimes \mathbf{I}_{M^2}$.
- Denoting by \tilde{D}_θ the denoising network trained on 3-channel images, the denoiser of \mathbf{x} is given by

$$D(\mathbf{x}) = \begin{cases} \frac{1}{3} \mathbf{H}' \tilde{D}_\theta(\mathbf{H}\mathbf{x}), & \mathbf{x} \in \mathbb{R}^{1 \times M \times M}, \\ \mathbf{H}' \tilde{D}_\theta(\mathbf{H}\mathbf{x}) \oslash [2 \ 1], & \mathbf{x} \in \mathbb{R}^{2 \times M \times M}, \end{cases} \quad (8.8)$$

where \oslash denotes channel-wise division.

We can express (8.8) succinctly as $D(\mathbf{x}) = \mathbf{H}^\dagger \tilde{D}_\theta(\mathbf{H}\mathbf{x})$, where \mathbf{H}^\dagger denotes the pseudoinverse of the operator \mathbf{H} . In general, we have $\mathbf{H}^\dagger = (\mathbf{K} \otimes \mathbf{I})^\dagger = \mathbf{K}^\dagger \otimes \mathbf{I}^\dagger = \mathbf{K}^\dagger \otimes \mathbf{I}$. In particular:

$$\begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}^\dagger = \frac{1}{3} \cdot [1 \ 1 \ 1], \quad \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}^\dagger = \begin{bmatrix} 0.5 & 0.5 & 0 \\ 0 & 0 & 1 \end{bmatrix}. \quad (8.9)$$

The first case corresponds to taking an average across all 3 channels to obtain a single-channel image, which is consistent with the first part of (8.8). The second case corresponds to averaging the first 2 channels to obtain the first output channel and setting the third channel to be the second output channel, which is consistent with the second part of (8.8). Hence, we have shown that we convert a 3-channel denoiser \tilde{D}_θ to a 1- or 2-channel denoiser using $D(\mathbf{x}) = \mathbf{H}^\dagger \tilde{D}_\theta(\mathbf{H}\mathbf{x})$.

Algorithm 6 summarizes the entire proposed algorithm with self-supervised diffusion model refining at reconstruction time. Because [47] observed irregularities at the high noise levels when using similar self-supervised diffusion algorithms, we adopt a similar approach as [45] where we initialize the image using a simple baseline $\mathbf{A}^\dagger \mathbf{y}$ and choose a much lower value for σ_T than would otherwise be optimal.

8.2.4 Theoretical Analysis of Majorizer Loss

This section shows that the majorizer loss is less sensitive to inaccuracies in the estimated score function for low noise levels. We return to the specific majorizer given by (8.6). Previous works analyzed theoretical properties of diffusion models assuming that the prior is a mixture of Gaussians [180, 113]. Each component of the mixture may represent one

Algorithm 6 SPAR Image Reconstruction

Input: $\sigma_1 < \sigma_2 < \dots < \sigma_T, \epsilon > 0, T, \mathbf{y}, K, \sigma'$
Initialize $\mathbf{x} \sim \mathcal{N}(\mathbf{A}^\dagger \mathbf{y}, \sigma_T^2 \mathbf{I})$
Initialize $\tilde{\mathbf{x}} \sim \mathcal{N}(0, \mathbf{I})$
for $t = T : 1$ **do**
 if $t \bmod K = 0$ **then**
 Compute $D_\theta(\mathbf{x}_t)$ using (8.2) with a random index i
 and apply (8.8) if \mathbf{x} has fewer than 3 channels.
 Compute $D_\theta(\mathbf{x}_t|\mathbf{y})$ using $D_\theta(\mathbf{x}_t)$ and (2.7).
 if $\sigma_t > \sigma'$ **then**
 Define $L(\boldsymbol{\theta}) = \|\mathbf{y} - \mathbf{A}D_\theta(\mathbf{x}_t|\mathbf{y})\|_2^2$.
 else
 Define $L(\boldsymbol{\theta}) = \langle \mathbf{A}'(\mathbf{A}\tilde{\mathbf{x}} - \mathbf{y}), \mathbf{x}_\theta \rangle + \frac{\epsilon}{2} \|\mathbf{x}_\theta - \tilde{\mathbf{x}}\|_2^2$,
 where $\mathbf{x}_\theta = D_\theta(\mathbf{x}_t|\mathbf{y})$.
 end if
 Update $\boldsymbol{\theta}$ by backpropagating $L(\boldsymbol{\theta})$.
 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 (8.2) with a random index i
 and apply (8.8) if \mathbf{x} has fewer than 3 channels.
 Solve (2.7) for $D(\mathbf{x}_t|\mathbf{y})$ and set it to $\tilde{\mathbf{x}}$.
 Set $\mathbf{s}_t = (D(\mathbf{x}_t|\mathbf{y}) - \mathbf{x}_t)/\sigma_t^2$.
 Set \mathbf{x}_{t-1} to $\mathbf{x}_t + \frac{\alpha_t}{2} \mathbf{s}_t + \sqrt{\alpha_t} \mathbf{z}$.
end for

image class in the dataset. When solving inverse problems, the measurements lack fine details but provide global information about the ground truth image structure that can help identify the image class. Thus, the mixture of Gaussians approximately reduces to a single Gaussian from which the ground truth belongs with high probability. As such, the following analysis assumes that the clean image distribution $p(\mathbf{x})$ follows a Gaussian with mean μ and covariance Σ . Furthermore, Algorithm 6 uses the majorizer loss only for small σ_t , so we assume $\sigma_t \ll 1$ for the remainder of the analysis. Finally, since the number of diffusion iterations is large, σ_t does not change significantly between consecutive sampling steps.

We begin by analyzing the following problem setup: suppose $\mathbf{X} \sim \mathcal{N}(\mu, \Sigma)$ and $\mathbf{Y} = \mathbf{X} + \sigma\mathbf{z}$, where $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$. Because \mathbf{X} and \mathbf{Y} are jointly Gaussian:

$$\begin{aligned}\mathbb{E}[\mathbf{X}|\mathbf{Y}] &= \mu + \text{Cov}(\mathbf{X}, \mathbf{Y}) \text{Cov}(\mathbf{Y})^{-1}(\mathbf{Y} - \mu) \\ \text{Cov}(\mathbf{X}, \mathbf{Y}) &= \text{Cov}(\mathbf{X}, \mathbf{X} + \sigma\mathbf{z}) = \text{Cov}(\mathbf{X}) = \Sigma \\ \text{Cov}(\mathbf{Y}) &= \text{Cov}(\mathbf{X}) + \text{Cov}(\sigma\mathbf{z}) = \Sigma + \sigma^2\mathbf{I} \\ \mathbb{E}[\mathbf{X}|\mathbf{Y}] &= \mu + \Sigma(\Sigma + \sigma^2\mathbf{I})^{-1}(\mathbf{Y} - \mu).\end{aligned}$$

Expanding out the inverse with a geometric series and assuming σ is small, this simplifies to

$$\mathbb{E}[\mathbf{X}|\mathbf{Y}] \approx \mathbf{Y} - \sigma^2\Sigma^{-1}(\mathbf{Y} - \mu). \quad (8.10)$$

Returning to the problem at hand, assume the measurement is noiseless so that $\mathbf{y} = \mathbf{A}\mathbf{x}_0$. For small σ_t , $\mathbb{E}[\mathbf{x}_0|\mathbf{x}_t, \mathbf{y}]$ is approximately a projection of $\mathbb{E}[\mathbf{x}_0|\mathbf{x}_t]$ onto the affine subspace defined by $\mathbf{y} = \mathbf{A}\mathbf{x}_0$; past diffusion inverse solvers [41, 181] have made similar assumptions. Using $\mathbf{X} \triangleq \mathbf{x}_0$ and $\mathbf{Y} \triangleq \mathbf{x}_t = \mathbf{x}_0 + \sigma_t\mathbf{z}$ in (8.10) yields $\mathbb{E}[\mathbf{x}_0|\mathbf{x}_t] = \mathbf{x}_t + O(\sigma^2)$. Thus,

$$\mathbb{E}[\mathbf{x}_0|\mathbf{x}_t, \mathbf{y}] - \mathbb{E}[\mathbf{x}_0|\mathbf{x}_t] \approx \mathbf{P}(\mathbf{x}_t + O(\sigma_t^2) - \mathbf{x}_0), \quad (8.11)$$

where \mathbf{P} is a linear operator that represents orthogonally projecting onto the row space of \mathbf{A} . Since $\mathbf{x}_t - \mathbf{x}_0 = \sigma_t\mathbf{z} = O(\sigma_t)$, we have $\mathbb{E}[\mathbf{x}_0|\mathbf{x}_t, \mathbf{y}] - \mathbb{E}[\mathbf{x}_0|\mathbf{x}_t] = O(\sigma_t)$.

Further assume that the trained denoiser satisfies $D_\theta(\mathbf{x}_t) = \mathbb{E}[\mathbf{x}_0|\mathbf{x}_t] + \Delta$, where Δ is a small error term. If this denoising neural network were simply the identity function, the error Δ would be $\mathbf{x}_t - \mathbf{x}_0$, which is $O(\sigma_t)$; hence, a well trained denoiser should have $\Delta < O(\sigma_t)$. In conclusion, the difference between $D_\theta(\mathbf{x}_t)$ and $\mathbb{E}[\mathbf{x}_0|\mathbf{x}_t, \mathbf{y}]$ is $O(\sigma_t)$.

Next we analyze the majorizer loss. From the setup above, we have $\mathbf{X} \triangleq \mathbf{x}_0$, $\mathbf{Y} \triangleq \mathbf{x}_0 + \sigma_t\mathbf{z}$, $\tilde{\mathbf{x}} = D_\theta(\mathbf{x}_{t+1}|\mathbf{y}) = \mathbb{E}[\mathbf{X}|\mathbf{Y}] + O(\sigma_t)$, using the conclusion from the previous paragraph in the final equality as well as assuming the noise scale is changing consecutively in small steps $\sigma_t \approx \sigma_{t+1}$. Note that $\mathbf{x}_t = \mathbf{x}_0 + \sigma_t\mathbf{z}$. Then $\tilde{\mathbf{x}} - \mathbf{x}_0 = \sigma_t\mathbf{z} - \sigma_t^2\Sigma^{-1}(\mathbf{x}_0 + \sigma_t^2\mathbf{z} - \mu) + O(\sigma_t)$. This shows that $\tilde{\mathbf{x}} - \mathbf{x}_0$ is approximately $O(\sigma_t)$, when σ_t is small. Now we rewrite (8.7) as

$$L(\theta) = (\tilde{\mathbf{x}} - \mathbf{x}_0)' \mathbf{A}' \mathbf{A} \mathbf{x} + \frac{L}{2} \|(\mathbf{x} - \mathbf{x}_0) + (\mathbf{x}_0 - \tilde{\mathbf{x}})\|_2^2. \quad (8.12)$$

Note that $\mathbf{x} = D_\theta(\mathbf{x}_t|\mathbf{y})$ (before any network refining) is the current prediction of \mathbf{x}_0 from a noisy version of \mathbf{x}_0 . We assume the noise level σ_t does not change significantly between consecutive sampling steps, so we can follow the same derivation above to get

that $\mathbf{x} - \mathbf{x}_0$ is $O(\sigma_t)$. This shows that the entire term inside of the squared Euclidean norm in (8.12) is $O(\sigma_t)$, so the term itself is $O(L\sigma_t^2)$. In contrast, the linear term in this loss is only $O(\|\mathbf{A}\|_2^2\sigma_t)$, where typically L is chosen to be on the same order as $\|\mathbf{A}\|_2^2$. Hence, when σ_t is small, the linear term dominates the loss.

Next we look at the gradients of the loss functions. For convenience, set $\mathbf{x}_\theta = D_\theta(\mathbf{x}_t|\mathbf{y})$ so that (7.4) becomes $L(\theta) = (1/2)\|\mathbf{A}\mathbf{x}_0 - \mathbf{A}\mathbf{x}_\theta\|_2^2$. Hence, the backpropagation step through the network uses the gradients

$$\frac{\partial L}{\partial \theta} = \frac{\partial L}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial \theta} = \mathbf{A}'\mathbf{A}(\mathbf{x}_\theta - \mathbf{x}_0) \frac{\partial \mathbf{x}}{\partial \theta}. \quad (8.13)$$

The gradient for the quadratic majorizer loss is

$$\frac{\partial L}{\partial \theta} = (\mathbf{A}'\mathbf{A}(\tilde{\mathbf{x}} - \mathbf{x}_0) + L(\mathbf{x}_\theta - \tilde{\mathbf{x}})) \frac{\partial \mathbf{x}}{\partial \theta}. \quad (8.14)$$

The inaccuracies in the approximate score function arise in \mathbf{x}_θ , and hence $\mathbf{x}_\theta - \mathbf{x}_0$ in (8.13) and $\mathbf{x}_\theta - \tilde{\mathbf{x}}$ in (8.14). Crucially, we showed earlier that the dominant term in (8.14) is $\mathbf{A}'\mathbf{A}(\tilde{\mathbf{x}} - \mathbf{x}_0)$. Therefore, while errors in the score function fully propagate themselves to the gradient when using the original loss (8.13), those errors have much less effect on the gradient of the majorizer (8.14). This analysis provides a theoretical justification why the majorizer loss is less sensitive to errors in the score function than the original loss when σ is small, and guided the design of Alg. 6.

8.3 Experiments

The pretraining used a small set of CT images and human face images only. For the CT images, we used the AAPM 2016 CT challenge data from [135]. We preprocessed the data [82] to get a total of 2304 training slices of size 256×256 . For the natural images, we used the CelebA-HQ dataset [126] with each image having size 256×256 . We then randomly selected a small subset of 3000 CelebA images for pre-training. We trained the patch-based network with patch sizes from 16×16 to 64×64 , and used a zero padding value of 64. The network architecture follows from [93] for both the patch-based networks and whole-image networks. All networks were trained on PyTorch using the Adam optimizer with two NVIDIA A40 GPUs.

TBL 8.1 – Comparison of quantitative results on four different medical imaging inverse problems. Results are averages across all images in the test dataset. Best results for practical use are in bold.

Method	PBCT, 60 Views		FBCT, 40 Views		512 × 512 CT		CS-MRI, 7×	
	PSNR↑	SSIM↑	PSNR↑	SSIM↑	PSNR↑	SSIM↑	PSNR↑	SSIM↑
Baseline	30.15	0.784	17.86	0.381	28.33	0.700	33.94	0.894
ADMM-TV	31.14	0.862	24.20	0.628	29.36	0.788	36.74	0.924
PnP-ADMM [196]	36.75	0.932	28.86	0.747	37.48	0.910	35.77	0.907
PaDIS+FC [82]	39.16	0.942	27.91	0.796	33.11	0.831	35.17	0.904
SCD [10]	41.16	0.962	21.28	0.463	–	–	–	–
Ours (SPAR)	42.72	0.972	36.11	0.918	38.81	0.929	39.15	0.949
Ideal*	42.82	0.973	36.34	0.923	38.94	0.930	39.42	0.953

*not available in practice with a single diffusion model

TBL 8.2 – Comparison of quantitative results on different natural imaging inverse problems: deblurring and super-resolution. The second column indicates whether the method can handle images from different modalities. Best results for practical use are in bold.

Method	Diff. Modality	Deblur, LSUN		Deblur, FFHQ		Superresolution	
		PSNR↑	SSIM↑	PSNR↑	SSIM↑	PSNR↑	SSIM↑
Baseline	✓	22.03	0.602	24.10	0.649	25.42	0.742
ADMM-TV	✓	27.01	0.801	28.44	0.772	25.22	0.729
PnP-ADMM [196]	✓	28.28	0.847	30.79	0.829	27.29	0.796
DPIR [204]	×	28.96	0.857	–	–	28.98	0.825
DiffPIR [210]	×	28.11	0.798	–	–	27.68	0.796
DDRM [95]	×	27.19	0.803	–	–	28.60	0.811
PaDIS+FC [82]	✓	27.17	0.809	28.91	0.774	27.23	0.782
SCD [10]	×	28.57	0.842	–	–	26.15	0.635
Ours (SPAR)	✓	29.49	0.877	31.43	0.843	29.10	0.842
Ideal*	✓	29.57	0.879	31.55	0.848	29.26	0.848

*not available in practice with a single diffusion model.

To evaluate our proposed SPAR method, we ran experiments on a wide variety of datasets with many different imaging applications, particularly on images that are out

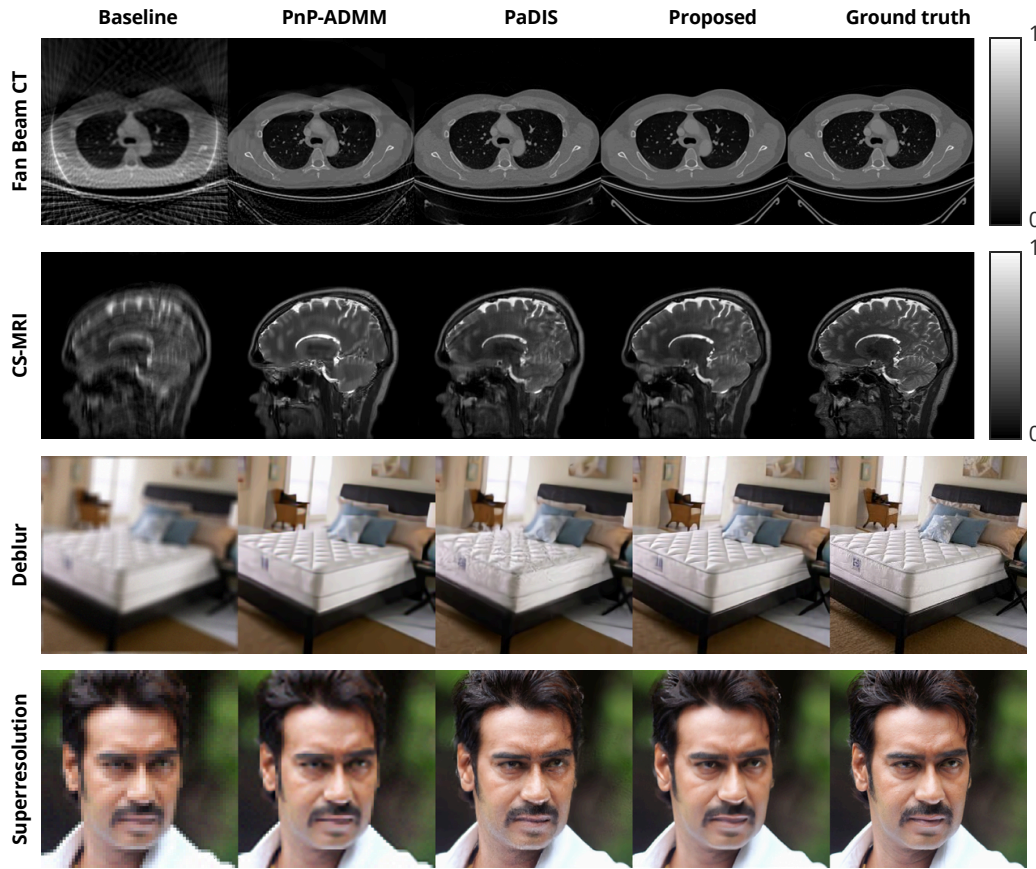


FIG 8.3 – Results of our proposed method and comparison methods for various inverse problems. MRI images are complex-valued so we show the pixel magnitudes.

of the training data distribution. We benchmarked the algorithm performance on: LIDC-IDRI dataset [6], fastMRI brain dataset from [100], the LSUN dataset [200], the FFHQ-512 dataset [94], in addition to the AAPM and CelebA datasets that we used for pre-training. For the LIDC-IDRI dataset, we first applied data preprocessing by setting the entire background of the volumes to zero. We rescaled the images in the XY-plane to have size 256×256 . For the fastMRI brain dataset, we took the complex-valued image space data and rescaled it so that the real parts had a maximum value of 1, and applied the same scale factor to the imaginary parts. We also resized the images and coil maps (provided in the dataset) to 256×256 . We also ran 512×512 CT experiments on the AAPM dataset by taking the original data from [135] without downsampling. For all experiments, we set aside 10 images in each dataset to use as test images and reported quantitative metrics that were averaged across these test images.

For CT imaging, we performed sparse-view reconstruction with both a parallel-beam projector (PBCT) and a fan-beam projector (FBCT), using 60 projection views in the former and 40 views in the latter. In both cases, we used a detector with size 512 pixels, with the forward and backward projectors implemented by the ODL package [173]. We also performed PBCT reconstruction experiments on the 512×512 sized CT images with 60 projection views. The PBCT experiments were performed with the AAPM dataset and the FBCT experiments were performed with the LIDC dataset.

For the MRI experiments, we performed compressed sensing MRI (CS-MRI) reconstruction. The sampling mask was chosen in a Cartesian pattern where each horizontal line within the center 10% of the image width was sampled, and all other horizontal lines were sampled with probability 0.03. Hence, the acceleration factor was $6.9\times$. We used 12 coils with the implementation for the forward and back projectors being the same as those in [75].

For the natural images, we ran deblurring and superresolution experiments. For the 256×256 images, we used a uniform blur kernel of size 9×9 and added white Gaussian noise with $\sigma = 0.01$ where the clean image was scaled between 0 and 1. We used the same settings but with a blur kernel of size 17×17 for the 512×512 images. For the superresolution experiments (only done on 256×256 images), we used a scaling factor of 4 with downsampling by averaging and added white Gaussian noise with $\sigma = 0.01$. Appendix A.5.5 contains additional experiments with noise level $\sigma = 0.05$.

For comparison, we applied the same patch-based diffusion network trained on the mixed dataset to solve the inverse problems, but without self-supervision. Hence, the checkpoint weights were held constant throughout the reconstruction process, similar to PaDIS [82]. We denote this method as PaDIS+FC (flexible channel) to clarify that in addition to [82], we applied (8.8) to handle images of different modalities. Previous work [81] showed that using conjugate gradient descent to enforce data fidelity yielded better results than DIS methods such as DDNM [181] and DPS [39]. Thus, we ran experiments using steerable controlled diffusion (SCD) [10] which also uses conjugate gradient descent; since this method does not support using a single network to reconstruct images of different modalities, we separately trained networks for it on the AAPM and CelebA datasets. For a more thorough comparison, we also compared with other SOTA methods for natural imaging: DPIR [204], DiffPIR [210], and DDRM [95]. Although these methods usually use a diffusion model trained for each dataset, for fair comparisons with SPAR we trained the networks on the CelebA dataset only.

We also compared our method with traditional image reconstruction methods: applying a non-iterative baseline, reconstructing using the total variation (TV) regularizer

TBL 8.3 – Comparison between diffusion inverse solvers using the same network for CT reconstruction and image deblurring.

Solver	FBCT (40 Views)		Deblur (LSUN)	
	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow
Langevin [166]	26.11	0.798	26.67	0.790
VE-DDNM [181]	18.05	0.424	26.40	0.777
VE-DPS [39]	23.03	0.714	28.36	0.847
Proposed (SPAR)	36.11	0.918	29.49	0.877

and solving the optimization problem using ADMM (ADMM-TV), and plug and play via ADMM (PnP-ADMM) [196]. The baseline methods for each of the inverse problems is as follows: applying the filtered back-projection (FBP) method to the measurement \mathbf{y} for CT, taking $A'\mathbf{y}$ for CS-MRI, taking the blurred image \mathbf{y} for deblurring, and upsampling the low resolution image \mathbf{y} using bilinear interpolation for superresolution. The implementation of ADMM-TV is in [74]. Finally, since our proposed method does not require a clean dataset for each type of image, we used the off-the-shelf denoiser BM3D [49] for PnP-ADMM. Appendix A.5.6 reports the values of all the parameters of the algorithms.

Tables 8.1 and 8.2 show the quantitative results of the medical and natural image experiments, respectively. In Table 8.1, the second column indicates whether or not the method directly supports using a single network to handle images of different modalities. For reference, we also show the results of the ideal case, where we trained a patch-based diffusion model on a large clean dataset of the same distribution as the test images and still applied the majorizer loss during reconstruction. Our proposed method outperformed the state-of-the-art benchmarks that do have access to a large training dataset in all the inverse problems across different modalities. Notably, our SPAR method performed well in CS-MRI reconstruction even though the diffusion network was not trained on any complex-valued or 2-channel images. Figure 8.3 shows the visual results of the experiments. Appendix A.5.5 reports LPIPS scores and runtimes for the comparison methods.

We further investigated the benefit of using the majorizer loss function (8.7). We ran a “simplified” version of Algorithm 6 with $\sigma' = 0$ (i.e., always using the loss in (7.4)) and evaluated the PSNR of the intermediate image throughout the CS-MRI reconstruction process. Figure 8.4 shows the results of this experiment compared with using the proposed Algorithm 6 that uses the majorizer loss (8.7). Both reconstruction algorithm versions consist of running the diffusion process in reverse, so the image quality improves

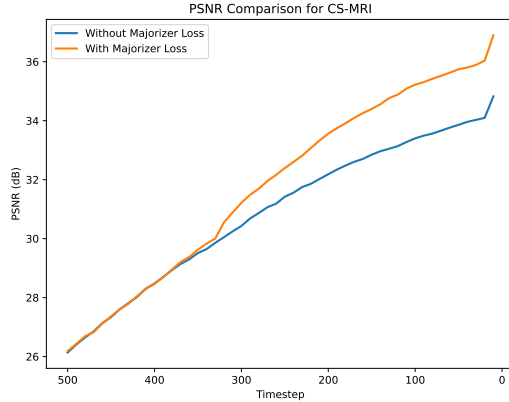


FIG 8.4 – PSNR over time for different loss types in CS-MRI.

as the timestep decreases. Although the diffusion process has 1000 steps, we truncated the graph at $t = 500$ to better illustrate the divergence in PSNR between the two different self-supervised loss functions. When using the simpler loss (7.4), due to inaccuracies in the learned score function, when the timestep (and hence the noise level of the image) becomes small, the improvement in image quality is lower compared to using the proposed majorizer loss (8.7). This experimentally reinforces the theoretical findings shown in Section 8.2.4 that the majorizer loss is less sensitive to inaccuracies in the learned score function.

We also investigated various other diffusion inverse solvers that leverage an unconditional prior. These methods assume that a diffusion model has been trained properly on a large sample of training images that belong to the same distribution as the test data. Therefore, subpar performance is expected when the test images belong to a different distribution from the training data. Since we trained our patch-based diffusion model on a mixture of AAPM CT images and CelebA images, we performed experiments on LIDC CT images and LSUN images to illustrate this point. In particular, we ran fan-beam CT with 40 projection views using the LIDC dataset and deblurring on the LSUN dataset.

Table 8.3 shows the results of these experiments. We compared Langevin dynamics with gradient descent for data fidelity [166], variance exploding DDNM [181] (as opposed to variance preserving, which would have required a differently trained network), and variance exploding DPS [39]. These results show that since the other inverse solvers do not account for a distribution shift, the results are substantially worse than our proposed method. Appendix A.5.3 contains more details as well as visualizations of these experiments. Finally, Appendix A.5.5 contains a discussion with some results about applying acceleration methods in conjunction with our method.

8.4 Conclusion

This chapter proposes SPAR (Single Pretrain Always Refine), a method that significantly reduces the pre-training overhead by leveraging adaptive self-supervision to refine the network at reconstruction time in many modalities. This flexibility allows SPAR to tackle inverse problems even in modalities absent from the initial training dataset. Extensive experiments demonstrated that using a single pretrained network, our SPAR method achieved consistently better image quality than existing methods. Limitations of this work include a lack of demonstration on nonlinear inverse problems as well as 3D and video inverse problems. The memory efficiency of patch-based modeling should facilitate such extensions. Additionally, for inverse problems such as inpainting with a large missing area, the proposed method may not be effective, as the pretrained image prior does not have any information from the measurement to guide the reconstruction in the missing area. This is in contrast to inpainting with random missing pixels, for which the proposed method should be more effective. Image priors such as the ones used in this paper have the potential to benefit society by reducing the scan time of MRI scans and the dose of CT scans.

CHAPTER 9

Future Work

Diffusion models have shown outstanding performance for various inverse problems spanning several modalities, but are limited by their computational and data requirements. This thesis covered several works that improve the flexibility of diffusion models and allow them to overcome these shortcomings. This chapter outlines some future directions and potential projects that build off this thesis. These projects advance our overall goal of developing diffusion models that can be used to solve large-scale image reconstruction problems, possibly with limited training data.

9.1 Improving Existing Works

9.1.1 Generalizing Patch-Based Diffusion Models

Chapter 4 described a general framework for patch-based diffusion models via (4.4) and proved that under certain conditions, a stochastic version of Langevin dynamics can be applied to sample from the prior. We also identified some previous works, as well as our own works which use prior models that fall under this framework. However, a few questions remain to be answered which could be the subject of future research.

- It was shown that regardless of the method by which patches are used to tile the whole image, provided the size of the patch is unchanging, the same method of training the network via performing denoising score matching on individual patches can be applied. However, different patch tiling methods do affect reconstruction results even when the same network is used, as shown in Chapter 4. A natural question then arises: which patch tiling methods are optimal and how does the application as well as diffusion inverse solver affect the results?

- The stochastic sampling theorem applies only for Langevin dynamics as the number of sampling steps tends to infinity. Are there analogous results for other sampling methods such as DDPM, DDIM, or other acceleration methods?
- Even when global information is utilized as in Chapter 6, for complex datasets such as natural imaging datasets, unconditional image generation can often lead to unrealistic results. This suggests the existence of a gap between the true whole-image prior and the patch-based image prior. How can the patch-based prior be improved to better reflect real-world image priors?
- More generally, one can not only learn the unconditional priors of patches but a prior of patches conditioned on other randomly chosen patches. In this case, the training process would consist of choosing a patch at random to denoise, while conditioning on some number of other randomly selected patches. This process could also help the model learn global context, reduce overfitting, and improve robustness.

9.1.2 Test-time Adaptation for 3D Problems

Chapters 7 and 8 tackled the problem of refining a diffusion network according to a measurement at reconstruction time for 2D images. Chung et al. [47] applied a similar method as Chapter 7 to perform test-time adaptation for 3D CT reconstruction. Since parallel-beam CT was used, different z -slices of the volume do not interfere with each other, and it is possible to choose slices of the volume (instead of the entire volume) and their corresponding CT projections to perform test-time adaptation. However, for more complicated forward models, such as cone-beam CT or spatiotemporal deblurring in dynamic imaging, it would be necessary to perform network backpropagation with the entire volume, which is computationally infeasible using either the loss of [47] or Chapter 8. Future work could consist of developing an alternative way to refine the network during reconstruction time which does not require backpropagating through the entire forward model and volume.

Separately, a natural question that arises during test-time adaptation is the characteristics of the network after reconstruction and hence finetuning. In the 2D setting, we have observed that for CT reconstruction, this finetuned network does not have the ability to unconditionally generate meaningful images. It may not be reasonable to expect the network to generate a fully diverse set of images from the measurement of only a single in-distribution image; nevertheless, [164] shows that under deterministic sampling via DDIM, small perturbations to the initial noise result in small perturbations in the generated image which still result in a realistic image. Hence, it may be possible to obtain a refined network that can regenerate the clean image and local variants of it.

9.1.3 Diffusion Bridge Network Refining

Diffusion bridges are a generalization of the traditional diffusion model framework. While the traditional diffusion model framework involves learning a neural network that can transform an image from the pure white Gaussian noise distribution to some target data distribution [73, 166], diffusion bridges use a network to transform an image between two arbitrary distributions. We use p_0 to denote the distribution of images degraded in some fashion (i.e., blurry images, low resolution images, etc.) and p_1 to denote the clean image distribution [120, 40]. Hence, the forward process consists of adding noise in a controlled fashion to go from the clean image to the degraded image, while the reverse process consists of recovering the clean image from the degraded image.

When using diffusion bridges to solve inverse problems, a separate network must be trained for each inverse problem. For example, for deblurring, the degraded image distribution consists of blurry images, and a network trained for this cannot then be applied to perform image superresolution. This is a disadvantage of diffusion bridge methods compared to traditional diffusion model approaches, where the unconditionally trained image prior can be used to solve any inverse problem on that image dataset. The advantage of diffusion bridge methods is that since they start from a distribution that is closer to the clean image than pure noise, the number of NFEs (neural function evaluations) needed at reconstruction time is much lower than traditional diffusion models. In practice, competitive performance has been observed using only around 10 NFEs with diffusion bridges [120] compared to 100-1000 for traditional diffusion models. Furthermore, the number of NFEs used can be freely changed at reconstruction time, and leads to an observable tradeoff between FID and PSNR that pushes the Pareto boundary [40].

Future work could consist of one possible solution that achieves the best of both worlds: solving many different inverse problems using one pretrained neural network in a small number of sampling steps. We achieve this by starting with a diffusion bridge network that is trained specifically for some inverse problem. Next, in a similar way as Chapter 7, if the inverse problem or the dataset differs from that of the pretrained network, we include a data consistency loss function to adapt the network to the different dataset and inverse problem. Table 9.1 summarizes properties of different DIS methods. Notably, this method would only require training a network for each modality. If we were to combine the ideas of patch-based diffusion models and channel operators from SPAR, it could even be possible to use a single neural network in conjunction with a diffusion bridge. Furthermore, the number of required diffusion iterations would be far smaller than previous unconditional DIS methods, significantly improving the reconstruction speed.

	Training	Patches?	Sampling speed
DPS/DDNM/standard DIS	For each dataset	Optional	Slow (2 min)
SPAR	Once	Yes	Very slow (3 min)
Diffusion bridge: I2SB, CDDB	For each inverse problem	No	Fast (5 sec)
Bridge + refine	For each modality	Optional	Somewhat fast?

TBL 9.1 – Comparison of different diffusion-based methods

9.1.4 Provable Posterior Sampling

There are various diffusion inverse solvers that leverage an unconditionally trained prior model; Chapter 4 uses DPS [39] and Chapters 5 and 6 use the conjugate gradient descent method of [41]. However, all of these methods only *approximately* sample from the posterior $p(\mathbf{y}|\mathbf{x})$ and in particular, approximate the intractable $\nabla \log p(\mathbf{y}|\mathbf{x}_t)$ during the reconstruction process. Methods that can perform exact posterior sampling have various limitations: Monte Carlo methods [32, 36] are only exact when the number of particles used tends to infinity, and are thus computationally expensive. Diffusion bridges [120, 40] provably sample from the posterior in fewer steps than traditional diffusion models, but require training a network specific to each inverse problem, limiting their flexibility. Finally, variational methods [60, 8] directly train a second network (in addition to the unconditional prior network) to approximate the posterior distribution and use this second network to quickly sample from the posterior, but the additional training step is costly. Future work could involve developing a method that still only uses a single unconditional diffusion model, yet can provably sample from the posterior reasonably quickly.

9.2 Leveraging Video Diffusion Models

Video diffusion models have recently risen to popularity, showing the ability to generate high quality videos [129, 124]. By treating the temporal dimension of videos as the third spatial dimension for 3D imaging problems, one can hope to use pretrained video diffusion models (possibly with finetuning) to solve 3D inverse problems. However, there are two primary difficulties with this approach. Firstly, video diffusion models generally use a network that consists of two modules: a spatial module that batchifies the temporal dimension and processes the spatial dimensions using a typical diffusion UNet or transformer, and a temporal module that batchifies the two spatial dimensions but applies an

attention layer in the temporal dimension [134, 202]. As a result, the network must process the entire video at once, which constrains the number of frames in the videos, usually to fewer than 40 even with modern GPUs. This is prohibitively limiting for CT and MRI reconstruction, where the third spatial dimension can often have size 256 or greater. Secondly, video generation papers frequently do not have a simple black box network whose input is solely the noisy video and whose output is the denoised video or its score function [160, 19], instead relying on a conditioning mechanism on lower resolution videos or previous frames to guide the generation process. This makes it less straightforward to leverage these models in conjunction with diffusion inverse solvers.

Existing works on video inverse problems often instead leverage a 2D diffusion model, but add similar noise across different frames during the reconstruction process so as to preserve temporal continuity [102, 103]. These methods generally outperform 3D image reconstruction methods, as the pretrained 2D diffusion models are expressive and powerful, capturing a diverse range of modes in natural videos. However, these methods still do not learn a data driven prior that captures the relationship between different frames. Bai et al. [8] uses a pretrained video prior and trains a network approximating the posterior distribution for a specific inverse problem using a variational approach. Furthermore, the frames are processed in batches with a shifting window, which allows the network to learn a prior across multiple frames while also bypassing the memory requirements of processing the entire video at once. Although this method is able to perform reconstruction of videos of 30 frames in only a second, it requires training separate posterior networks for each inverse problem. Future work would consist of developing a method that can achieve the following objectives:

- Leverage a pretrained video diffusion model, such as stable video diffusion [19], with only a relatively cheap network finetuning step for each dataset
- Learn correlations between different slices in a data-driven way
- Process volumes of varying numbers of slices in the z direction
- Have a direct black box representation of the score of the volume which can then be applied in conjunction with a wide range of diffusion inverse solvers

APPENDIX A

Appendix

This section contains the appendices for Chapters 3-8.

A.1 Appendix for "Learning Image Priors through Patch-based Diffusion Models for Solving Inverse Problems"

A.1.1 Ablation studies

We performed four ablation studies to evaluate the impact of different parameters on the performance of our proposed method. Similar to Table 4.1, we ran the experiments on all the images in the test dataset and computed the average metric. Section A.1.2 shows visualizations of these studies.

Effect of patch size. We investigated the effect of the patch size P used at reconstruction time for the 20-view CT reconstruction problem. We continued to augment the training with smaller patch sizes when possible so as to be consistent with the main experiments (patch size of 56 but also trained with patch sizes of 32 and 16), while using the same neural network architecture. Different amounts of zero padding were needed for each of the experiments per (4.1). App. A.4.3 provides the full details. At reconstruction time, the same patch size was used throughout the entire algorithm. Using a “patch size” of 256 corresponds to training a diffusion model on the whole image (without zero padding).

Table A.1 shows that careful selection of the patch size is required to obtain the best results for a given training set size. If the patch size is too small, the network has trouble capturing global information across the image. Although the positional information helps in this regard, there may be some inconsistencies between patches, so the learned image prior is suboptimal although the patch priors may be learned well. At the other extreme, very large patch sizes and the whole image diffusion model require more memory to train and run. The image quality drops in this case as limited training data prevents the network from learning the patch prior well.

TBL A.1 – Effect of patch size P on CT reconstruction

P	PSNR \uparrow	SSIM \uparrow
8	32.57	0.844
16	32.57	0.829
32	32.72	0.853
56	33.57	0.854
96	33.36	0.854
256	32.84	0.835

TBL A.2 – Dataset size effect on CT reconstruction

Dataset size	Patches		Whole image	
	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow
144	32.28	0.841	29.12	0.804
288	32.43	0.837	31.09	0.829
576	33.03	0.846	31.81	0.835
1152	33.01	0.849	31.36	0.834
2304	33.57	0.854	32.84	0.835

Effect of training dataset size. A key motivation of this work is large-scale inverse problems having limited training data. To investigate the effects of using small datasets on our proposed method, compared to standard whole image models, we trained networks on random subsets of the CT dataset. Table A.2 summarizes the results. Crucially, although the reconstruction quality tends to drop as the dataset size decreased for both the patch-based model and the whole image model, the drop is much more sharp and noticeable for the whole image model, particularly when the dataset is very small. This behavior is consistent with the observations of previous works where large datasets consisting of many thousands of images were used to train traditional diffusion models from scratch.

Effect of positional encoding. High quality image generation via patch-based models that lack positional encoding information would be impossible, as no global information about the image could be learned at all. We demonstrate that positional information is also crucial for solving inverse problems with patch-based models. We examined the results of performing CT reconstruction for trained networks *without* positional encoding as an input compared to networks *with* positional encoding. According to [45], when solving inverse problems in some settings, it can be beneficial to initialize the image with some baseline image instead of with pure noise (as is traditionally done). To allow the network that did not learn positional information to possibly use a better initialization with patches roughly in the correct positions, we also ran experiments by initializing with the baseline. Table A.3 shows that in both cases, the network completely failed to learn the patch-based prior and the reconstructed results were very low quality. Hence, positional information is crucial to learning the whole image prior well.

Sampling methods. One benefit of our proposed method is it provides a black box image prior for the entire image that can be computed purely through neural network operations on image patches. We demonstrate the versatility of this method by pairing a variety of different sampling and inverse problem solving algorithms with our patch-based image prior, along with comparisons with a whole-image prior. The implemented sampling methods include Langevin dynamics [166] with a gradient descent term for enforcing data fidelity step and the predictor-corrector method for solving SDEs [46]. Since we observed better stability and results with Langevin dynamics, we also combined this sampling method with nullspace methods that rely on hard constraints [181] and DPS [39]. To use the same neural network checkpoint across these implementations, we used the variance exploding SDE [168] method as the backbone for both training and reconstruction. DPS [39] and DDNM [181] were originally implemented with networks trained under

the VP-SDE framework; here, we implemented those methods with the VE-SDE framework. Table A.4 shows that generally, VE-DPS performed the best and that the patch-based method consistently outperformed the whole image method. However, the patch-based method still obtained reasonable results for all the implemented methods, showing that the learned image prior is indeed flexible enough to be paired with a variety of sampling algorithms. App. A.4.3 provides more details about the implemented algorithms.

TBL A.3 – Positional encoding effect for CT reconstruction

	PSNR↑	SSIM↑
no position enc.	23.25	0.459
no position+init	24.51	0.518
with position	33.57	0.854

TBL A.4 – Dataset size effect on CT reconstruction

Method	Patch-based		Whole image	
Metric	PSNR↑	SSIM↑	PSNR↑	SSIM↑
Langevin dynamics	33.03	0.846	30.92	0.813
Predictor-corrector	32.35	0.820	18.95	0.149
VE-DDNM	31.98	0.861	29.49	0.830
VE-DPS	33.57	0.854	32.84	0.835

A.1.2 Ablation study images

Figure A.1 shows the results of applying PaDIS to two example test images with different patch sizes. The main results, i.e., those shown in Table 4.1, used $P = 56$. For some of the other patch sizes, some artifacts can be seen in the images. Namely, the smooth parts of the image become riddled with "fake" features for small patch sizes and some of the sharp features become more blurred. The fake features in the right half of the image in the top row are especially apparent when applying the whole-image model. The runtime for different patch sizes were fairly similar, with $P = 8$ taking notably longer than the others due to the large number of patches required. The image size for these experiments was small enough so that the score function of all the patches could be computed in parallel; however, for larger scale problems such as high resolution 2D images or 3D images, large patch sizes become infeasible due to memory constraints.

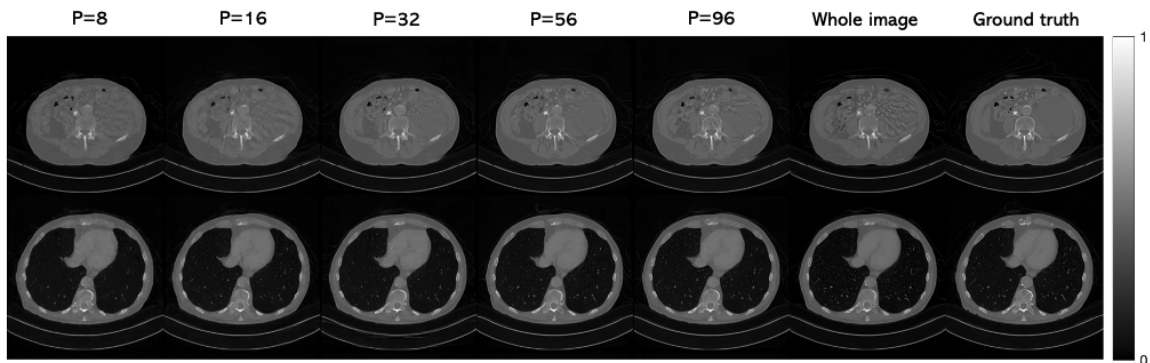


FIG A.1 – Results of PaDIS for 20 view CT reconstruction with different sized patches.

Figure A.2 shows the results of applying our proposed method and the whole image diffusion model to 20-view CT reconstruction for varying sizes of the training dataset. The image quality for PaDIS remains visually consistent as the size of the training dataset shrinks, as each image contains thousands of patches which helps avoid overfitting and memorization. However, the drop in quality for the whole image model is much more visible: in particular, the sharp features of the image are lost and the image becomes blurry. Hence, for applications where data is even more limited, such as medical imaging, our method can potentially have a greater benefit.

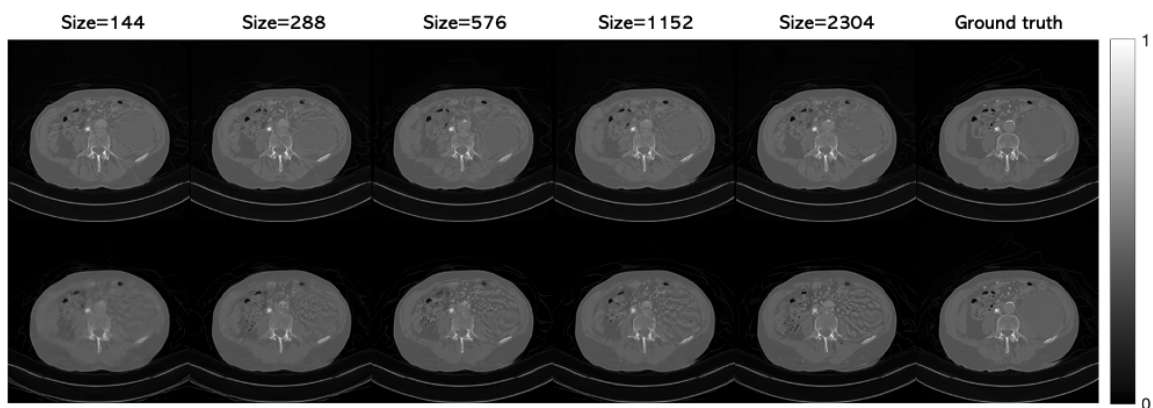


FIG A.2 – Results for 20 view CT reconstruction with different dataset sizes. Top row shows recon performed by PaDIS; bottom row shows recon performed with the whole-image model.

Figure A.3 demonstrates the importance of adding positional encoding information into the patch-based network on two different images. When positional information is not included, the network simply learns a mixture of all patches, resulting in a very blurry image with many artifacts resulting from the data fidelity term. Even when a better initialization of the image is provided, the same blurriness remains.

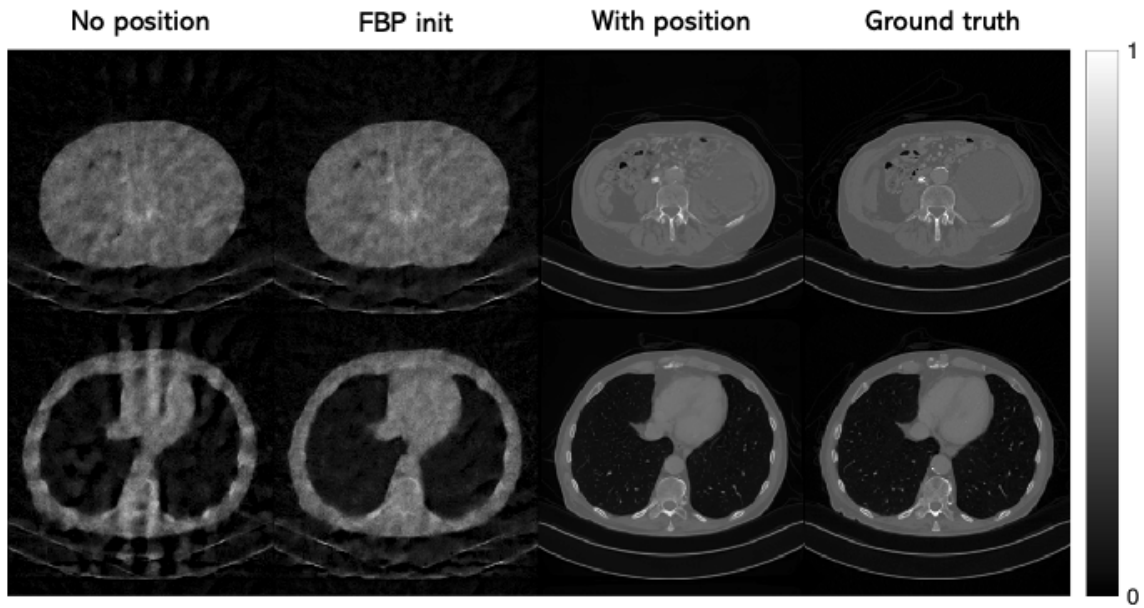


FIG A.3 – Results of PaDIS for 20 view CT reconstruction for different positional encoding methods.

Figure A.4 shows the results of using our proposed method compared with the whole-image diffusion model with different sampling and inverse problem solving algorithms. The predictor-corrector algorithm fails completely when using the whole-image model, indicating that this model could not be well-trained in this limited data setting. Quantitatively, DPS performs the best for PaDIS; visually, all of the methods obtain reasonable results, although some more minor artifacts are present in the first four methods. Nevertheless, this shows that the patch-based prior is flexible and can be used with a variety of existing algorithms.

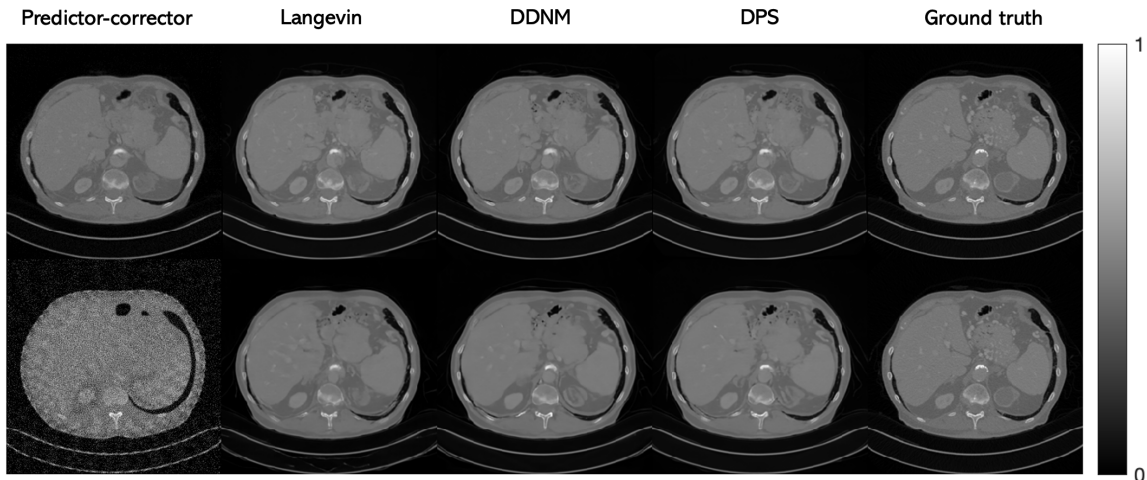


FIG A.4 – Results of PaDIS for 20 view CT reconstruction using different sampling and inverse problem solving algorithms. Top row is with PaDIS and bottom row is with the whole-image model.

A.1.3 Experiment parameters

We trained the patch-based networks and whole-image networks following [93]. Since images were scaled between 0 and 1, we chose a maximum noise level of $\sigma = 40$ and a minimum noise level of $\sigma = 0.002$. We used the same UNet architecture for all the patch-based networks consisting of a base channel multiplier size of 128 and 1, 2, 2, and 2 channels per resolution for the four 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 patches to avoid overfitting. Finally, the learning rate was chosen to be $2 \cdot 10^{-4}$ and the batch size for the main patch size was 128, although batch sizes of 256 and 512 were used for the two smaller patch sizes. The entire model had around 60 million weights. For the whole image model, we kept all the parameters the same, but increased the number of

channels per resolution in the fourth layer to 4 so that the model had around 110 million weights. The batch size in this case was 8.

For image generation and solving inverse problems, 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, as follows:

- CT with 20 views: $\sigma_{\max} = 10, \sigma_{\min} = 0.002$
- CT with 8 views: $\sigma_{\max} = 10, \sigma_{\min} = 0.003$
- Deblurring: $\sigma_{\max} = 40, \sigma_{\min} = 0.005$
- Superresolution: $\sigma_{\max} = 40, \sigma_{\min} = 0.01$.

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

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

where $\operatorname{TV}(\mathbf{x})$ represents the L1 norm total variation of \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.

Ablation study details. For each patch size, we trained with the main patch size along with smaller patches whenever possible. However, since we did not modify the network architecture, and the architecture consists of downsampling the image 3 times by a factor of 2, it was necessary for the input dimension to be a multiple of 8. Furthermore, we followed a patch scheduling method similar to that of the main experiments unless otherwise noted. Finally, to avoid excessive zero padding, for larger patch sizes, we used patch sizes that were smaller than the next power of 2 such that the main image could still be fully covered by the same number of patches. The details are as follows.

- $P = 8$: This was trained only with this patch size as no smaller sizes could be used.
- $P = 16$: Trained with patch sizes of 8 and 16 with probabilities of 0.3 and 0.7 respectively.
- $P = 32$: Trained with patch sizes of 8, 16, and 32 with probabilities of 0.2, 0.3, and 0.5 respectively.
- $P = 56$: Trained with patch sizes of 16, 32, and 56 with probabilities of 0.2, 0.3, and 0.5 respectively. Zero padding width was set to $5 \cdot 56 - 256 = 24$.

- $P = 96$: Trained with patch sizes of 32, 64, and 96 with probabilities of 0.2, 0.3, and 0.5 respectively. Zero padding width was set to $3 \cdot 96 - 256 = 32$.

A.1.4 Comparison algorithms

This section provides pseudocode for the implemented alternative sampling algorithms whose results are shown in Table A.4. Here, for brevity, we show the versions using the whole-image diffusion model; the versions with our proposed method are readily implemented by computing $\mathbf{s} = \mathbf{s}(\mathbf{x}, \sigma_i)$ through the procedure illustrated in Alg. 2.

Algorithm 7 Image Recon via Langevin Dynamics

Require: $\sigma_1 < \sigma_2 < \dots < \sigma_T, \epsilon > 0, \zeta_i > 0, \mathbf{y}$
Initialize $\mathbf{x} \sim \mathcal{N}(0, \sigma_T^2 \mathbf{I})$
for $i = T : 1$ **do**
 Sample $\mathbf{z} \sim \mathcal{N}(0, \sigma_i^2 \mathbf{I})$
 Set $\alpha_i = \epsilon \cdot \sigma_i^2$
 Apply neural network to get $D = D_\theta(\mathbf{x}, \sigma_i)$
 Set $\mathbf{s} = (D - \mathbf{x}) / \sigma_i^2$
 Set \mathbf{x} to $\mathbf{x} + \zeta_i \mathcal{A}^T(\mathbf{y} - \mathcal{A}(\mathbf{x}))$
 Set \mathbf{x} to $\mathbf{x} + \frac{\alpha_i}{2} \mathbf{s} + \sqrt{\alpha_i} \mathbf{z}$
end for
Return \mathbf{x} .

Algorithm 8 Image Recon via Predictor-Corrector Sampling

Require: $\sigma_1 < \sigma_2 < \dots < \sigma_T, \epsilon > 0, \zeta_i > 0, r, \mathbf{y}$
Initialize $\mathbf{x} \sim \mathcal{N}(0, \sigma_T^2 \mathbf{I})$
for $i = T : 1$ **do**
 Set \mathbf{x} to $\mathbf{x} + (\sigma_{i+1}^2 - \sigma_i^2) \mathbf{s}_\theta(\mathbf{x}, \sigma_{i+1})$
 Set \mathbf{x} to $\mathbf{x} + \zeta_i \mathcal{A}^T(\mathbf{y} - \mathcal{A}(\mathbf{x}))$
 Sample $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$
 Set \mathbf{x} to $\mathbf{x} + \sqrt{\sigma_{i+1}^2 - \sigma_i^2} \mathbf{z}$
 Sample $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$
 Set $\epsilon_i = 2r \frac{\|\mathbf{z}\|_2}{\|\mathbf{s}_\theta(\mathbf{x}, \sigma_i)\|_2}$
 Set $\mathbf{s} = \mathbf{s}_\theta(\mathbf{x}, \sigma_i)$
 Set \mathbf{x} to $\mathbf{x} + \epsilon_i \mathbf{s} + \sqrt{2\epsilon_i} \mathbf{z}$
 Set \mathbf{x} to $\mathbf{x} + \zeta_i \mathcal{A}^T(\mathbf{y} - \mathcal{A}(\mathbf{x}))$
end for
Return \mathbf{x} .

Algorithm 9 DDNM

Require: $\sigma_1 < \sigma_2 < \dots < \sigma_T, \epsilon > 0, \zeta_i > 0, \mathbf{y}$ Initialize $\mathbf{x} \sim \mathcal{N}(0, \sigma_T^2 \mathbf{I})$ **for** $i = T : 1$ **do** Sample $z \sim \mathcal{N}(0, \sigma_i^2 \mathbf{I})$ Set $\alpha_i = \epsilon \cdot \sigma_i^2$ Apply neural network to get $D = D_\theta(\mathbf{x}, \sigma_i)$ Set $D = \mathcal{A}^\dagger \mathbf{y} + D - \mathcal{A}^\dagger \mathcal{A}(D)$ Set $\mathbf{s} = (D - \mathbf{x}) / \sigma_i^2$ Set \mathbf{x} to $\mathbf{x} + \frac{\alpha_i}{2} \mathbf{s} + \sqrt{\alpha_i} z$ **end for**Return \mathbf{x} .

In all cases, we used the same noise schedule as the main 20 view CT reconstruction experiment. For Langevin dynamics and DDNM, we set $\epsilon = 1$; the final results were not sensitive with respect to this parameter. For Langevin dynamics and predictor-corrector sampling, we took $\zeta_i = 0.3 / \|\mathbf{y} - \mathcal{A}(\mathbf{x})\|_2$, similar to the step size selection of DPS. Following the work of [46], we chose $r = 0.16$ for PC-sampling. The same parameters were used for the patch-based and whole image methods.

Table A.5 shows the average runtimes of each of the implemented methods when averaged across the test dataset for 20 view CT reconstruction.

TBL A.5 – Average runtimes of different methods across images in the test dataset for 20 view CT recon.

Method	Runtime (s) ↓
Baseline	0.1
ADMM-TV	1
PnP-ADMM	8
PnP-RED	22
Whole image diffusion	172
Langevin dynamics	98
Predictor-corrector	189
VE-DDNM	105
PaDIS (VE-DPS)	195

A.1.5 Markov random field interpretation

Markov random fields (MRF) are a tool used to represent certain image distributions and are particularly applicable to patch-based diffusion models. Describing the connection between MRF and this work requires some notation: let $\mathbf{x} = \{x_s : s \in \mathcal{S}\}$ denote the random field, where the index \mathcal{S} denotes the sites. The neighborhood system \mathcal{N} is defined as $\mathcal{N} = \{\mathcal{N}_s : s \in \mathcal{S}\}$. A model for $\mathbf{x} \in \mathcal{X}$ is a MRF on \mathcal{S} with respect to the neighborhood system \mathcal{N} if

$$p(x_s | x_{\mathcal{S}-\{s\}}) = p(x_s | x_{\mathcal{N}_s}), \forall \mathbf{x} \in \mathcal{X}, \forall s \in \mathcal{S}. \quad (\text{A.2})$$

Therefore, the distribution of each site (normally chosen to be a pixel) conditioned on the rest of the pixels depends only on the neighboring pixels.

By the Hammersley-Clifford theorem [16], such a MRF satisfying $p(\mathbf{x}) > 0$ everywhere can also be rewritten as $p(\mathbf{x}) = \frac{1}{Z} e^{-U(\mathbf{x})}$, where Z is a normalizing constant. In this case $U(\mathbf{x})$ is called the energy function and has the form $U(\mathbf{x}) = \sum_{c \in \mathcal{C}} V_c(\mathbf{x})$, which is a sum of clique potentials $V_c(\mathbf{x})$ over all all possible cliques. Thus, the score function for a MRF model is:

$$\mathbf{s}(\mathbf{x}) = \nabla \log p(\mathbf{x}) = -\nabla U(\mathbf{x}) = -\sum_c \nabla V_c(\mathbf{x}). \quad (\text{A.3})$$

If we let the neighborhood system be the patches of the image, then V_c corresponds to the clique potential for the c th patch of an image, and $-\nabla V_c(\mathbf{x})$ denotes the score function of that patch. Denoting by \mathbf{G}_c the wide binary matrix that extracts the pixels corresponding to the c th patch from the whole image, we define $V_c(\mathbf{x}) = V(\mathbf{G}_c \mathbf{x}, c^*)$, where c^* denotes the positional encoding method used for the c th patch, and now we simply have one (patch) clique function V . Finally, the overall score function under this model becomes

$$\mathbf{s}(\mathbf{x}) = \sum_c \mathbf{G}'_c \mathbf{s}_V(\mathbf{G}_c \mathbf{x}, c^*), \quad (\text{A.4})$$

where $\mathbf{s}_V(\mathbf{v}, c^*) \triangleq -\nabla_{\mathbf{v}} V(\mathbf{v}, c^*)$ is the shared score function of each of the patches with a positional encoding input.

In this work, we approximate \mathbf{s}_V with a neural network parameterized by θ , and we use denoising score matching to train the network via the loss function

$$L(\theta) = \mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma_t^2 \mathbf{I})} \left\| \mathbf{s}_\theta(\mathbf{x} + \boldsymbol{\epsilon}, \sigma_t) + \boldsymbol{\epsilon} / \sigma_t^2 \right\|_2^2 \quad (\text{A.5})$$

$$= \mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma_t^2 \mathbf{I})} \left\| \sum_c \mathbf{G}'_c \mathbf{s}_V(\mathbf{G}_c \mathbf{x}, c^*; \theta) + \boldsymbol{\epsilon} / \sigma_t^2 \right\|_2^2. \quad (\text{A.6})$$

This derivation makes no assumptions on the patches; in particular, this method to train the score function would still hold if the patches overlapped for each iteration. However, such overlap would make it costly to train the network, as the loss function would need to be propagated through the sum over all patches every training iteration. We circumvent this problem by using non-overlapping patches (within a given reconstruction iteration; i.e., our approach only uses patches that “overlap” only across different iterations). For non-overlapping patches, the sum can be rewritten as follows:

$$L(\theta) = \mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma_t^2 \mathbf{I})} \left\| \sum_c \mathbf{G}'_c \mathbf{s}_V(\mathbf{G}_c \mathbf{x}, c^*; \theta) + \mathbf{G}'_c \mathbf{G}_c \boldsymbol{\epsilon} / \sigma_t^2 \right\|_2^2 \quad (\text{A.7})$$

$$= \mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma_t^2 \mathbf{I})} \mathbb{E}_{\text{random } c} \left\| \mathbf{s}_V(\mathbf{G}_c \mathbf{x}, c^*; \theta) + \mathbf{G}_c \boldsymbol{\epsilon} / \sigma_t^2 \right\|_2^2. \quad (\text{A.8})$$

This loss function is much easier to compute, as we can now randomly select patches and perform denoising score matching on the individual patches, as opposed to considering the entire image at once, and is equivalent to (A.5).

A.1.6 Acceleration methods

Although diffusion models are capable of generating high quality images, the iterative generative process typically requiring around 1000 neural function evaluations (NFEs) [73, 168] is a major disadvantage. In recent years, significant work has been done to improve sampling speed of diffusion models [165, 93, 128]. To reuse the same trained network, one can first derive an algebraic relationship between the score function $\mathbf{s}(\mathbf{x}, \sigma)$ (readily computed using the denoiser network $D(\mathbf{x}, \sigma)$ trained via (4.3)) and the residual function $\boldsymbol{\epsilon}(\mathbf{x}, t)$ which is approximated with a neural network in papers such as [73]. The score matching network $\mathbf{s}_\theta(\mathbf{x}, \theta)$ learns to map $\mathbf{x} + \sigma\boldsymbol{\epsilon}$ to \mathbf{x} , whereas the residual network learns to map $\mathbf{x}_t = \sqrt{\alpha_t}\mathbf{x}_0 + \sqrt{1 - \alpha_t}\boldsymbol{\epsilon}$ to $\boldsymbol{\epsilon}$ where $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I})$. Hence, we may input $\frac{\mathbf{x}_t}{\sqrt{\alpha_t}} = \mathbf{x}_0 + \frac{\sqrt{1 - \alpha_t}}{\sqrt{\alpha_t}}\boldsymbol{\epsilon}$ as the noisy image into the denoising score matching network so that the correct output becomes \mathbf{x}_0 . Then the corresponding noise level is $\sigma_t = \frac{\sqrt{1 - \alpha_t}}{\sqrt{\alpha_t}}$. Finally, the outputs of the network must be scaled via $\mathbf{s}_\theta(\mathbf{x}) = -\boldsymbol{\epsilon}(\mathbf{x})/\sigma$. Thus, by using this transformation, the network trained via (4.3) may also be applied to sampling algorithms requiring $\boldsymbol{\epsilon}_\theta$.

Using these ideas, we implemented the EDM sampler, a second-order solver for SDEs, according to [93], which can produce high fidelity images in 18 iterations (equating to 36 NFEs as each iteration requires two NFEs). We also implemented DDIM [165] using 50 sampling steps. Figure A.5 shows the results of using these methods with the proposed patch-based prior along with Langevin dynamics with 300 NFEs. The EDM sampler produces images that have clear boundary artifacts and the images from the DDIM method also have some discontinuous parts. This behavior is due to the stochastic method of computing the patch-based prior: according to Algorithm 2, we randomly choose integers i and j with which to partition the zero padded image and compute the score function according to this partition. Hence, accelerated sampling algorithms that attempt to remove large amounts of noise at each step tend to fare worse at removing boundary artifacts. This limitation of our method makes it difficult to run accelerated sampling algorithms, so is a direction for future research.

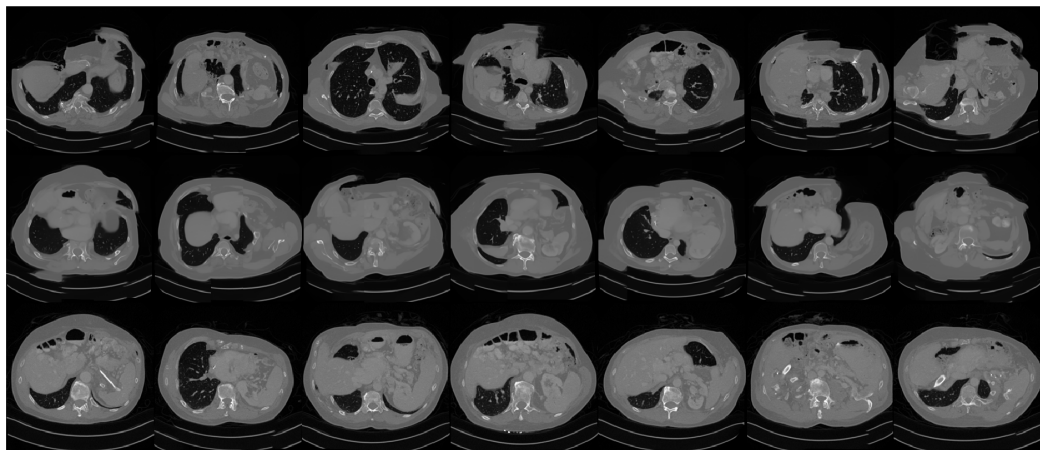


FIG A.5 – Generation of CT images with various acceleration algorithms. Top row shows generation with the EDM sampler [93], middle row uses DDIM [165], bottom row is our proposed method.

A.1.7 Additional figures

Figure A.6 shows the visual results of the extra inverse problems in Table 4.2.

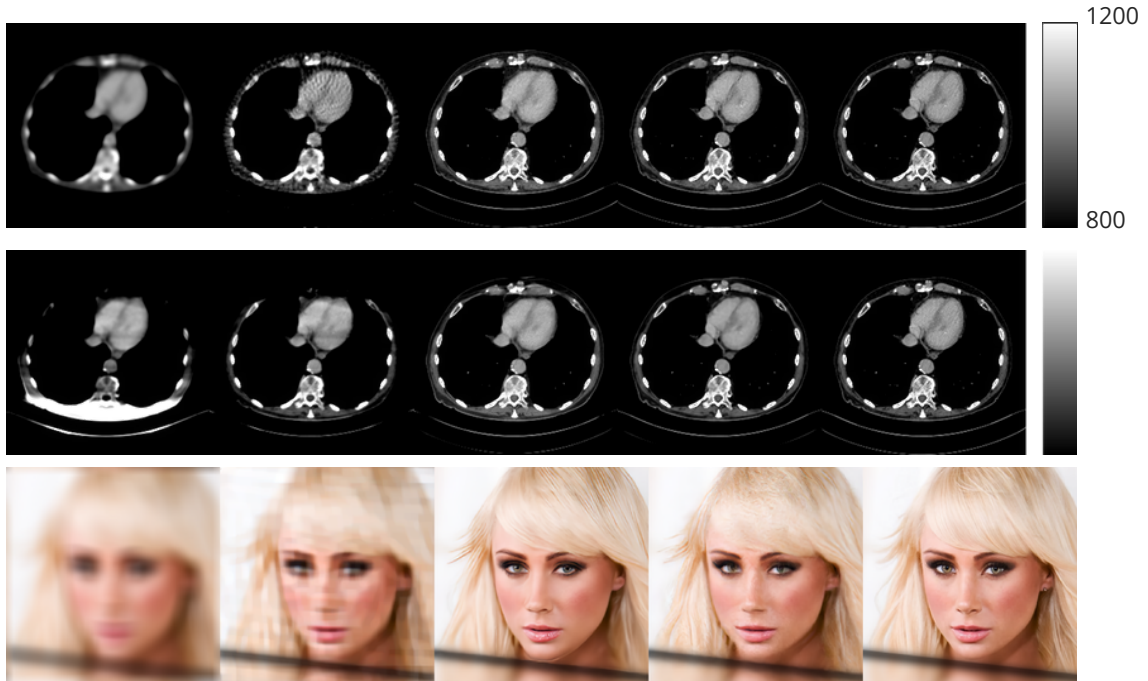


FIG A.6 – Results of extra inverse problem experiments. From top to bottom: 60 view CT, fan beam CT, heavy deblurring. From left to right: baseline, ADMM-TV, whole image diffusion, PaDIS, ground truth.

To further explore the different methods of assembling patches to form the whole image, we looked at unconditionally generated CT images according to the methods of [142] and [152]. Unlike with our proposed method, for these two methods, the patch locations are fixed throughout all of the timesteps. Furthermore, overlapping patches must be used to avoid boundary artifacts. The main difference is the way in which the methods handle the overlapping pixels: [152] *overrides* the overlapped areas with the new patch update while [142] *averages* over the overlapped area. We use the same network checkpoint trained on CT images as the main experiments and a patch size of 56 with overlap of 8 for the experiments. Figure A.7 shows the generated images: while both methods are able to avoid boundary artifacts, the overall structure of the generations are of lower quality than the images generated by our proposed method (see Table 4.4). This suggests that our proposed method most effectively combines the patch priors to form a prior for the entire image.

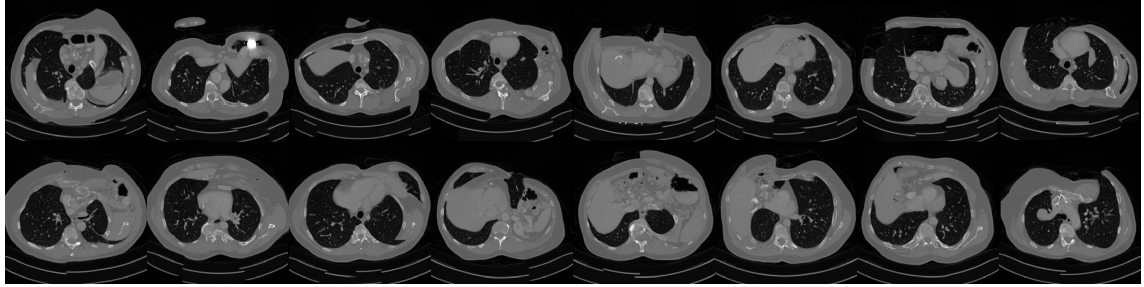


FIG A.7 – Unconditionally generated CT images using patch stitching [152] for the top row and patch averaging [142] for the bottom row. Compare this to the images generated by our proposed method in Figure 4.4.

Figures A.8 and A.9 show the PSNR of each individual image in the test dataset when using the whole-image model versus our proposed method. The plots show that our method exhibited reasonably consistent performance improvements over the test dataset.

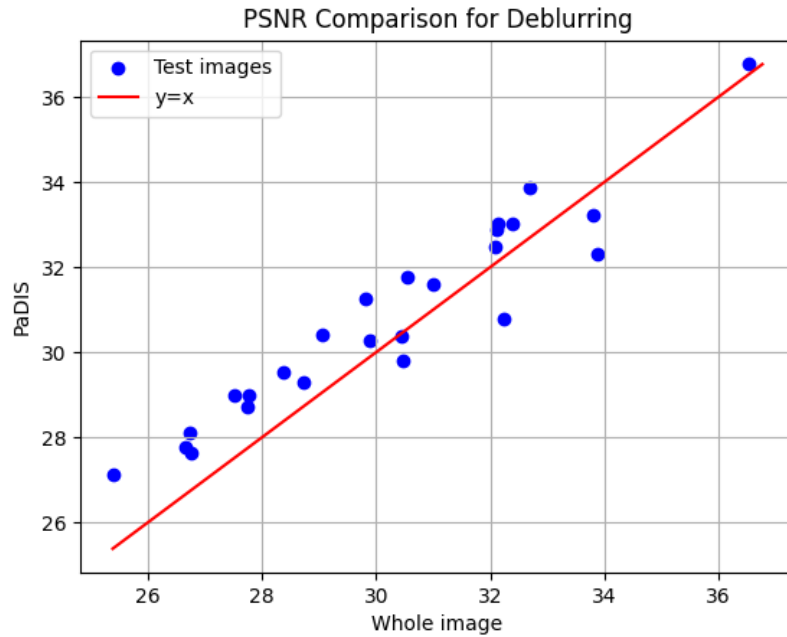


FIG A.8 – Comparison between PSNR of deblurring between whole-image model and proposed method for each image in the test dataset.

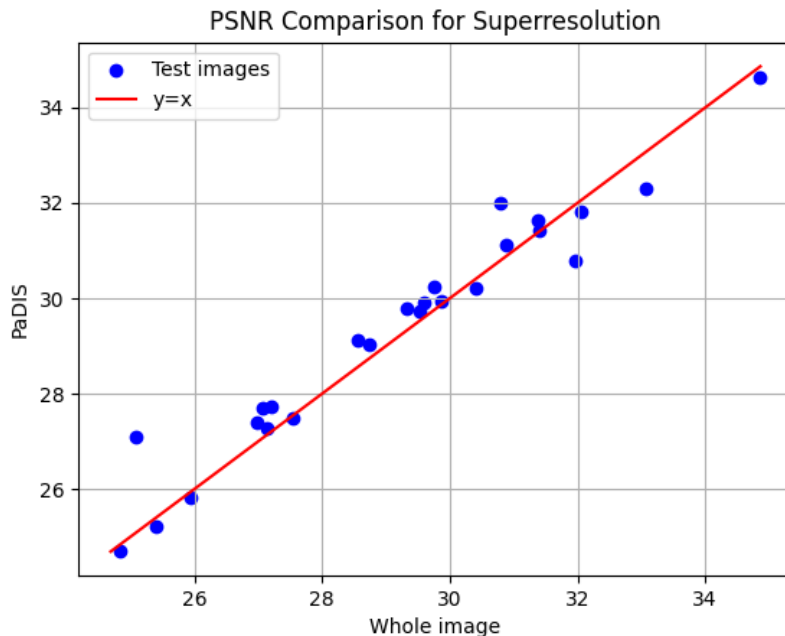


FIG A.9 – Comparison between PSNR of superresolution between whole-image model and proposed method for each image in the test dataset.

A.2 Appendix for "DiffusionBlend: Learning 3D Image Prior through Position-aware Diffusion Score Blending for 3D Computed Tomography Reconstruction"

A.2.1 Score matching derivations for DiffusionBlend++

We show how the score matching method described in [166] can be simplified in the case of assumptions such as the ones described in Table A.11.

Product of distributions. Suppose first that the distribution of interest can be expressed as $p(\mathbf{x}) = q(\mathbf{x})^a r(\mathbf{x})^b / Z$ for density functions q and r , constant positive scalars a and b , and a scaling factor Z . Following [166], to learn the score function, we can minimize the loss function

$$\mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\mathbf{y} \sim \mathcal{N}(\mathbf{x}, \sigma_t^2 I)} \left\| \mathbf{s}_\theta(\mathbf{y}, \sigma_t) - \frac{\mathbf{y} - \mathbf{x}}{\sigma_t^2} \right\|_2^2, \quad (\text{A.9})$$

where s_θ represents a neural network. Denoting the score functions of p, q , and r by s, s_p , and s_q , we have $s(\mathbf{x}) = a s_q(\mathbf{x}) + b s_r(\mathbf{x})$. Hence, if we instead use neural networks to learn s_p and s_q , we could minimize the loss function

$$L_1 = \mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\mathbf{y} \sim \mathcal{N}(\mathbf{x}, \sigma_t^2 I)} \left\| a \mathbf{s}_{q, \theta}(\mathbf{y}, \sigma_t) + b \mathbf{s}_{r, \theta}(\mathbf{y}, \sigma_t) - \frac{\mathbf{y} - \mathbf{x}}{\sigma_t^2} \right\|_2^2. \quad (\text{A.10})$$

However, this loss function is computationally expensive to work with, as backpropagation through both networks is necessary. Thus, it would be ideal to derive a simpler form of this loss function.

Toward these ends, for simplicity we define $X = a \mathbf{s}_{q, \theta}(\mathbf{y}, \sigma_t) - \frac{a}{a+b} \cdot \frac{\mathbf{y} - \mathbf{x}}{\sigma_t^2}$ and $Y = b \mathbf{s}_{r, \theta}(\mathbf{y}, \sigma_t) - \frac{b}{a+b} \cdot \frac{\mathbf{y} - \mathbf{x}}{\sigma_t^2}$, where all images have been vectorized. Now

$$\|X - Y\|_2^2 = \|X\|_2^2 + \|Y\|_2^2 - 2\langle X, Y \rangle \geq 0. \quad (\text{A.11})$$

Thus, rearranging the inequality and adding $\|X\|_2^2 + \|Y\|_2^2$ to both sides yields $\|X + Y\|_2^2 \leq 2\|X\|_2^2 + 2\|Y\|_2^2$.

Returning to the original loss function, we have

$$L_1 = \mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\mathbf{y} \sim \mathcal{N}(\mathbf{x}, \sigma_t^2 I)} \|X + Y\|_2^2. \quad (\text{A.12})$$

By applying the inequality proven above, we get

$$L_1 \leq 2 \mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\mathbf{y} \sim \mathcal{N}(\mathbf{x}, \sigma_t^2 I)} \left\| a \cdot \mathbf{s}_{q, \theta}(\mathbf{y}, \sigma_t) - \frac{a}{a+b} \cdot \frac{\mathbf{y} - \mathbf{x}}{\sigma_t^2} \right\|_2^2 \quad (\text{A.13})$$

$$+ 2 \mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\mathbf{y} \sim \mathcal{N}(\mathbf{x}, \sigma_t^2 I)} \left\| b \cdot \mathbf{s}_{r, \theta}(\mathbf{y}, \sigma_t) - \frac{b}{a+b} \cdot \frac{\mathbf{y} - \mathbf{x}}{\sigma_t^2} \right\|_2^2. \quad (\text{A.14})$$

For the special case of $a = b = \frac{1}{2}$, this inequality is rewritten as

$$L_1 \leq \mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\mathbf{y} \sim \mathcal{N}(\mathbf{x}, \sigma_t^2 I)} \left\| \mathbf{s}_{q, \theta}(\mathbf{y}, \sigma_t) - \frac{\mathbf{y} - \mathbf{x}}{\sigma_t^2} \right\|_2^2 \quad (\text{A.15})$$

$$+ \mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\mathbf{y} \sim \mathcal{N}(\mathbf{x}, \sigma_t^2 I)} \left\| \mathbf{s}_{r, \theta}(\mathbf{y}, \sigma_t) - \frac{\mathbf{y} - \mathbf{x}}{\sigma_t^2} \right\|_2^2. \quad (\text{A.16})$$

Note that each of the two individual terms in the sum precisely represents the score matching equation for learning the score functions s_p and s_q . Hence, to train the networks $s_{q, \theta}$ and $s_{r, \theta}$ by minimizing L_1 , we may instead minimize the upper bound of L_1 by separately training these two networks.

In practice, we may opt to use the same network for $\mathbf{s}_{q,\theta}$ and $\mathbf{s}_{r,\theta}$ but with an additional input specifying which distribution between q and r to use. In this case, at each training iteration, we randomly choose from one of the two distributions and perform backpropagation using this distribution. More precisely, we redefine our network $\mathbf{s}_\theta(\mathbf{x}, \sigma_t, v)$ with v being either 0 or 1. When $v = 0$, $\mathbf{s}_\theta(\mathbf{x}, \sigma_t, v) = \mathbf{s}_{q,\theta}(\mathbf{x}, \sigma_t)$ and when $v = 1$, $\mathbf{s}_\theta(\mathbf{x}, \sigma_t, v) = \mathbf{s}_{r,\theta}(\mathbf{x}, \sigma_t)$. Thus the loss bound becomes

$$L_1 \leq \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\mathbf{y} \sim \mathcal{N}(\mathbf{x}, \sigma_t^2 I)} \mathbb{E}_{v \in \{0,1\}} \left\| \mathbf{s}_\theta(\mathbf{y}, \sigma_t, v) - \frac{\mathbf{y} - \mathbf{x}}{\sigma_t^2} \right\|_2^2. \quad (\text{A.17})$$

Finally, this derivation easily extends to the more general case where the distribution of interest can be expressed as

$$p(\mathbf{x}) = \prod_{i=1}^k p_i(\mathbf{x})^{1/k} / Z. \quad (\text{A.18})$$

In this case, the similarly defined score matching loss function L_1 can be upper bounded by an expression similar to (A.17), but with v being randomly selected from k possible values.

In summary, we have shown that for the case of a decomposable distribution $p(\mathbf{x})$, the score function of $p(\mathbf{x})$ can be learned simply through the score function of the individual components $p_i(\mathbf{x})$. In the special case when each of the components have equal weight, it suffices to randomly choose one of the components and backpropagate through the score matching loss function according to that component.

Separable distributions. Next, we show how the score matching method is simplified for distributions of the form $p(\mathbf{x}) = \prod_{i=1}^r p(\mathbf{x}[:, :, \mathcal{S}_i]) / Z$, where the same notation as Table A.11 is used and $\mathcal{S} = \mathcal{S}_1 \cup \dots \cup \mathcal{S}_r$ denotes an arbitrary partition of $\{1, 2, \dots, H\}$. The score function of $p(\mathbf{x})$ can be written as

$$\mathbf{s}(\mathbf{x}) = \sum_{i=1}^H \nabla \log p(\mathbf{x}[:, :, \mathcal{S}_i]) = \sum_{i=1}^H \mathbf{s}_i(\mathbf{x}[:, :, \mathcal{S}_i]), \quad (\text{A.19})$$

where \mathbf{s}_i represents the score function of the slices of \mathbf{x} corresponding to \mathcal{S}_i . Then (A.9) becomes

$$L = \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\mathbf{y} \sim \mathcal{N}(\mathbf{x}, \sigma_t^2 I)} \left\| \sum_{i=1}^H \mathbf{s}_{\theta,i}(\mathbf{x}[:, :, \mathcal{S}_i]) - \frac{\mathbf{y} - \mathbf{x}}{\sigma_t^2} \right\|_2^2. \quad (\text{A.20})$$

Since each of the \mathcal{S}_i 's are disjoint, this can be broken up and rewritten as

$$L = \sum_{i=1}^H \mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\mathbf{y} \sim \mathcal{N}(\mathbf{x}, \sigma_t^2 I)} \left\| \mathbf{s}_{\theta, i}(\mathbf{x}[:, :, \mathcal{S}_i]) - \frac{\mathbf{y}[:, :, \mathcal{S}_i] - \mathbf{x}[:, :, \mathcal{S}_i]}{\sigma_t^2} \right\|_2^2. \quad (\text{A.21})$$

Thus, after replacing the outer sum with an expectation over i , this is equivalent to randomly choosing one of the partitions \mathcal{S}_i and performing denoising score matching on only $\mathbf{x}[:, :, \mathcal{S}_i]$.

A very similar derivation holds for the general case where the 3D volume \mathbf{x} can be partitioned into an arbitrary number of smaller volumes of any shape $\mathbf{x} = \mathbf{x}_1 \cup \mathbf{x}_2 \cup \dots \cup \mathbf{x}_H$ and $p(\mathbf{x}) = \prod_{i=1}^H p(\mathbf{x}_i) / Z$. For this case, training consists of randomly selecting one of the partitions at each iteration and performing score matching on that partition. For example, when $\mathbf{x}_i = \mathbf{x}[:, :, i]$, it is common to select 2D slices from the training volumes and learn a two dimensional diffusion model on those slices [42, 109].

Applying to DiffusionBlend++. When $p(\mathbf{x})$ follows the distribution in DiffusionBlend++ we can combine the results of the previous two sections to show how to perform score matching. In the first part of this section, we showed how to perform score matching for a distribution expressed as a product of “simpler” distributions by performing score matching on the individual distributions. DiffusionBlend++ follows this assumption where

$$p_i(\mathbf{x}) = \left(\prod_{j=1}^r p(\mathbf{x}[:, :, \mathcal{S}_j]) \right) / Z_i. \quad (\text{A.22})$$

Here, i represents an index that can iterate through the ways of partitioning $\mathcal{S} = \mathcal{S}_1 \cup \dots \cup \mathcal{S}_r$. The input v to the network specifying which of the simpler distributions is used is embedded as the relative position encoding for each of the partitions as described in Section 3. Finally, to learn the score function of $p_i(\mathbf{x})$, we can use the loss function in (A.21).

A.2.2 Additional Algorithms

The reconstruction algorithm for DiffusionBlend is provided below.

The training algorithms for DiffusionBlend and DiffusionBlend++ are provided below.

Algorithm 10 DiffusionBlend

Require: \mathbf{A} , M , $\zeta_i > 0$, j , \mathbf{y}

Initialize $\mathbf{x}_T \sim \mathcal{N}(0, \sigma_T^2 \mathbf{I})$

for $t = T : 1$ **do**

 For each i compute $\epsilon_\theta(\mathbf{x}_t[:, :, i] | \mathbf{x}_t[:, :, i-j : i-1], \mathbf{x}_t[:, :, i+1 : i+j])$

 Compute $\mathbf{s} = \nabla \log p(\mathbf{x}_t)$ using (5.3)

 Compute $\hat{\mathbf{x}}_t = \mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t]$ using Tweedie's formula

 Set $\hat{\mathbf{x}}'_t = \text{CG}(\mathbf{A}^* \mathbf{A}, \mathbf{A}^* \mathbf{y}, \hat{\mathbf{x}}_t)$

 Sample \mathbf{x}_{t-1} using $\hat{\mathbf{x}}'_t$ and \mathbf{s} via DDIM sampling

end for

Return \mathbf{x} .

Algorithm 11 DiffusionBlend training

repeat

 Select $\mathbf{x} \sim p(\mathbf{x})$

 Select $t \sim \text{Uniform}[1, T]$

 Set $\mathbf{y} \sim \mathcal{N}(\mathbf{x}, \sigma_t^2 I)$

 Select $i \sim \text{Uniform}[1, H]$

 Take gradient descent step on $\nabla_\theta \| (D_\theta(\mathbf{y}[:, :, i-j : i+j], \sigma_t) - \mathbf{x}[:, :, i]) / \sigma_t^2 \|_2^2$

until converged

Return D_θ

Algorithm 12 DiffusionBlend++ training

repeat

 Select $\mathbf{x} \sim p(\mathbf{x})$

 Select $t \sim \text{Uniform}[1, T]$

 Set $\mathbf{y} \sim \mathcal{N}(\mathbf{x}, \sigma_t^2 I)$

 Select a partition $\mathcal{S} = \mathcal{S}_1 \cup \dots \cup \mathcal{S}_r$

 Select $i \sim \text{Uniform}[1, r]$

 Take gradient descent step on $\nabla_\theta \| (D_\theta(\mathbf{y}[:, :, \mathcal{S}_i], \sigma_t) - \mathbf{x}[:, :, \mathcal{S}_i]) / \sigma_t^2 \|_2^2$

until converged

Return D_θ

A.2.3 Ablation studies

We run the following ablation studies to examine each of the individual components of our DiffusionBlend++ method. Firstly, we examine the performance gain of adding adjacency slice blending (DiffusionBlend+) and adding cross-slice blending. Next, we examine the effect of including the positional encoding as an input to the network. Then we look at the quantitative metrics of the reconstructed images when applying different numbers of NFEs for the comparison methods. Finally, we examine the effect of choosing different slices for each partition.

Adjacency	Cross	PSNR \uparrow	SSIM \uparrow
		34.85	0.954
✓		36.02	0.965
✓	✓	36.48	0.968

TBL A.6 – Effectiveness of Blending Modules, Sagittal view performance on AAPM

Effectiveness of adjacency-slice blending and cross-slice blending We demonstrate that both the adjacency-slice blending and the cross-slice blending module are instrumental to a better reconstruction quality. Table A.6 demonstrates the effectiveness of adding blending modules to the reverse sampling. Given the pretrained diffusion prior over slice patches, we observe that adding the adjacency-slice blending module improves the PSNR over a fixed partition by 1.17dB, and adding an additional cross-slice blending module further improves the PSNR by 1.63dB. Fig. 5.5 demonstrates that adding the cross-slice blending module removes artifacts and recovers sharper edges.

Robust performance with low NFEs. Since DDS and DiffusionBlend++ both use the DDIM sampler for acceleration, we performed experiments with both of these methods using different NFEs. The left of Fig. A.10 shows graphs of these two methods and Table A.7 shows the quantitative results. DDS is very sensitive to the number of NFEs used and there is a sharp dropoff in PSNR if too few or too many NFEs are used. On the other hand, DiffusionBlend++ performs the best for the highest number of NFEs due to the slice blending strategy while still obtaining superior results for 50 NFEs. Also, this method is much more robust to varying NFEs, displaying only 1.4dB of variance in the shown results compared to 2.4dB of variance for DDS. For fair comparisons, we use 200 NFEs for all the main experiments.

Frequency of applying slice jumps. To demonstrate the use of jump slice partitions at reconstruction time, we performed experiments varying the frequency of applying these jump slices. For instance, for a frequency of 8, the reconstruction algorithm consisted of

TBL A.7 – Axial PSNR for 8 view SVCT recon for different NFEs

Method	50	100	200	334
DDS	30.8	32.2	33.2	32.7
DiffusionBlend++	34.5	35.0	35.7	35.9

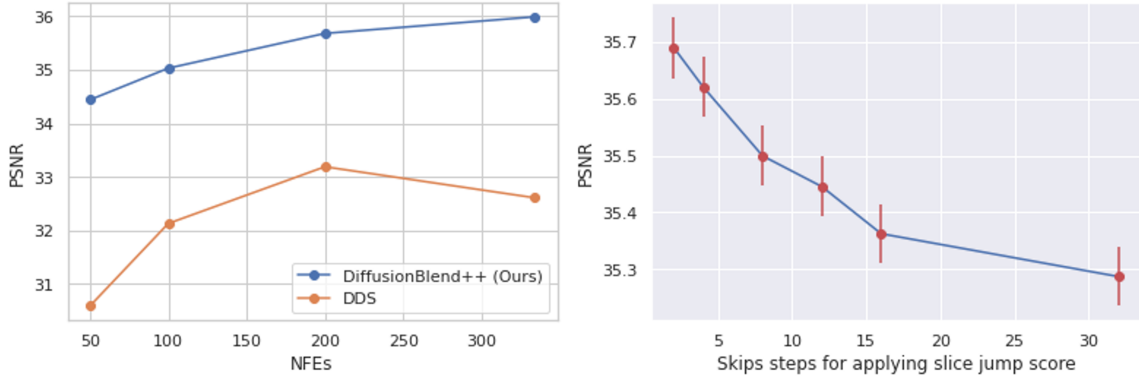


FIG A.10 – Quantitative results (axial view) of CT reconstruction with 8 views on AAPM dataset for different NFEs and slice blending methods.

updating the volume using jump slices for iteration numbers that are a multiple of 8 and updating using adjacent slices for all other iterations. The right of Fig. A.10 shows a graph of the results for different frequencies and the quantitative results are presented in Table A.8. The best results are obtained when the frequency is 2, corresponding to alternating updates with adjacent slices and jump slices, and the PSNR decreases monotonically as the frequency increases. This indicates that the jump slices capture more nonlocal information across a volume and help to improve the image quality.

TBL A.8 – Axial PSNR for 8 view SVCT recon for different slice jump frequencies

Frequency	2	4	8	12	16	32
PSNR	35.69	35.62	35.50	35.45	35.37	35.28

TBL A.9 – Wall times of various methods for 8 view 3D CT reconstruction

Method	NFEs	Wall time (min)
DiffusionMBIR	2000	1400
TPDM	2000	1200
DDS	200	48
DiffusionBlend	200	70
DiffusionBlend++	200	32

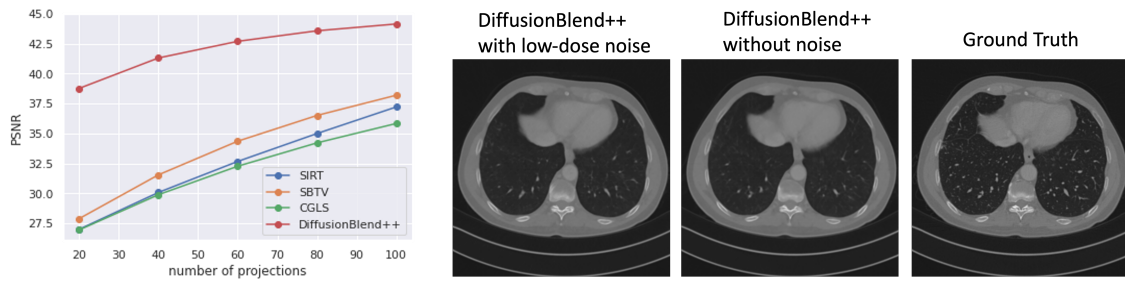


FIG A.11 – Left: Performance of DiffusionBlend++ on more angles, Right: Reconstruction of DiffusionBlend++ with low-dose noise

A.2.4 Additional Results

Classical Baselines and more Projection Angles We provide additional results with classical baselines (without deep learning) such as SIRT, SBTv, and CGLS for the LDCT dataset. Results show that our method outperforms the baselines significantly for every angle we evaluated on. DiffusionBlend++ starts to reconstruct images very close to the ground truth with 20 projections or more, but other baselines such as SIRT and CGLS

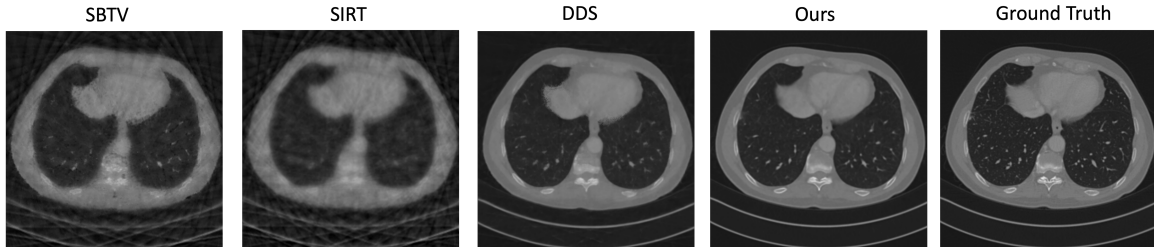


FIG A.12 – Comparison of DiffusionBlend++ with classical methods

Method	Sparse-View CT Reconstruction on AAPM						Sparse-View CT Reconstruction on LIDC					
	8 views		6 views		4 views		8 Views		6 Views		4 Views	
	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow
FBP	2.64	0.16	2.88	0.26	2.91	0.20	2.57	0.17	2.78	0.21	2.84	0.16
FBP-UNet	3.46	0.30	3.34	0.28	3.04	0.31	3.42	0.31	3.23	0.29	3.14	0.27
DiffusionMBIR	1.93	0.08	1.48	0.13	1.28	0.11	1.89	0.09	1.56	0.12	1.31	0.14
TPDM	-	-	-	-	-	-	2.32	0.14	2.02	0.16	2.52	0.19
DDS 2D	1.76	0.07	2.04	0.10	2.73	0.11	1.85	0.09	2.13	0.13	2.67	0.18
DDS 2D	1.76	0.07	2.04	0.10	2.73	0.11	1.85	0.09	2.13	0.13	2.67	0.18
DDS 2D	1.68	0.06	1.96	0.09	2.65	0.11	1.84	0.09	2.10	0.12	2.64	0.18
DiffusionBlend (Ours)	1.67	0.06	1.78	0.08	1.98	0.09	1.70	0.07	2.03	0.11	2.54	0.16
DiffusionBlend++ (Ours)	1.50	0.06	1.65	0.08	1.71	0.10	1.60	0.09	1.68	0.10	1.82	0.11

TBL A.10 – Standard Deviation of Performance on Sparse-View CT Reconstruction on Sagittal View for AAPM and LIDC datasets. Best results are in bold.

still struggle to get a satisfying reconstruction with 60 views or more. Fig. A.11 shows the reconstruction performance on the coronal plane of different methods. We observe that DiffusionBlend++ (ours) has a significant margin above baselines for every view. Note that DiffusionBlend++ still outperforms DDS2D significantly with >40 views, which demonstrates that our 3D prior is still very useful even with much more views. We also simulate low-dose noise to the reconstruction, which showing our algorithm is robust to noise by a minor decrease in reconstruction performance. Our method (DiffusionBlend++) is shown to outperforms all baselines at every angle as in Fig. A.12.

Error Bars We demonstrate the standard deviation of the results with sparse-view CT reconstruction on AAPM and LIDC dataset here to demonstrate that the result is statistically significant.

Fig. A.13 shows the visual results for SVCT reconstruction with 8 views on the LIDC dataset.

Fig. A.14 shows the visual results for SVCT reconstruction with 6 views on the LIDC dataset.

Fig. A.15 shows the visual results for SVCT reconstruction with 4 views on the LIDC dataset.

Fig. A.16 shows the visual results for LACT reconstruction on the LIDC dataset.

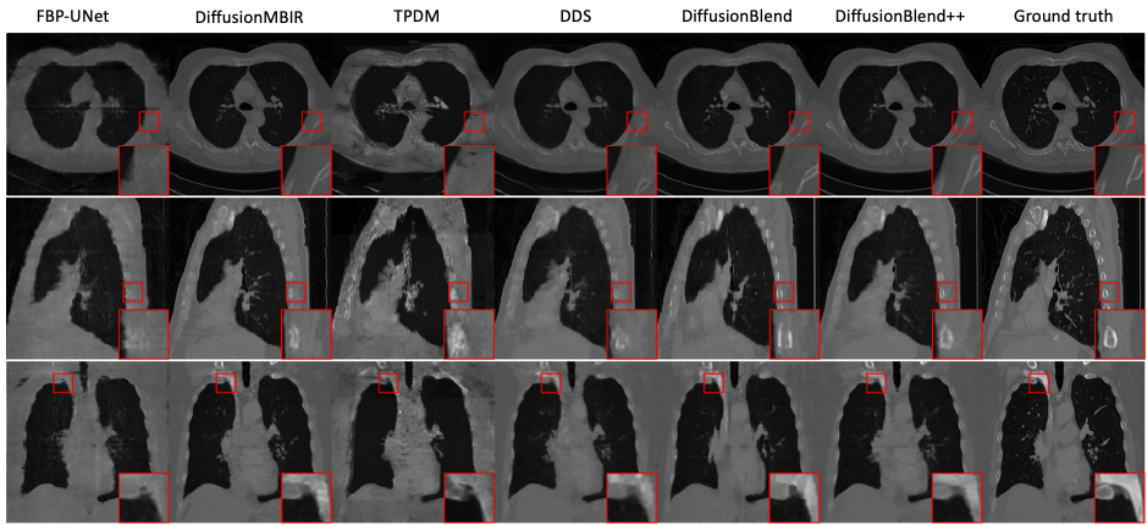


FIG A.13 – Results of 3D CT reconstruction with 8 views on LIDC dataset. Top row is axial view, middle row is sagittal view, bottom row is coronal view.

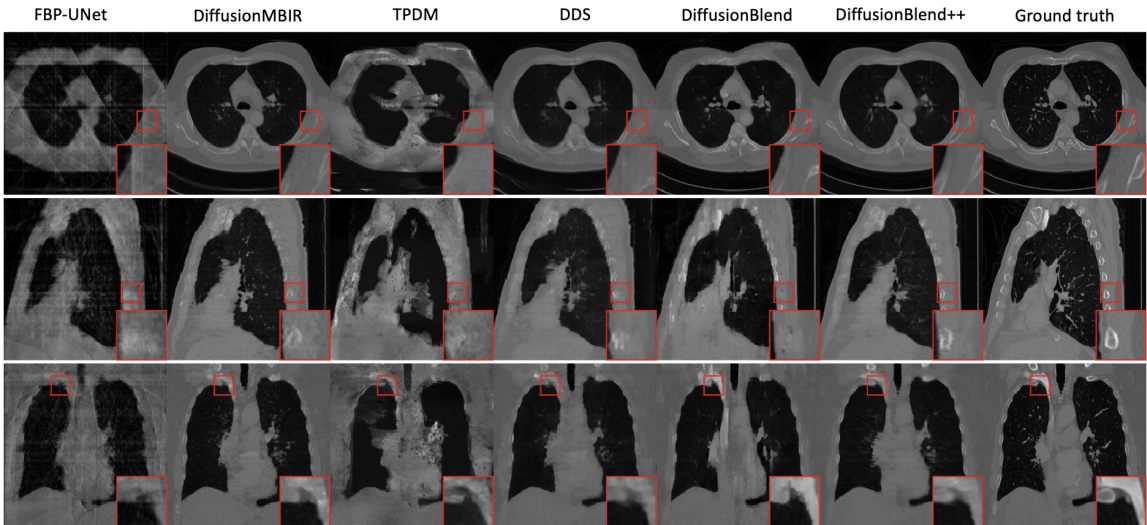


FIG A.14 – Results of 3D CT reconstruction with 6 views on LIDC dataset. Top row is axial view, middle row is sagittal view, bottom row is coronal view.

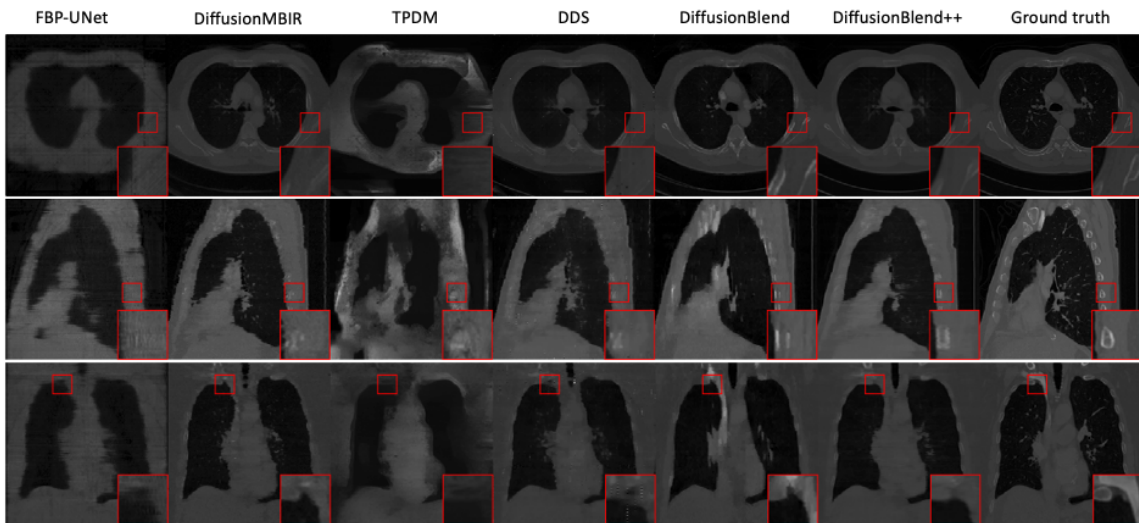


FIG A.15 – Results of 3D CT reconstruction with 4 views on LIDC dataset. Top row is axial view, middle row is sagittal view, bottom row is coronal view.

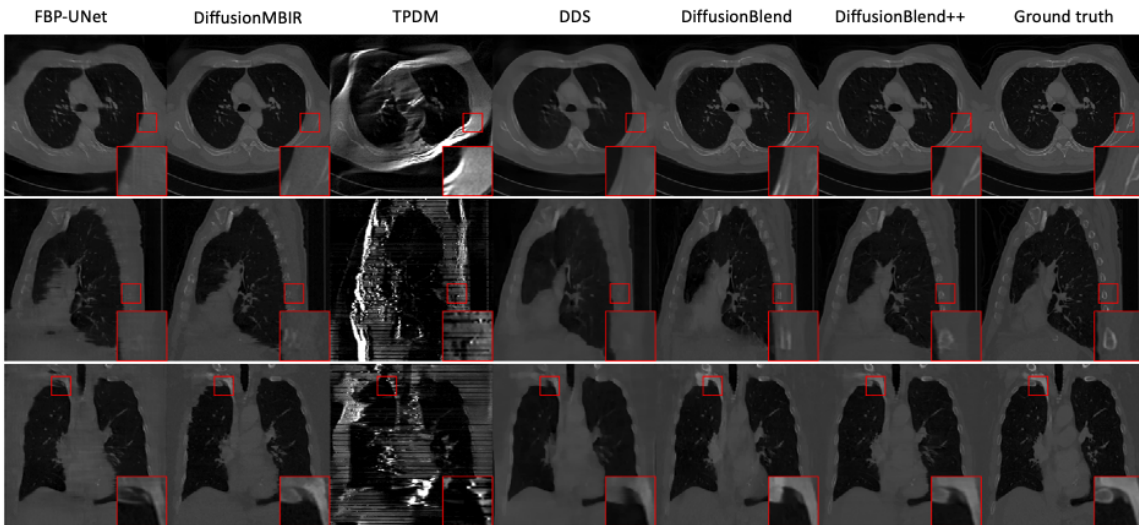


FIG A.16 – Results of limited angle 3D CT reconstruction on LIDC dataset. Top row is axial view, middle row is sagittal view, bottom row is coronal view.

TBL A.11 – 3D prior modeling methods

Method	Distribution Model
DiffusionMBIR [42]	$\prod_{i=1}^H p(\mathbf{x}[:, :, i])/Z$
TPDM [109]	$\left(\prod_{i=1}^N q_{\theta}(\mathbf{x}[:, :, i])^{\alpha}\right) \left(\prod_{j=1}^N q_{\phi}(\mathbf{x}[j, :, :])^{\beta}\right) / Z$
DiffusionBlend	$\prod_{i=1}^H p(\mathbf{x}[:, :, i] \mathbf{x}[:, :, i-j : i-1], \mathbf{x}[:, :, i+1 : i+j]) / Z$
DiffusionBlend++	$\prod_{i=1}^r p(\mathbf{x}[:, :, \mathcal{S}_i]) / Z$

A.2.5 Experiment parameters

Since axial slices belonging to the same volume that are far apart have limited correlation, DiffusionBlend++ selects only partitions of \mathcal{S} for training where slices belonging to the same partition are fairly close to one another. Then the same range of possible partition schemes are used during reconstruction time. More precisely, we take the size of each \mathcal{S}_i to be 3 and first repetition pad the volume so that the number of axial slices is a multiple of 9. Then we consider the following partitions:

- $\mathcal{S}_1 = \{1, 2, 3\}$, $\mathcal{S}_2 = \{4, 5, 6\}$, $\mathcal{S}_3 = \{7, 8, 9\}$. Furthermore, for all integers $k > 1$, $\mathcal{S}_k = \mathcal{S}_{k-3} \oplus 9 \lfloor (k-1)/3 \rfloor$, where \oplus represents adding the same number to each element of the set. For example, $\mathcal{S}_4 = \{10, 11, 12\}$, $\mathcal{S}_5 = \{13, 14, 15\}$, $\mathcal{S}_6 = \{16, 17, 18\}$.
- $\mathcal{S}_1 = \{1, 4, 7\}$, $\mathcal{S}_2 = \{2, 5, 8\}$, $\mathcal{S}_3 = \{3, 6, 9\}$. Furthermore, for all integers $k > 1$, $\mathcal{S}_k = \mathcal{S}_{k-3} \oplus 9 \lfloor (k-1)/3 \rfloor$.

A.2.6 Comparison experiment details

FBP-UNet. We used the same neural network architecture as the original paper [89]. Individual networks were trained for each of the 8 view, 6 view, 4 view, and LACT experiments for each of the datasets. Each of the networks were trained from scratch with a batch size of 32 for 150 epochs.

DiffusionMBIR. We separately trained networks for the AAPM and LIDC datasets by fine-tuning the original checkpoint provided in [42] for 100 and 10 epochs, respectively. The batch size was set to 4. For reconstruction, we used the same set of hyperparameters

for all of the experiments: $\lambda = 0.04$, $\rho = 10$, and $r = 0.16$ for the sampling algorithm. 2000 NFEs were used for the diffusion process.

TPDM. We fine-tuned the axial and sagittal checkpoints provided in [109] on the LIDC dataset for 10 epochs. For reconstruction, we used 2000 NFEs and alternated between updating the volume using the axial checkpoint and sagittal checkpoint, with each checkpoint being used equally frequently. The DPS step size parameter was set to $\zeta = 0.5$.

DDS. We separately trained networks for the AAPM and LIDC datasets by fine-tuning the original checkpoint provided in [41] for 100 and 10 epochs, respectively. We used 100 NFEs at reconstruction as this was observed to give the best performance. The reconstruction parameters were set to $\eta = 0.85$, $\lambda = 0.4$, and $\rho = 10$. Five iterations of conjugate gradient descent were run per diffusion step. For DDS 2D, the parameters were left unchanged with the exception of using $\rho = 0$ to avoid enforcing the TV regularizer between slices.

SBTV. We implement this algorithm with variables splitting of 3D anisotropic TV regularization (Dz, Dx, and Dy). We first check number of iterations, note that the performance converges with around 30 iterations. We did a grid search of hyperparameters on 9 validation images (not in the test set) for every projection angles.

SIRT. This algorithm iteratively updates the reconstruction based on the residual between projection of the reconstruction and the GT. It only has the number of iterations as its hyper-parameter. We note that during inference, PSNR increases with more iterations, but saturates later. So we set the total number iterations to be 1000, with an early stopping threshold of $1e-6$ between two consecutive iterations.

CGLS. This algorithm uses conjugate gradient for solving least square problems. In our case, we use $CG(A^T A + \rho x^T x, A^T y)$, ρ is set to be $1e-4$ based on grid search for numerical stability. We tune the number of iterations on validation set, and find that performance saturates at around 25 iterations.

A.3 Appendix for Local Patches Meet Global Context: Scalable 3D Diffusion Priors for Computed Tomography Reconstruction

A.3.1 Training Algorithm

Algorithm 13 summarizes the training algorithm for 3D patch diffusion.

Algorithm 13 Training 3D patch diffusion model

Require: image size $N_x \times N_y \times N_z$, patch size P

1: **repeat**

2: $t \sim \text{Uniform}(\{1, \dots, T\})$

3: $c \sim \text{Uniform}(\{0, \dots, (N_x + P)(N_y + P)(N_z + P)\})$

4: $\epsilon \sim \mathcal{N}(0, \mathbf{I})$

5: $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x} + \sqrt{1 - \bar{\alpha}_t} \epsilon$

6: $\mathbf{u} = \mathbf{G}_c \mathbf{x}_t$

▷ Patchify

7: $\mathbf{v} = \mathbf{D} \mathbf{x}_t$

▷ Downsample

8: Apply gradient descent on $\nabla_{\theta} \|D_{\theta}(\mathbf{u}, \mathbf{v}) - \mathbf{G}_c \epsilon\|^2$

9: **until** converged

A.3.2 Experiment Parameters

For the network used in the proposed method, we extended the 2D UNet architecture used in [165] to 3D by replacing all 2D operators with their 3D counterparts. The model was optimized by Adam, and we used an exponential moving average with a decay rate of 0.999 to stabilize the training. The diffusion process used 1000 timesteps with a linearly increasing noise variance schedule from 0.0001 to 0.02. All of our proposed models were trained with the same architecture. Table A.12 summarizes the architecture and hyperparameter settings.

A.3.3 Limitations

Limitation on reducing sampling steps In Eq. 6.2, we model our 3D prior as the product of all non-overlapping patches for a given initial point, and then take the product over all possible initial points. However, during each sampling step, we approximate the 3D prior as the product of all non-overlapping patches from one random initial point [82]

TBL A.12 – Model hyperparameters.

Base channel width	64
Channel multipliers	[1, 2, 4, 4]
# of input channel	5
# of output channel	1
Attention resolution	[8, 8]
# of residual block	2
Learning rate	2×10^{-5}

to reduce the computation cost. Consequently, the model cannot be sampled with smaller step sizes, as it relies on a random initial point at each step to reduce the boundary artifacts. Figure A.17 shows that with smaller sampling steps, the boundary artifacts of patches still remains.

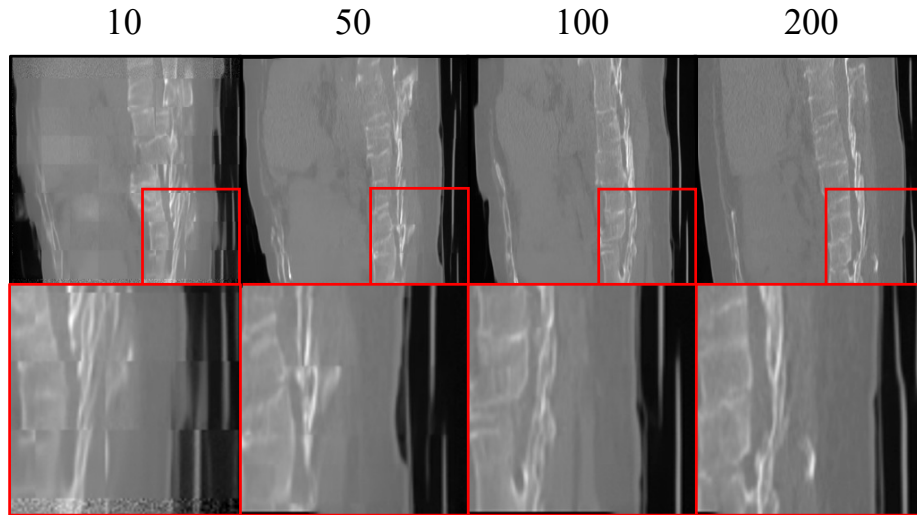


FIG A.17 – Sagittal slices of generated 3D volumes for different numbers of different sampling steps. Generated image quality degrades if the number of steps is reduced too much.

Limitation on boundary region of volume Our 3D prior model zero pads each 3D volume. The zero padding can create unwanted artifacts in the top few and bottom few slices of generated volumes. Figure A.18 shows the first and last four axial slices of a representative generated volume. Future work could consider other padding approaches.

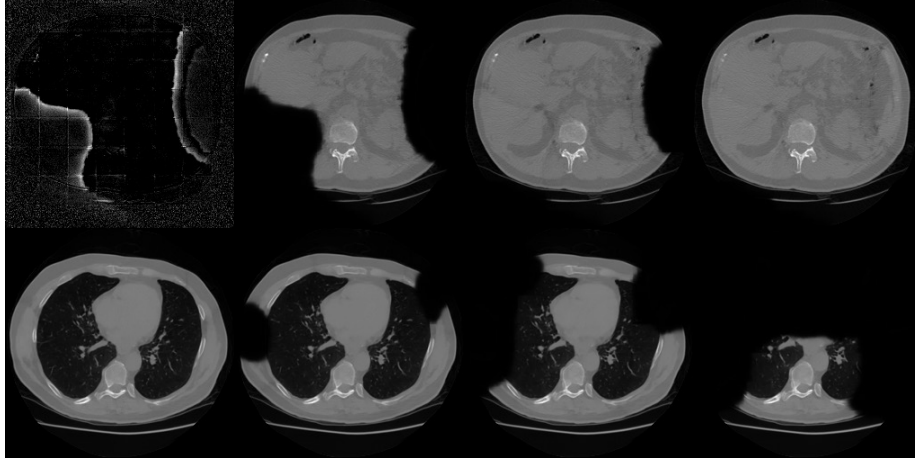


FIG A.18 – Top row shows the first four axial slices, and the bottom row shows the last four axial slices of the generated volume

Low generation quality on large volume Although our proposed method achieves SOTA performance on reconstructing $512 \times 512 \times 256$ CT volume, it fails to generate fine-grained details as shown in Figures A.25, A.26, A.27. We speculate that because the number of model parameters remains the same when training priors on $256 \times 256 \times 256$ and $512 \times 512 \times 256$ sized volumes, it is difficult to achieve comparable generative performance for the larger volume size.

A.3.4 Trade-offs

In this section, we discuss the trade-offs between the different patch sizes. All experiments were conducted on the $256 \times 256 \times 256$ LIDC-IDRI prior for the 8-view SVCT task. Here, the patch size plays a critical role in training efficiency: when the patch size doubles, then the batch size decreases by a factor of eight. However, if we reduce the patch sizes too much, the performance drops drastically, as shown in Figure A.13. Therefore, we choose a patch size that is a factor of eight of the full image resolution as our baseline, balancing training efficiency and reconstruction quality.

TBL A.13 – Trade-offs on different patch sizes

Patch size (factor)	PSNR \uparrow	training time	sampling time
$64 \times 64 \times 64$ (4)	31.74	8 days	30 min
$32 \times 32 \times 32$ (8)	33.06	8 days	20 min
$16 \times 16 \times 16$ (16)	27.87	3 days	10 min

A.3.5 Additional Figures

Figure A.19 shows visual results of 8-view CT reconstruction on the AAPM dataset with our proposed method and some comparison methods. Compared to the result with DiffusionBlend (fine-tuned from a checkpoint pretrained on natural images), our method is able to preserve more details, such as bone artifacts and blood vessels of a lung.

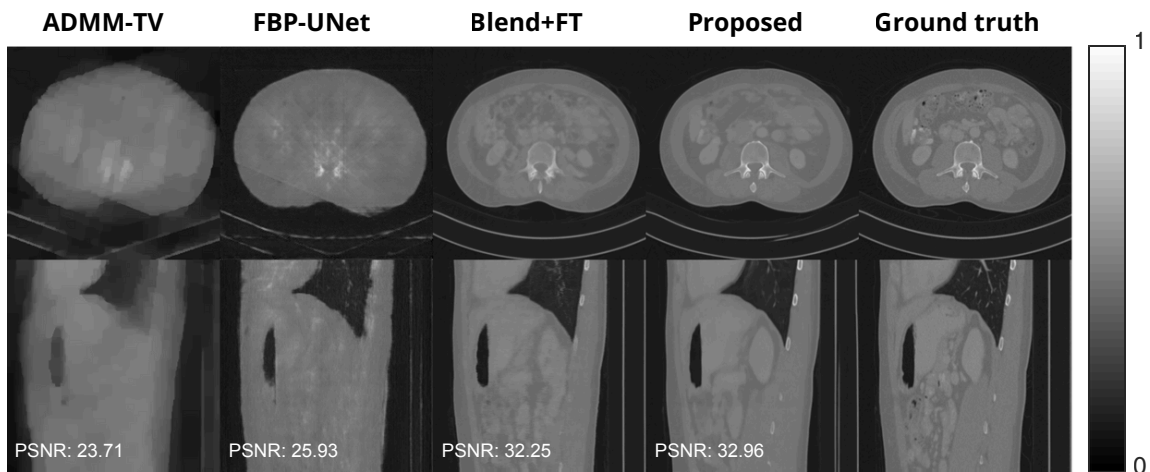


FIG A.19 – Results of our proposed method and comparison methods for 8-view 3D CT reconstruction on AAPM dataset. The top row shows an axial slice and the bottom row shows a sagittal slice from each reconstructed volume.

Figure A.20 shows visual results of 60-view CT reconstruction on the $256 \times 256 \times 256$ LIDC dataset with our proposed method and some comparison methods. We see that even in this fairly sparse setting, the proposed method is able to obtain a very high quality reconstruction that is absent of any hallucinations or artifacts.

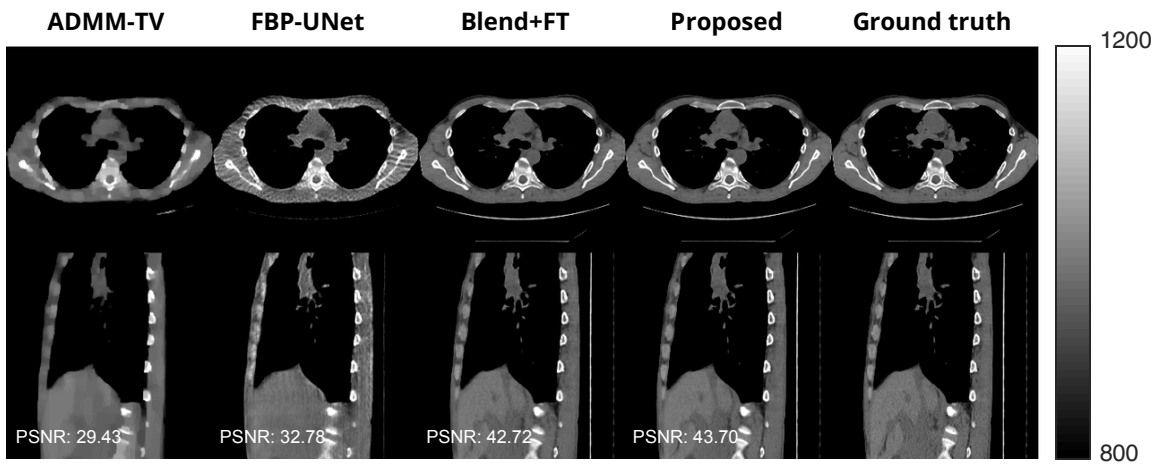


FIG A.20 – Results of our proposed method and comparison methods for 60 view CT recon on $256 \times 256 \times 256$ LIDC dataset. Plots are shown in Hounsfield units. The top row shows the axial slice and the bottom row shows the sagittal slice from the reconstructed volume.

Figure A.21, A.22 shows visual results of 20-view and 60-view CT reconstruction on the $512 \times 512 \times 256$ LIDC dataset with our proposed method and some comparison methods. To the best of our knowledge, we are the first work to leverage unconditional diffusion priors to perform CT reconstruction on CT volumes of size 512.

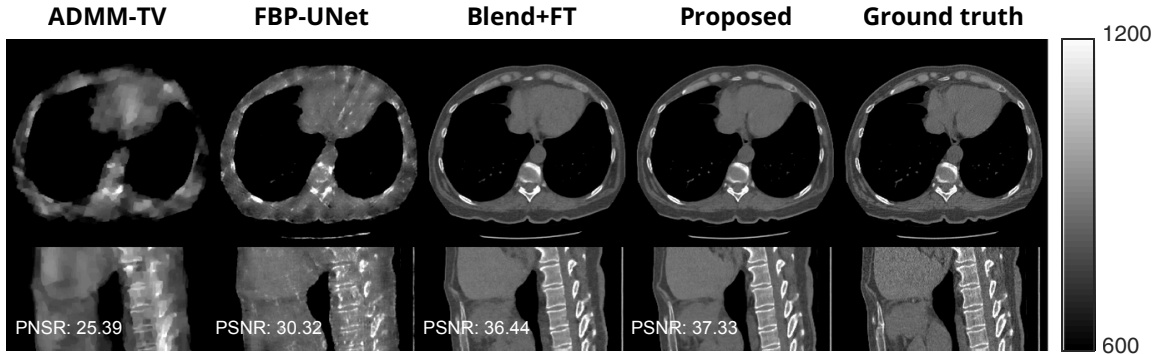


FIG A.21 – Results of our proposed method and comparison methods for 20 view CT recon on $512 \times 512 \times 256$ LIDC dataset. Plots are shown in Hounsfield units. The top row shows the axial slice and the bottom row shows the sagittal slice from the reconstructed volume.

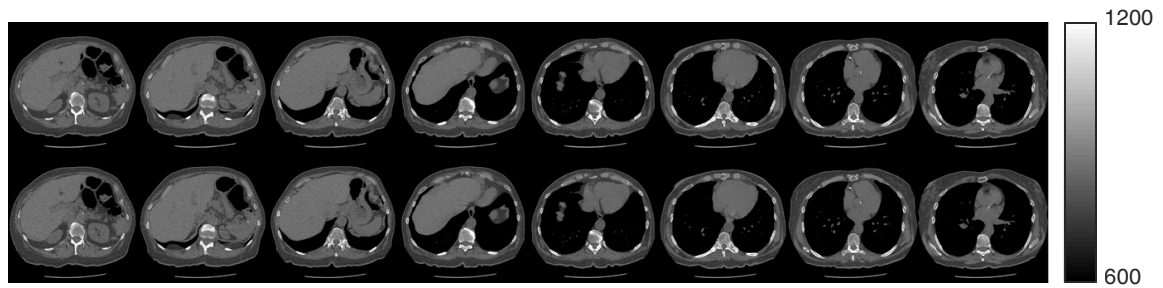


FIG A.22 – Results of our proposed method for 60 view CT recon on $512 \times 512 \times 256$ LIDC dataset. Plots are shown in Hounsfield units. The top row shows our proposed method and the bottom row shows the ground truth. All slices are axial slices from the same volume.

Figure A.23, A.24 shows the selected slice of unconditionally sampled volume from 256×256 LIDC-IDRI prior and $256 \times 256 \times 256$ AAPM prior. We calculated its nearest neighbor from the training dataset and showed the corresponding slices below. This further solidifies that our model could generate volumes instead of memorizing from the training dataset.

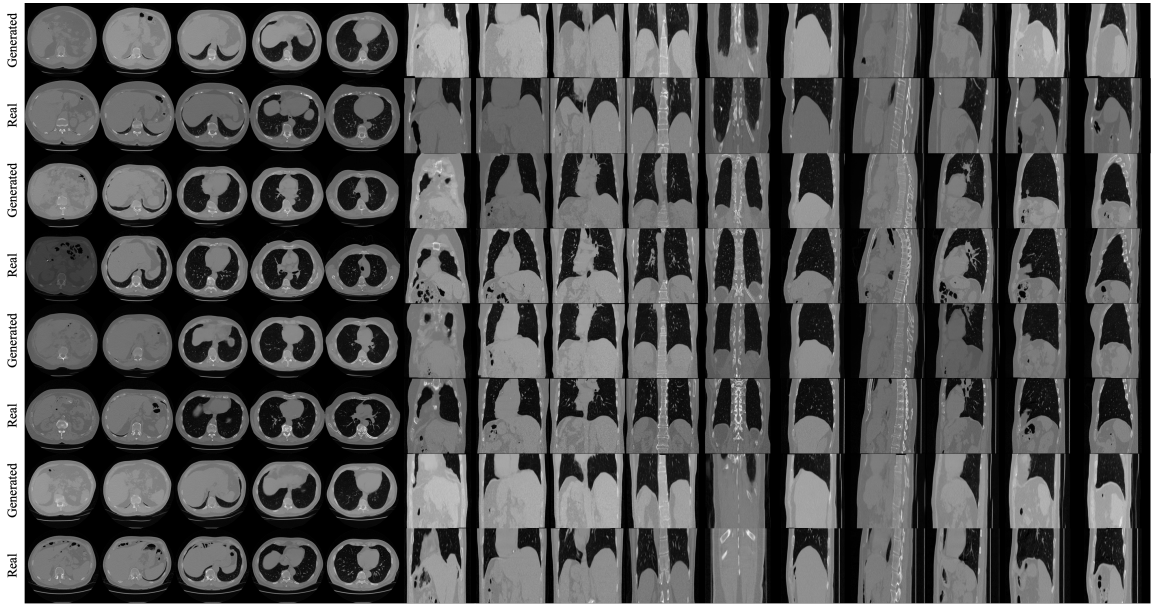


FIG A.23 – Unconditional sampling result on $256 \times 256 \times 256$ LIDC-IDRI prior and its nearest neighbor from the training dataset sliced by axial, coronal, and sagittal view.

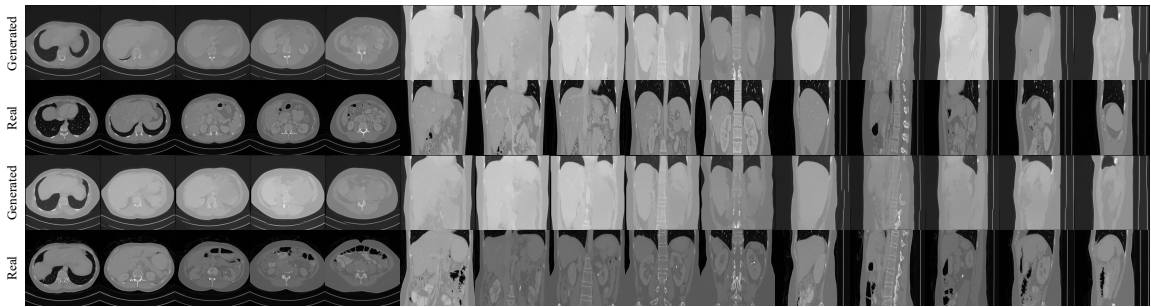


FIG A.24 – Unconditional sampling result on $256 \times 256 \times 256$ AAPM prior and its nearest neighbor from the training dataset sliced by axial, coronal, and sagittal view.

Figures A.25, A.26, A.27 show unconditional sampling results of $512 \times 512 \times 256$, $256 \times 256 \times 256$ LIDC-IDRI prior and $256 \times 256 \times 256$ AAPM prior. Here, our proposed method is capable of generating realistic volumes, even from the $256 \times 256 \times 256$ AAPM dataset that consists of only 9 training volumes.

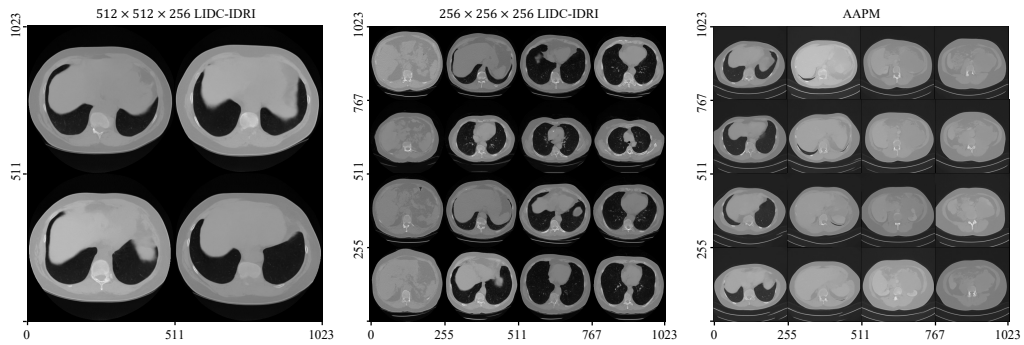


FIG A.25 – Unconditional sampling examples for three different priors (axial slices of 3D volumes).

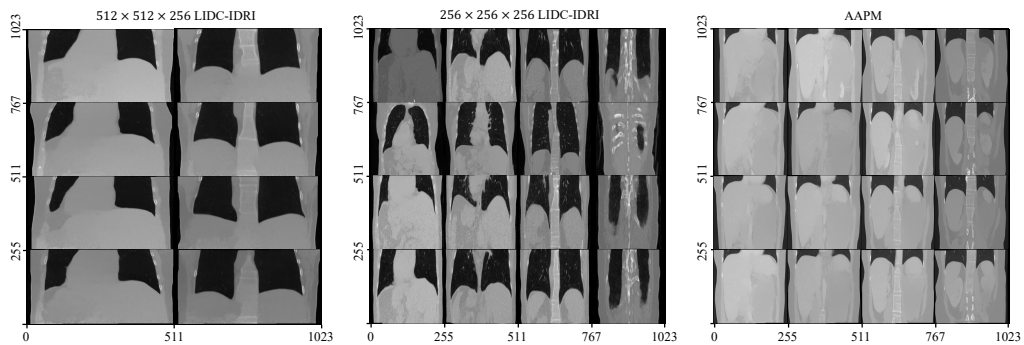


FIG A.26 – Unconditional sampling examples for three different priors (coronal slices of 3D volumes).

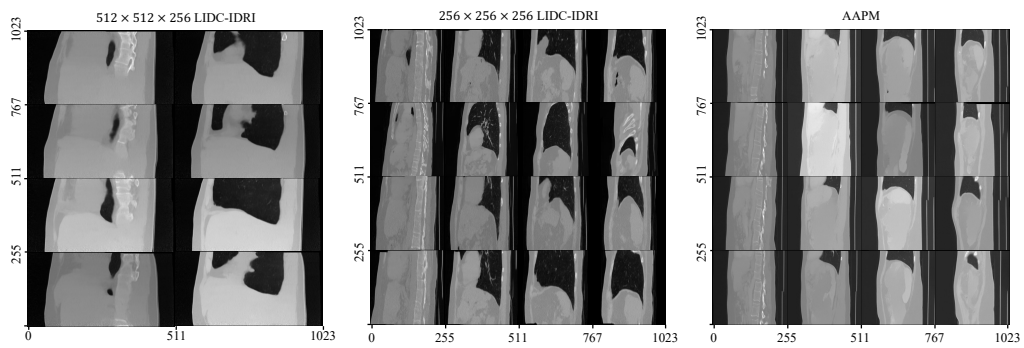


FIG A.27 – Unconditional sampling examples for three different priors (sagittal slices of 3D volumes).

Figure A.28, A.29, A.30 shows the continuous axial, coronal and sagittal slices of a volume that is unconditionally sampled from $256 \times 256 \times 256$ LIDC-IDRI prior.

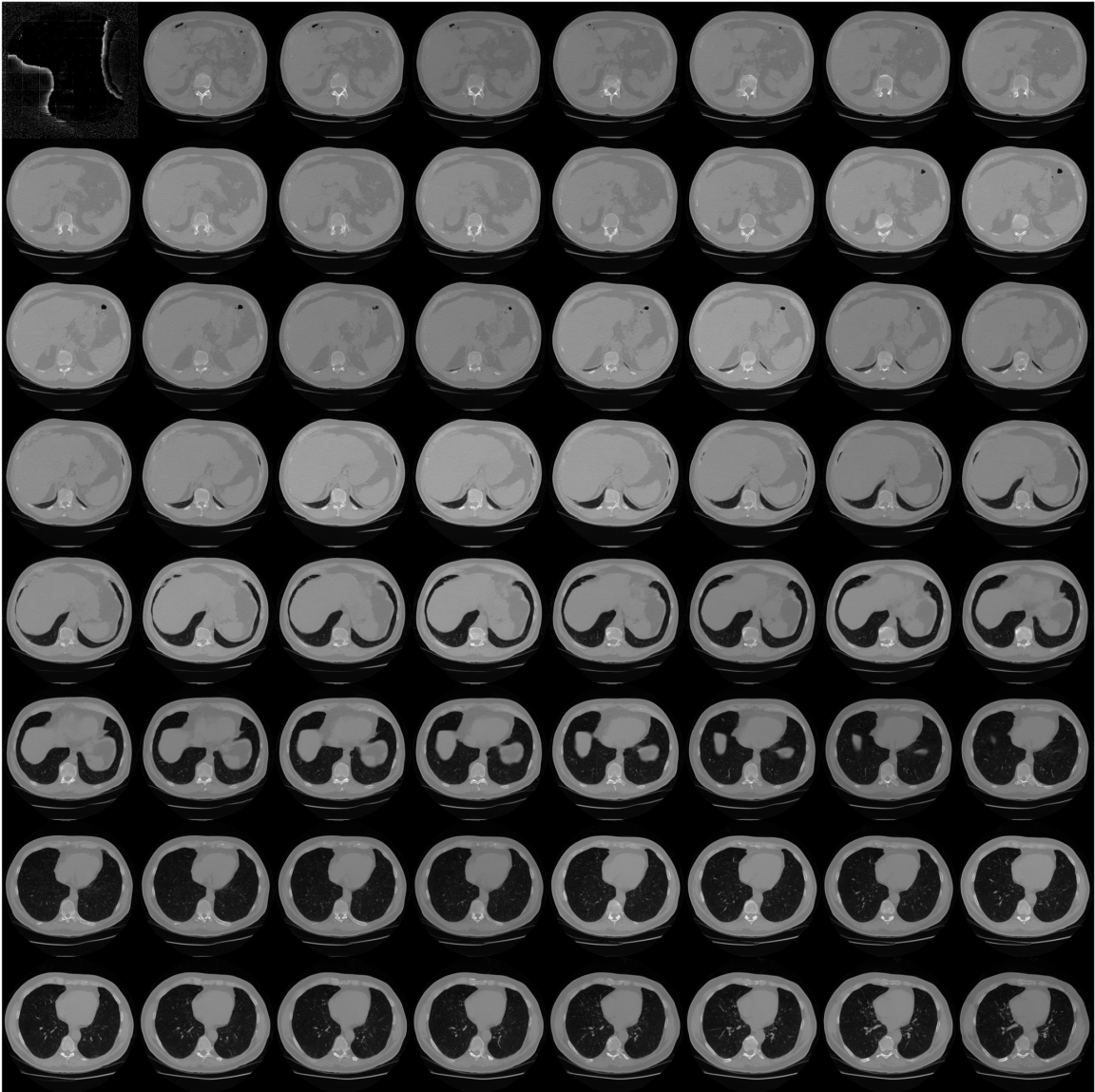


FIG A.28 – Continuous axial slice of a volume sampled from $256 \times 256 \times 256$ LIDC-IDRI prior; every fourth slice is shown

A.4 Appendix for "Test-Time Adaptation Improves Inverse Problem Solving with Patch-Based Diffusion Models "

A.4.1 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 when using the self-supervised loss, [10] 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 [10] and [47], we only used the LoRA module for attention and convolution layers. We also allowed the biases of the network to be changed.

Tables A.14 and A.15 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 resulted in worse reconstructions than simply fine-tuning the entire network. The visual results are especially apparent in Figure A.32: the reconstructed image became oversmoothed when using LoRA and artifacts became 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 of the mismatched network is not sufficient to represent the new distribution and thus the self-supervised loss function results in smoothed images.

Effect of network size. When applying the self-supervised loss, another potential method to avoid overfitting is to use a smaller network. We trained networks with differing numbers of base channels (but no other modifications) on the ellipse phantom dataset and then used Algorithm 5 to perform 60-view CT reconstruction with test-time adaptation. Table A.16 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 A.33

TBL A.14 – 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↑	SSIM ↑	PSNR↑	SSIM ↑
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

TBL A.15 – Performance of deblurring using self-supervised network refining with LoRA module. Best results are in bold.

Rank Parameters (%)		Patches		Whole image	
		PSNR↑	SSIM ↑	PSNR↑	SSIM ↑
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

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.

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 the self supervised loss). Table A.17 shows the results of these experiments, where we also included the results of using the in-distribution network trained on the entire 5000

TBL A.16 – Performance of 60 view CT recon using test-time adaptation 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

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 [82], which showed that patch-based models outperform whole image models in cases of limited data, and the results here. Since the networks in [82] were trained from scratch, more data was required: the smallest datasets used in [82] 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 7.6 shows that we fine-tuned a patch-based model in only about 2 hours, whereas [82] required 12-24 hours to train the patch-based models from scratch. Thus, our results complement the work of [82] by showing that, compared to whole-image models, patch-based diffusion models are easier to train from scratch in settings of limited data, and they are also easier to fine-tune when data is very limited.

TBL A.17 – 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 for the self-supervised loss. 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 A.35 and A.36 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.

TBL A.18 – Performance of Algorithm 5 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↑	SSIM↑	PSNR↑	SSIM↑
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.4.2 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 A.37 shows some of the sample phantoms.

A.4.3 Experiment parameters

We applied the framework of [93] 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 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 5, 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 A.15 and A.14), 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

$$\operatorname{argmax}_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - A\mathbf{x}\|_2^2 + \lambda \operatorname{TV}(\mathbf{x}), \quad (\text{A.23})$$

where $\operatorname{TV}(\mathbf{x})$ represents the L1 norm total variation of \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

$$\operatorname{argmax}_{\mathbf{x}} f(\mathbf{x}) + (\rho/2) \|\mathbf{x} - (\mathbf{z} - \mathbf{u})\|_2^2, \quad (\text{A.24})$$

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

$$\mathbf{x} \leftarrow \mathbf{x} + \mu (\nabla f - \lambda(\mathbf{x} - D(\mathbf{x}))), \quad (\text{A.25})$$

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 A.19 shows the average runtimes of each of the implemented methods when averaged across the test dataset for 60 view CT reconstruction.

TBL A.19 – 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

A.5 Appendix for "SPAR: Refine a Single Pretrained Diffusion Model to Solve Inverse Problems in Many Modalities"

A.5.1 Bound of majorizer loss

We begin by showing that the quadratic majorizer loss (8.6) is indeed a majorizer. We prove a slightly more general version of this fact, stated in the following theorem.

Theorem: Suppose $\Psi(\mathbf{x}) : \mathbb{F}^N \rightarrow \mathbb{R}$ is smooth and its gradient is \mathbf{S} -Lipschitz continuous, for an invertible $N \times N$ matrix \mathbf{S} , i.e.,

$$\|\mathbf{S}^{-1}(\nabla\Psi(\mathbf{x}) - \nabla\Psi(\mathbf{z}))\|_2 \leq \|\mathbf{S}'(\mathbf{x} - \mathbf{z})\|_2, \quad \forall \mathbf{x}, \mathbf{z} \in \mathbb{R}^N, \quad (\text{A.26})$$

where \mathbb{F} is either \mathbb{R} or \mathbb{C} . Then the majorizer for $\Psi(\mathbf{x})$ defined as

$$\phi_k(\mathbf{x}) = \Psi(\mathbf{x}_k) + \text{real}\{\langle \nabla\Psi(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle\} + \frac{1}{2} \|\mathbf{S}'(\mathbf{x} - \mathbf{x}_k)\|_2^2 \quad (\text{A.27})$$

satisfies $\phi_k(\mathbf{x}) \geq \Psi(\mathbf{x})$ for all \mathbf{x} . Clearly $\phi_k(\mathbf{x}_k) = \Psi(\mathbf{x}_k)$ by construction. Thus ϕ_k is a majorizer of Ψ [105].

Proof: Expanding $\Psi(\mathbf{x})$ around \mathbf{z} and applying Taylor's theorem gives

$$\begin{aligned}
\Psi(\mathbf{x}) &= \Psi(\mathbf{z}) + \text{real} \left\{ \int_0^1 \langle \nabla \Psi(\mathbf{z} + \tau(\mathbf{x} - \mathbf{z})), \mathbf{x} - \mathbf{z} \rangle d\tau \right\} \\
&= \Psi(\mathbf{z}) + \text{real} \{ \langle \nabla \Psi(\mathbf{z}), \mathbf{x} - \mathbf{z} \rangle \} + \text{real} \\
&\quad \left\{ \int_0^1 \langle \nabla \Psi(\mathbf{z} + \tau(\mathbf{x} - \mathbf{z})) - \nabla \Psi(\mathbf{z}), \mathbf{x} - \mathbf{z} \rangle d\tau \right\} \\
&= \Psi(\mathbf{z}) + \text{real} \{ \langle \nabla \Psi(\mathbf{z}), \mathbf{x} - \mathbf{z} \rangle \} + \text{real} \left(\int_0^1 \right. \\
&\quad \left. \langle \mathbf{S}^{-1}(\nabla \Psi(\mathbf{z} + \tau(\mathbf{x} - \mathbf{z})) - \nabla \Psi(\mathbf{z})), \mathbf{S}'(\mathbf{x} - \mathbf{z}) \rangle d\tau \right).
\end{aligned}$$

Thus we have

$$|\Psi(\mathbf{x}) - \Psi(\mathbf{z}) - \text{real}\{\langle \nabla \Psi(\mathbf{z}), \mathbf{x} - \mathbf{z} \rangle\}| = \left| \text{real} \right. \quad (\text{A.28})$$

$$\left. \left\{ \int_0^1 \langle \mathbf{S}^{-1}(\nabla \Psi(\mathbf{z} + \tau(\mathbf{x} - \mathbf{z})) - \nabla \Psi(\mathbf{z})), \mathbf{S}'(\mathbf{x} - \mathbf{z}) \rangle d\tau \right\} \right| \quad (\text{A.29})$$

$$\leq \int_0^1 \left| \langle \mathbf{S}^{-1}(\nabla \Psi(\mathbf{z} + \tau(\mathbf{x} - \mathbf{z})) - \nabla \Psi(\mathbf{z})), \mathbf{S}'(\mathbf{x} - \mathbf{z}) \rangle \right| d\tau \quad (\text{A.30})$$

$$\leq \int_0^1 \|\mathbf{S}^{-1}(\nabla \Psi(\mathbf{z} + \tau(\mathbf{x} - \mathbf{z})) - \nabla \Psi(\mathbf{z}))\|_2 \|\mathbf{S}'(\mathbf{x} - \mathbf{z})\|_2 d\tau \quad (\text{A.31})$$

$$= \|\mathbf{S}'(\mathbf{x} - \mathbf{z})\|_2 \int_0^1 \|\mathbf{S}^{-1}(\nabla \Psi(\mathbf{z} + \tau(\mathbf{x} - \mathbf{z})) - \nabla \Psi(\mathbf{z}))\|_2 d\tau \quad (\text{A.32})$$

$$\leq \|\mathbf{S}'(\mathbf{x} - \mathbf{z})\|_2 \int_0^1 \|\tau \mathbf{S}'(\mathbf{x} - \mathbf{z})\|_2 d\tau \quad (\text{A.33})$$

$$= \|\mathbf{S}'(\mathbf{x} - \mathbf{z})\|_2^2 \int_0^1 \tau d\tau = \frac{1}{2} \|\mathbf{S}'(\mathbf{x} - \mathbf{z})\|_2^2, \quad (\text{A.34})$$

where we used the triangle inequality in going from the second to the third line, the Cauchy-Schwarz inequality in going from the third to the fourth line, and \mathbf{S} -Lipschitz continuity in going from the fifth to the last line. Thus:

$$\Psi(\mathbf{x}) - \Psi(\mathbf{z}) - \text{real}\{\langle \nabla \Psi(\mathbf{z}), \mathbf{x} - \mathbf{z} \rangle\} \leq \frac{1}{2} \|\mathbf{S}'(\mathbf{x} - \mathbf{z})\|_2^2, \quad (\text{A.35})$$

$$\implies \Psi(\mathbf{x}) \leq \Psi(\mathbf{z}) + \text{real}\{\langle \nabla \Psi(\mathbf{z}), \mathbf{x} - \mathbf{z} \rangle\} + \frac{1}{2} \|\mathbf{S}'(\mathbf{x} - \mathbf{z})\|_2^2. \quad (\text{A.36})$$

Finally, setting $\mathbf{z} = \mathbf{x}_k$ yields the desired result. ■

The most common use of this theorem is with $\mathbf{S} = (1/\sqrt{L})\mathbf{I}$, where L is the Lipschitz constant of the gradient of the cost function Ψ .

When applying this majorizer with the self-supervised loss (7.4), the network refining loss

$$L(\boldsymbol{\theta}; \tilde{\boldsymbol{x}}) = \langle \mathbf{A}'(\mathbf{A}\tilde{\boldsymbol{x}} - \mathbf{y}), \mathbf{x}_\theta - \tilde{\boldsymbol{x}} \rangle + \frac{L}{2} \|\mathbf{x}_\theta - \tilde{\boldsymbol{x}}\|_2^2, \quad (\text{A.37})$$

where $\mathbf{x}_\theta = D_\theta(\mathbf{x}_t|\mathbf{y})$. We choose $\tilde{\boldsymbol{x}}$ to be the denoised image from the previous diffusion iteration, and we let L be a tunable parameter. We drop the constants that do not depend on \mathbf{x}_θ .

A.5.2 Analysis of Channel Operator

This section provides a theoretical justification for the denoising score matching procedure in Section 8.2.2 used to learn the score function of single-channel images with a 3-channel network. Recall that, as described in Section 8.2.2, the operator $\mathbf{H} : \mathbb{R}^{1 \times M \times M} \rightarrow \mathbb{R}^{3 \times M \times M}$ consists of stacking a grayscale image three times along the channel dimension. We first derive a useful property about this operator. From the equality $\mathbf{H} = \mathbf{K} \otimes \mathbf{I}$ it follows that the spectral norm of \mathbf{H} (and thus of \mathbf{H}') is the same as that of $\mathbf{K} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$ which is $\sqrt{3}$.

Next, let $\boldsymbol{x} \in \mathbb{R}^{1 \times M \times M}$ be a grayscale image. Then the denoising score matching (DSM) loss is given by $\hat{L} = \|D(\boldsymbol{z}) - \boldsymbol{x}\|_2^2$, where \boldsymbol{z} is the clean image \boldsymbol{x} with some added Gaussian noise. Following (8.8), we have

$$\begin{aligned} \hat{L} &= \left\| \frac{1}{3} \mathbf{H}' \tilde{D}_\theta(\mathbf{H}\boldsymbol{z}) - \boldsymbol{x} \right\|_2^2 = \left\| \frac{1}{3} \mathbf{H}' \tilde{D}_\theta(\mathbf{H}\boldsymbol{z}) - \frac{1}{3} \mathbf{H}' \mathbf{H}\boldsymbol{x} \right\|_2^2 \\ &\leq \frac{1}{9} \|\mathbf{H}'\|_2^2 \left\| \tilde{D}_\theta(\mathbf{H}\boldsymbol{z}) - \mathbf{H}\boldsymbol{x} \right\|_2^2 = \frac{1}{3} \left\| \tilde{D}_\theta(\mathbf{H}\boldsymbol{z}) - \mathbf{H}\boldsymbol{x} \right\|_2^2, \end{aligned}$$

where we have applied the spectral norm of \mathbf{H} to get the inequality.

Because \mathbf{K} and \mathbf{H} have full column rank, $\mathbf{H}^\dagger \mathbf{H} = \mathbf{I}$, so more generally:

$$\begin{aligned} \hat{L} &= \left\| \mathbf{H}^\dagger \tilde{D}_\theta(\mathbf{H}\boldsymbol{z}) - \boldsymbol{x} \right\|_2^2 = \left\| \mathbf{H}^\dagger \left(\tilde{D}_\theta(\mathbf{H}\boldsymbol{z}) - \mathbf{H}\boldsymbol{x} \right) \right\|_2^2 \\ &\leq \left\| \mathbf{H}^\dagger \right\|_2^2 \left\| \tilde{D}_\theta(\mathbf{H}\boldsymbol{z}) - \mathbf{H}\boldsymbol{x} \right\|_2^2 = \left\| (\mathbf{K}'\mathbf{K})^{-1} \right\|_2 \left\| \tilde{D}_\theta(\mathbf{H}\boldsymbol{z}) - \mathbf{H}\boldsymbol{x} \right\|_2^2, \end{aligned}$$

where $\|(\mathbf{K}'\mathbf{K})^{-1}\|_2 = 1/3$ for the grayscale case.

In particular, the final loss on the RHS implies that we can input $\mathbf{H}\boldsymbol{z}$ (the noisy grayscale image stacked three times) into the network that accepts three-channel images, and then enforce a loss function that allows this network to denoise all three channels simultaneously to obtain $\mathbf{H}\boldsymbol{x}$. Then since the original DSM loss \hat{L} is upper bounded by the

modified loss on the RHS, by minimizing the modified loss and observing that it becomes small in practice, we have also obtained a small upper bound for the original loss.

A.5.3 Ablation studies

We performed several ablation studies to evaluate the effect of different components of our method.

Training dataset diversity. We examined the effect of pretraining the network on the dataset consisting of a mixture of CelebA images and AAPM CT images. An alternative is to pretrain the network on only CelebA images and use the same channel concatenation strategy in (8.8) at reconstruction time to handle grayscale images. Table A.20 shows the results of this experiment for three different CT settings. Note that the FBCT experiments were performed on the LIDC CT images and the mixed dataset only contained 256×256 AAPM CT images. Therefore, both the FBCT and 512×512 experiments use test data that belongs to a different distribution than the training data when using the mixed dataset. Nevertheless, for all three cases, pretraining with the mixed dataset resulted in better image quality.

TBL A.20 – Comparison of using the proposed SPAR method while training only with CelebA images versus using a mixed dataset. Best results are in bold.

Method	PBCT, 60 views		FBCT, 40 views		512 × 512 CT	
	PSNR↑	SSIM↑	PSNR↑	SSIM↑	PSNR↑	SSIM↑
CelebA only	41.93	0.969	33.67	0.869	38.57	0.927
Mixed dataset	42.72	0.972	36.11	0.918	38.81	0.929

Diffusion Inverse Solvers. Table 8.3 showed that previous DIS methods that do not account for distribution shift between training and test data performed significantly worse than methods that account for this shift. Figure A.38 shows the visual results of these experiments. Particularly, artifacts are visible for the CT experiments because the pretrained network did not have access to the LIDC CT dataset that was used at reconstruction time.

Majorizer Loss Threshold. We showed experimentally in Figure 8.4 and theoretically in Section 8.2.4 that using the majorizer self-supervised loss for small noise levels improved image quality. We perform an additional study to evaluate the impact of different threshold values σ' in Algorithm 6. Table A.21 shows the results of this experiment for

parallel beam CT with 60 views. Since a threshold of $\sigma' = 0.02$ gave the best results, we chose this value for all our main experiments as well.

TBL A.21 – Comparison between different thresholds σ' when using SPAR image reconstruction. Best results are in bold.

σ'	PBCT, 60 Views	
	PSNR \uparrow	SSIM \uparrow
0	41.92	0.967
0.01	42.41	0.971
0.02	42.72	0.972
0.05	42.39	0.971
0.1	42.36	0.970
0.2	42.35	0.970

Effect of Majorizer Loss. Figure 8.4 shows the benefit of using the majorizer loss for MRI reconstruction; we analyzed the effect of using this majorizer loss for all the inverse problems. Table A.22 shows that for each of the inverse problems examined in this paper, the majorizer loss increased the reconstructed image quality.

TBL A.22 – Comparison of using the traditional loss versus the proposed majorizer loss for each of the inverse problems.

Inverse Problem	No Majorizer		Majorizer	
	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow
PBCT, 60 views	41.92	0.967	42.72	0.972
FBCT, 40 views	35.84	0.913	36.11	0.918
512 \times 512 CT	37.54	0.909	38.81	0.929
CS-MRI, 7 \times	37.54	0.927	39.15	0.949
Deblur, LSUN	29.15	0.861	29.49	0.877
Deblur, FFHQ	30.89	0.819	31.43	0.843
Superresolution	28.10	0.827	29.10	0.842

Frequency of Network Finetuning. In Algorithm 6, the parameter K represents the frequency at which the network is finetuned during the reconstruction process. We

performed an ablation study on the effect of K on the reconstruction results for 60 view parallel-beam CT reconstruction. Table A.23 shows that $K = 10$ yielded the best result and thus we used this value for all the main experiments.

TBL A.23 – Comparison of using different K for 60 view parallel beam CT recon. Average runtime for each image is also shown.

K	PSNR (dB)	SSIM	Runtime (s)
3	41.62	0.964	697
5	42.07	0.968	465
10	42.72	0.972	288
20	42.44	0.970	197
50	42.04	0.969	146

A.5.4 Algorithm Pseudocode

Algorithm 14 shows the pseudocode for pretraining the patch-based diffusion model for our SPAR method.

Algorithm 14 SPAR Pretraining

repeat
 Select z randomly from mixed dataset with zero padding and use H as necessary to get a 3 channel image z'
 Select a random patch x from z'
 Select $t \sim \text{Uniform}[1, T]$
 Set $\epsilon \sim \mathcal{N}(0, \sigma_t^2 \mathbf{I})$
 Take gradient descent step on $\nabla_{\theta} \|D_{\theta}(x + \epsilon, \sigma_t) - x\|_2^2$
until converged
 Return D_{θ}

TBL A.24 – 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	76
PaDIS	125
Langevin dynamics	122
VE-DDNM	132
VE-DPS	241
SCD	291
Proposed (SPAR)	288

A.5.5 Additional Results

Table A.24 shows the average runtime of various methods for one image. Although our method incurs extra runtime to refine the network during reconstruction time, the experimental results show that this network refining process is crucial to obtaining a high-quality reconstruction that is free of artifacts.

Table A.25 shows the LPIPS score using the VGG network [206] for various natural image experiments.

TBL A.25 – LPIPS scores for deblurring on LSUN dataset and superresolution on CelebA dataset.

Method	Deblur	Superresolution
Baseline	0.483	0.530
ADMM-TV	0.308	0.542
PNP-ADMM	0.257	0.701
PaDIS	0.305	0.465
SCD++	0.235	0.264
Proposed (SPAR)	0.215	0.179

Table A.26 shows the results of our method and various comparison methods for natural images with a higher noise level of $\sigma = 0.05$. Our proposed method outperforms the state of the art in a variety of image modalities and settings.

TBL A.26 – Performance of deblurring on LSUN dataset and superresolution on CelebA dataset with noise level $\sigma = 0.05$. Networks were trained with CelebA images. Best results are in bold.

Method size	Deblur		Superresolution	
	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow
Baseline	20.53	0.319	22.77	0.479
DPIR [208]	25.31	0.746	27.21	0.764
DiffPIR [210]	25.09	0.699	25.61	0.663
DDRM [95]	23.85	0.683	27.33	0.772
SCD++	21.55	0.592	21.87	0.445
Proposed (SPAR)	27.16	0.797	27.65	0.789

Acceleration methods. Previous diffusion works in the literature have investigated various acceleration methods to reduce the number of required sampling steps. For solving inverse problems in particular, DDIM [165] is the most common sampler, being used in [181, 95, 210], among others. DDIM requires predicting the clean image x_0 during each sampling step and using it as a guidance direction to converge toward the desired distribution faster. However, when the network has not been trained on the same distribution of images as is currently being tested, this predicted x_0 will tend to be far away from the ground truth, decreasing the effectiveness of DDIM. We illustrate this point below by running experiments on deblurring with the LSUN dataset using DiffPIR [210] while also incorporating our network finetuning strategy. Table A.27 shows that through different selections of hyperparameters (where K is the same K as in Algorithm 6), the finetuning method fails to substantially improve performance over vanilla DiffPIR and that our proposed method still obtained the best results.

A.5.6 Experiment parameters

We applied the framework of [82] to train the patch-based networks. Since images were scaled between 0 and 1 for both grayscale images and RGB channels, we chose a maximum

TBL A.27 – Results from using DiffPIR to accelerate reconstruction process while applying our network finetuning method for deblurring.

Iterations	K	PSNR (dB)	SSIM
100	5	27.72	0.773
100	10	28.01	0.797
200	2	28.33	0.814
200	10	28.38	0.818
100	None	28.11	0.798
200	None	28.36	0.819
SPAR	10	29.49	0.877

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.

We used a learning rate of $2 \cdot 10^{-4}$ to train the networks from scratch. 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 of 32×32 and 16×16 . The probabilities of using these three patch sizes were 0.5, 0.3, and 0.2 respectively.

For inverse problem solving, we used a geometrically spaced descending noise level that was fine tuned to optimize the performance for each type of problem. The values are as follows:

- CT: $\sigma_{\max} = 0.2, \sigma_{\min} = 0.005$
- MRI: $\sigma_{\max} = 0.4, \sigma_{\min} = 0.005$
- Deblurring: $\sigma_{\max} = 0.3, \sigma_{\min} = 0.005$
- Superresolution: $\sigma_{\max} = 0.4, \sigma_{\min} = 0.01$.

When running Algorithm 6, 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} . The value of σ' was set to 0.02 for all experiments.

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

$$\operatorname{argmax}_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \operatorname{TV}(\mathbf{x}), \quad (\text{A.38})$$

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

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

$$\operatorname{argmax}_{\mathbf{x}} f(\mathbf{x}) + (\rho/2) \|\mathbf{x} - (\mathbf{z} - \mathbf{u})\|_2^2, \quad (\text{A.39})$$

where ρ is a constant. The values for ρ we used for CT reconstruction, MRI, deblurring, and superresolution were 0.05, 0.07, 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.

Acknowledgement

Part of this thesis is based upon work supported by the National Science Foundation under Award No. IIS-2435746. Any opinions, findings, and conclusions or recommendations expressed in this thesis are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.

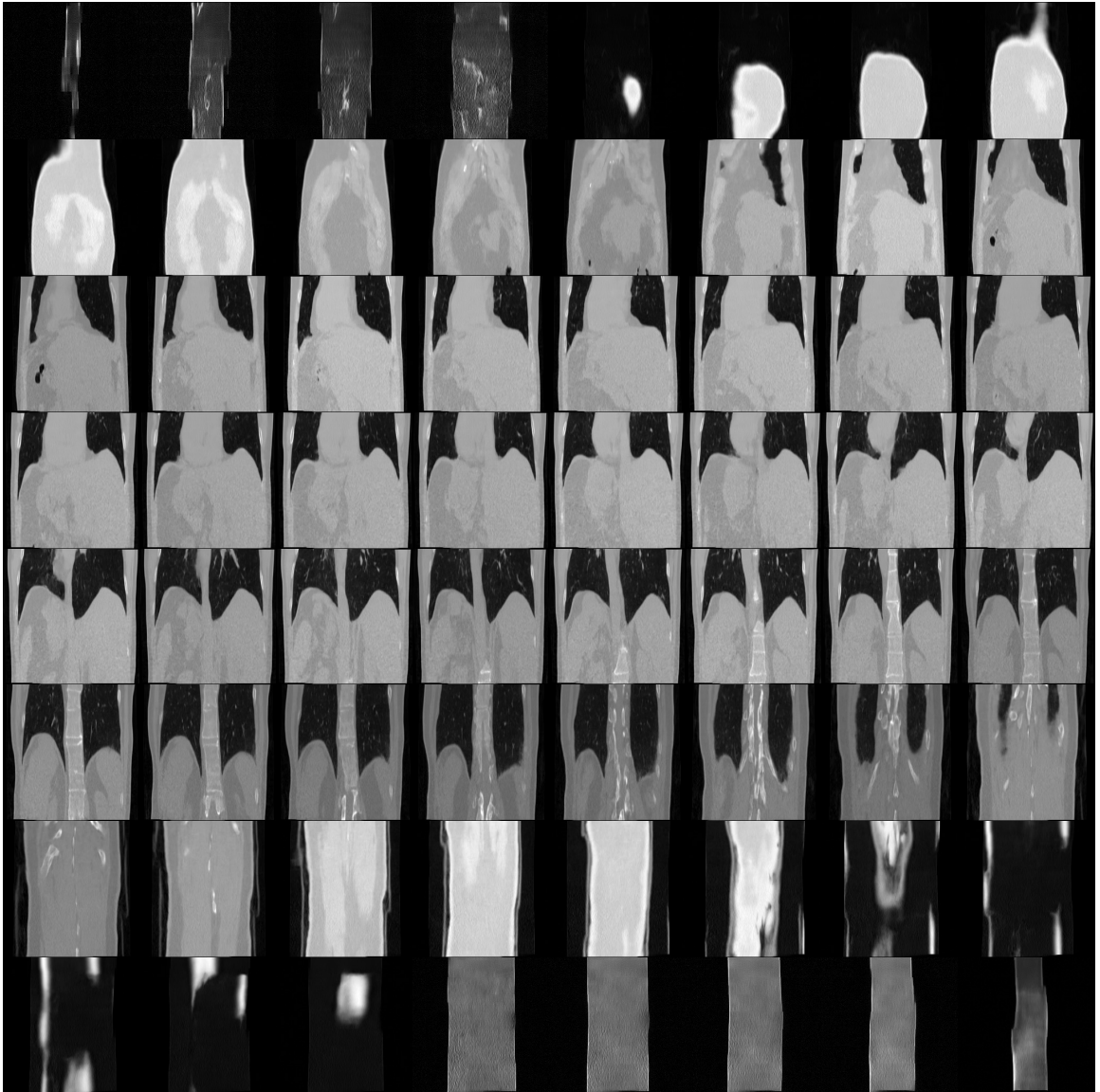


FIG A.29 – Continuous coronal slice of a volume sampled from $256 \times 256 \times 256$ LIDC-IDRI prior; every fourth slice is shown

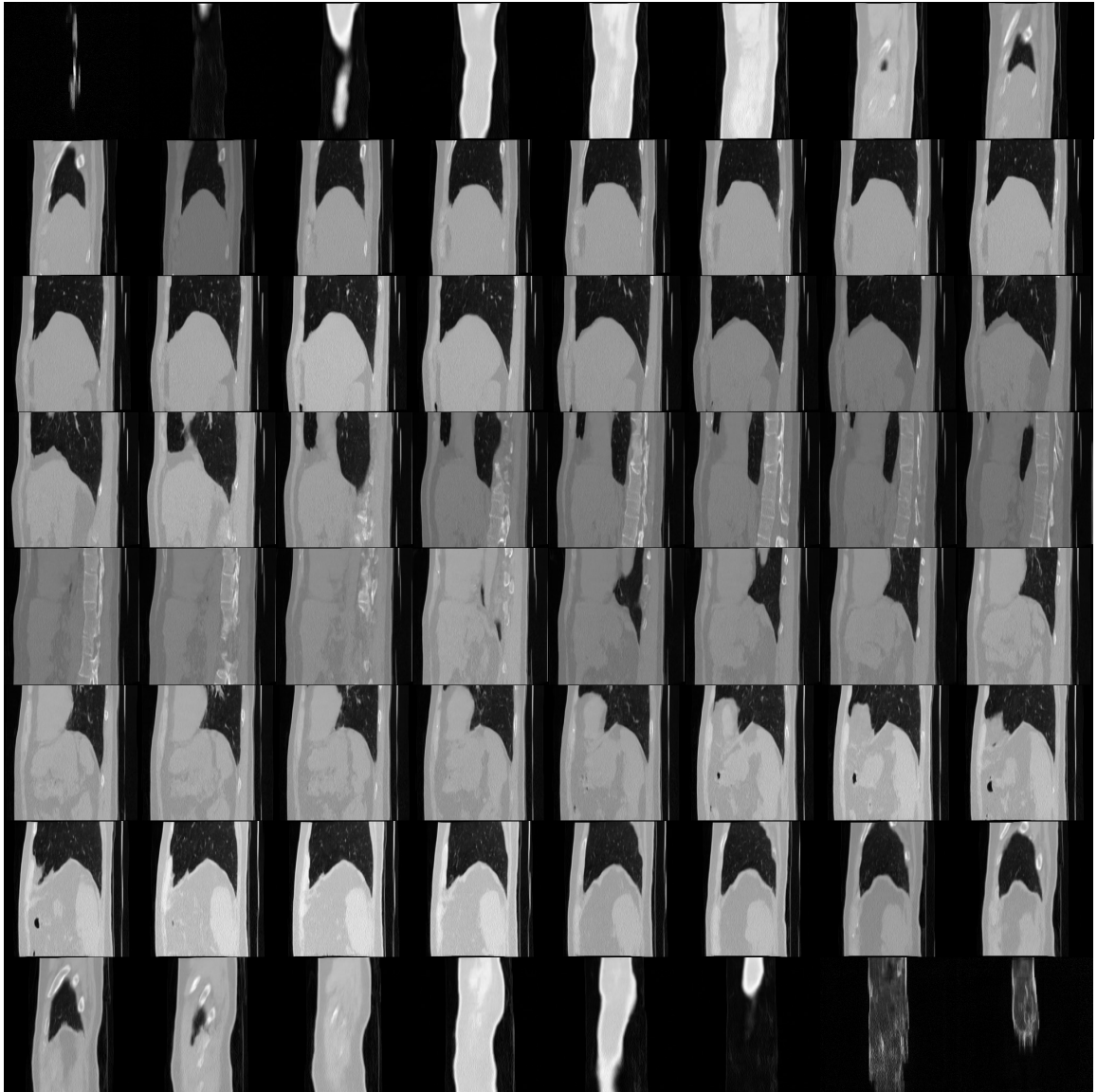


FIG A.30 – Continuous sagittal slice of a volume sampled from $256 \times 256 \times 256$ LIDC-IDRI prior; every fourth slice is shown

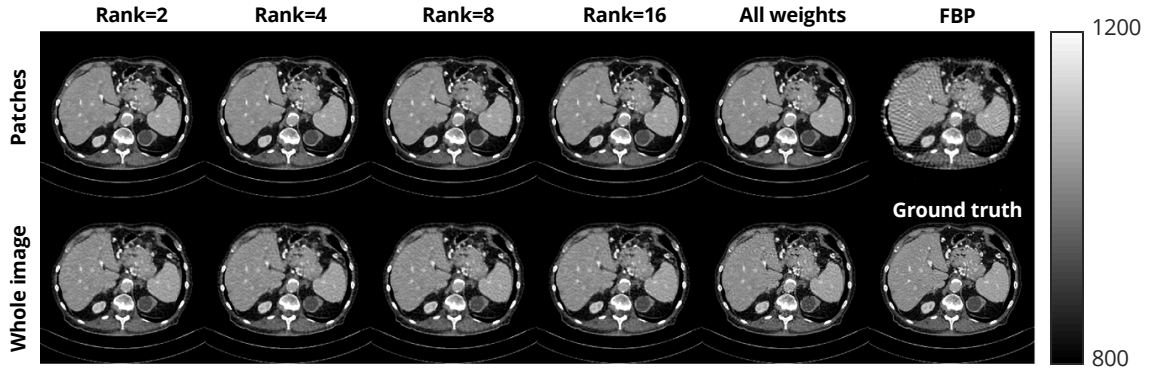


FIG A.31 – Results of using LoRA module for 60 view CT reconstruction in a single measurement setting. All weights refers to adjusting all the weights of the network at reconstruction time.

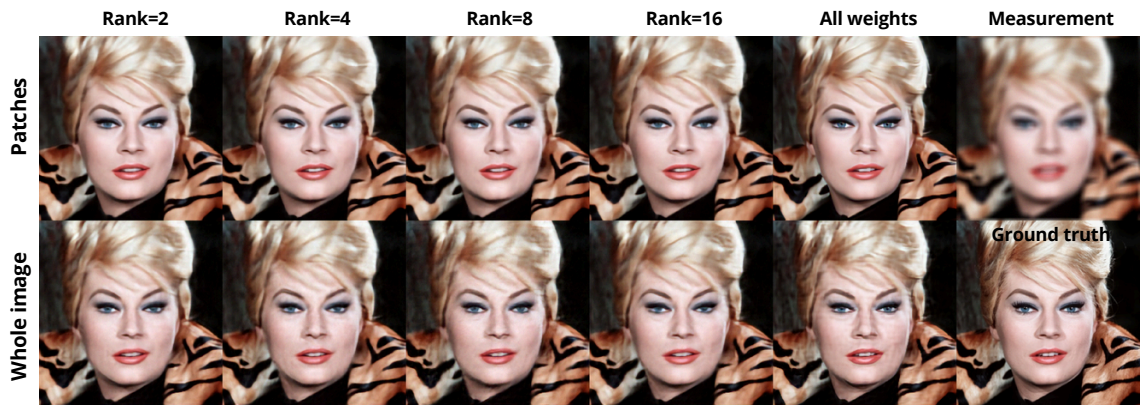


FIG A.32 – Results of using LoRA module for deblurring in a single measurement setting. All weights refers to adjusting all the weights of the network at reconstruction time.

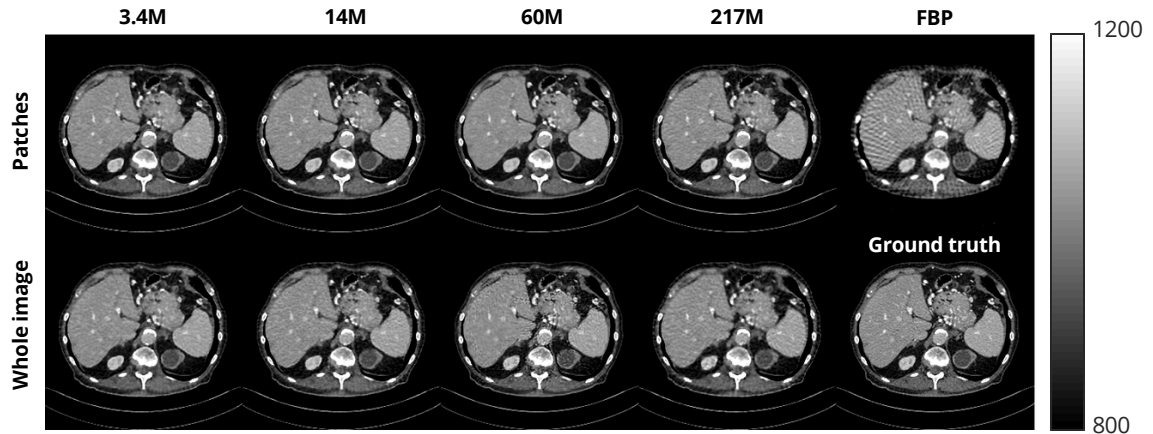


FIG A.33 – Results of 60 view CT recon using networks with different numbers of parameters in the single-measurement setting. The top numbers show the number of total parameters in the network.

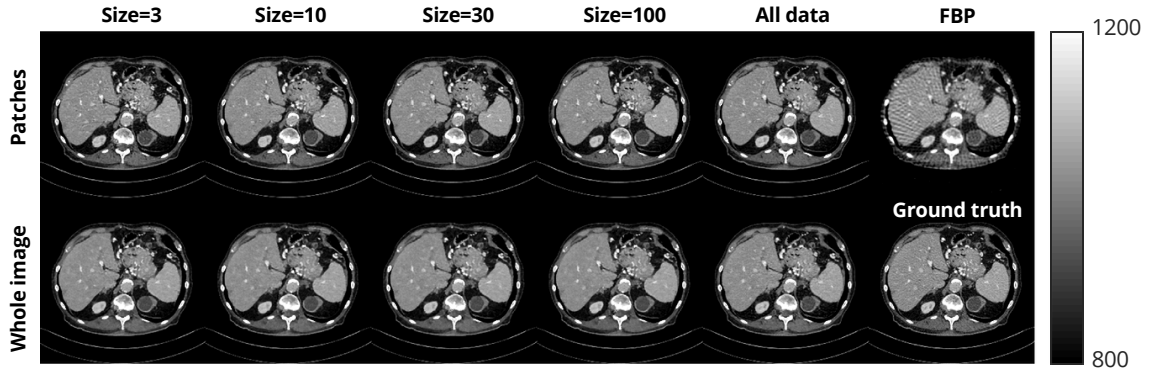


FIG A.34 – Results of 60 view CT recon in the small dataset setting where the size of the small dataset is varied.

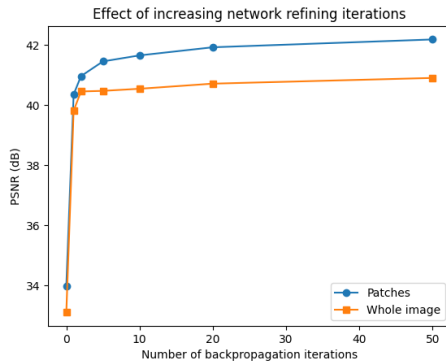


FIG A.35 – Comparison of PSNR between patch-based model and whole-image model for number of network refining iterations in single measurement setting.

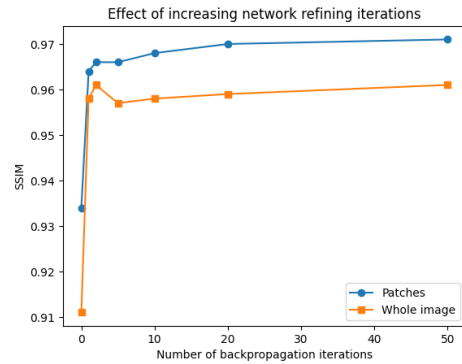
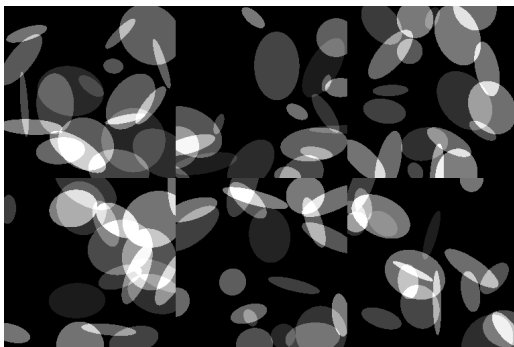
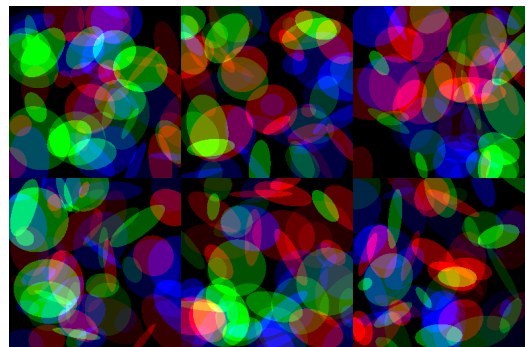


FIG A.36 – Comparison of SSIM between patch-based model and whole-image model for number of network refining iterations in single measurement setting.



(a) Six grayscale phantoms



(b) Six colored phantoms

FIG A.37 – Six sample grayscale phantoms and colored phantoms used to train the mismatched distribution diffusion models

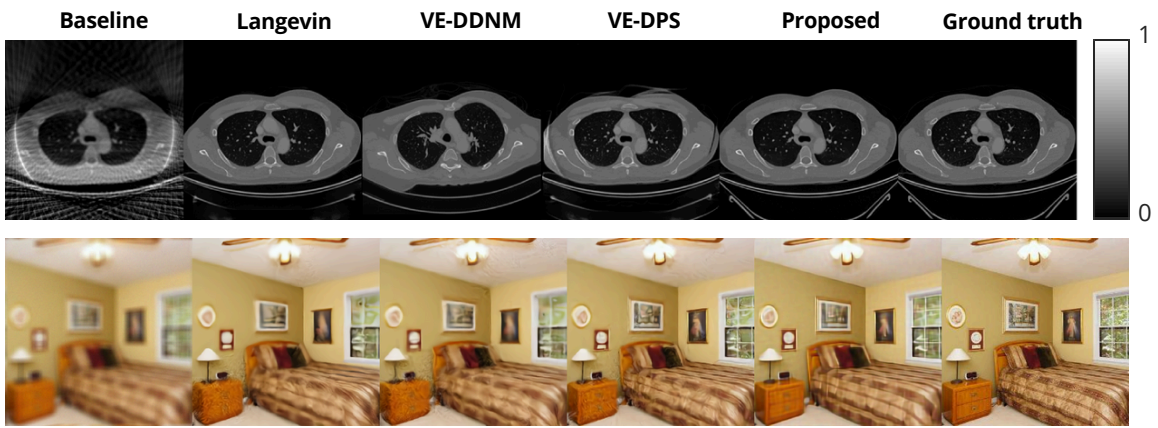


FIG A.38 – Visual results of using various DIS to solve FBCT with 40 views (top) and deblurring (bottom). The same single pretrained network was used in all the experiments. Arrows point to some artifacts in the comparison methods.

BIBLIOGRAPHY

- [1] R. Ahmad, C. A. Bouman, G. T. Buzzard, S. Chan, S. Liu, E. T. Reehorst, and P. Schniter. “Plug-and-play methods for magnetic resonance imaging: Using denoisers for image recovery”. In: *IEEE Signal Processing Magazine* 37.1 (2020), pp. 105–116 (cit. on p. 10).
- [2] A. Aksac, D. Demetrick, T. Ozyer, *et al.*. “BreCaHAD: a dataset for breast cancer histopathological annotation and diagnosis”. In: *MC Res Notes* 12.82 (2019). DOI: <https://doi.org/10.1186/s13104-019-4121-7> (cit. on pp. 22–24, 26–28, 32).
- [3] B. D. Anderson. “Reverse-time diffusion equation models”. In: *Stochastic Processes and their Applications* 12.3 (1982), pp. 313–326. ISSN: 0304-4149. DOI: [https://doi.org/10.1016/0304-4149\(82\)90051-5](https://doi.org/10.1016/0304-4149(82)90051-5) (cit. on p. 10).
- [4] V. Antun, F. Renna, C. Poon, B. Adcock, and A. C. Hansen. “On instabilities of deep learning in image reconstruction and the potential costs of AI”. In: *Proceedings of the National Academy of Sciences* 117.48 (2020), pp. 30088–30095 (cit. on p. 10).
- [5] Y. Arefeen, B. Levac, B. Patel, C. Ho, and J. I. Tamir. “Diffusion probabilistic generative models for accelerated, in-NICU permanent magnet neonatal MRI”. In: *Mag. Res. Med.* 94.4 (Oct. 2025), 1546–62. DOI: [10.1002/mrm.30585](https://doi.org/10.1002/mrm.30585) (cit. on p. 101).
- [6] S. G. Armato, G. McLennan, L. Bidaut, M. F. McNitt-Gray, and C. R. Meyer. “The lung image database consortium (LIDC) and image database resource initiative (IDRI): A completed reference database of lung nodules on CT scans.” In: *Medical Physics* 38 (2011), pp. 915–931 (cit. on pp. 61, 76, 96, 112).
- [7] M. Asim, M. Daniels, O. Leong, A. Ahmed, and P. Hand. “Invertible generative models for inverse problems: mitigating representation error and dataset bias”. In: *Proceedings of the 37th International Conference on Machine Learning* 119 (2020) (cit. on p. 16).
- [8] W. Bai, S. Xu, Y. Ren, J. Hao, M. Sun, W. Chen, and H. Sun. *InstantViR: Real-Time Video Inverse Problem Solver with Distilled Diffusion Prior*. 2025. arXiv: [2511.14208](https://arxiv.org/abs/2511.14208) [cs.CV] (cit. on pp. 120, 121).

- [9] O. Bar-Tal, L. Yariv, Y. Lipman, and T. Dekel. *MultiDiffusion: Fusing Diffusion Paths for Controlled Image Generation*. 2023. arXiv: 2302.08113 [cs.CV] (cit. on p. 12).
- [10] R. Barbano, A. Denker, H. Chung, T. H. Roh, S. Arrdige, P. Maass, B. Jin, and J. C. Ye. *Steerable Conditional Diffusion for Out-of-Distribution Adaptation in Imaging Inverse Problems*. 2023. arXiv: 2308.14409 [cs.CV] (cit. on pp. 14, 83, 86, 93, 101, 111, 113, 163).
- [11] R. Barbano, J. Leuschner, M. Schmidt, A. Denker, A. Hauptmann, P. Maass, and B. Jin. “An Educated Warm Start for Deep Image Prior-Based Micro CT Reconstruction”. In: *IEEE Transactions on Computational Imaging* 8 (2022), pp. 1210–1222. ISSN: 2573-0436. DOI: 10.1109/tci.2022.3233188 (cit. on p. 14).
- [12] D. A. Barmherzig and J. Sun. “Towards practical holographic coherent diffraction imaging via maximum likelihood estimation”. In: *Opt. Express* 30.5 (Feb. 2022), pp. 6886–6906. DOI: 10.1364/OE.445015 (cit. on pp. 23, 27).
- [13] D. A. Barmherzig, J. Sun, E. J. Candes, T. Lane, and P.-N. Li. “Dual-Reference Design for Holographic Phase Retrieval”. In: *2019 13th International conference on Sampling Theory and Applications (SampTA)*. 2019, pp. 1–4. DOI: 10.1109/SampTA45681.2019.9030848 (cit. on p. 8).
- [14] D. A. Barmherzig, J. Sun, P.-N. Li, T. Lane, and E. J. Candès. “Holographic phase retrieval and reference design”. In: *Inverse problems* 35.9 (), pp. 94001–. DOI: 10.1088/1361-6420/ab23d1 (cit. on p. 8).
- [15] J. Batson and L. Royer. “Noise2Self: Blind Denoising by Self-Supervision”. In: *Proceedings of the 36th International Conference on Machine Learning*. Ed. by K. Chaudhuri and R. Salakhutdinov. Vol. 97. Proceedings of Machine Learning Research. PMLR, June 2019, pp. 524–533 (cit. on pp. 16, 23, 28).
- [16] J. Besag. “Spatial Interaction and the Statistical Analysis of Lattice Systems”. In: *Journal of the Royal Statistical Society* 36.2 (1974), pp. 192–236 (cit. on p. 133).
- [17] L. Bian, J. Suo, J. Chung, X. Ou, C. Yang, F. Chen, and Q. Dai. “Fourier ptychographic reconstruction using Poisson maximum likelihood and truncated Wirtinger gradient”. In: *Nature Sci. Rep.* 6.1 (2016). DOI: 10.1038/srep27384 (cit. on p. 8).
- [18] F. Bieder, J. Wolleb, A. Durrer, R. Sandkuehler, and P. C. Cattin. “Denoising diffusion models for memory-efficient processing of 3D medical images”. In: *Proc. Mach. Learning Res.* 227 (2024), 552–67 (cit. on p. 12).

- [19] A. Blattmann, T. Dockhorn, S. Kulal, D. Mendeleevitch, M. Kilian, D. Lorenz, Y. Levi, Z. English, V. Voleti, A. Letts, V. Jampani, and R. Rombach. *Stable Video Diffusion: Scaling Latent Video Diffusion Models to Large Datasets*. 2023. arXiv: 2311.15127 [cs.CV] (cit. on p. 121).
- [20] A. Bora, A. Jalal, E. Price, and A. G. Dimakis. “Compressed Sensing using Generative Models”. In: *Proceedings of the 34th International Conference on Machine Learning*. Vol. 70. Proceedings of Machine Learning Research. PMLR, 2017, pp. 537–546 (cit. on p. 16).
- [21] E. Bostan, M. Soltanolkotabi, D. Ren, and L. Waller. “Accelerated Wirtinger Flow for Multiplexed Fourier Ptychographic Microscopy”. In: *2018 25th IEEE International Conference on Image Processing (ICIP)*. 2018, pp. 3823–3827. DOI: 10.1109/ICIP.2018.8451437 (cit. on p. 9).
- [22] K. Bredies and M. Holler. “Regularization of linear inverse problems with total generalized variation”. In: *Journal of Inverse and Ill-Posed Problems* 22.6 (2014), pp. 871–913. DOI: 10.1515/jip-2013-0068 (cit. on p. 16).
- [23] K. Bredies and T. Valkonen. “Inverse Problems with Second-Order Total Generalized Variation Constraints”. In: *arXiv preprint arXiv:2005.09725* (2020). DOI: 10.48550/arxiv.2005.09725 (cit. on p. 16).
- [24] T. M. Buzug. “Computed tomography”. In: *Springer handbook of medical technology*. Springer, 2011, pp. 311–342 (cit. on p. 9).
- [25] J.-F. Cai, M. Huang, D. Li, and Y. Wang. “Nearly optimal bounds for the global geometric landscape of phase retrieval”. In: *Inverse Problems* 39.7 (June 2023), p. 075011. DOI: 10.1088/1361-6420/acdab7 (cit. on p. 32).
- [26] T. T. Cai, X. Li, and Z. Ma. “Optimal Rates of Convergence for Noisy Sparse Phase Retrieval via Thresholded Wirtinger Flow”. In: *Annals Stat.* 44.5 (2016), pp. 2221–2251. DOI: 10.1214/16-AOS1443 (cit. on p. 9).
- [27] E. J. Candes, Y. C. Eldar, T. Strohmer, and V. Voroninski. “Phase retrieval via matrix completion”. In: *SIAM J. Imaging Sci.* 6.1 (2013), 199–225. DOI: 10.1137/110848074 (cit. on pp. 9, 31).
- [28] E. Candes, X. Li, and M. Soltanolkotabi. “Phase Retrieval via Wirtinger Flow: Theory and Algorithms”. In: *IEEE Trans. Info. Theory* 61.4 (Apr. 2015), pp. 1985–2007. DOI: 10.1109/TIT.2015.2399924 (cit. on pp. 9, 18, 23, 27, 31).
- [29] E. J. Candes and M. B. Wakin. “An Introduction To Compressive Sampling”. In: *IEEE Signal Processing Magazine* 25.2 (2008), pp. 21–30. DOI: 10.1109/MSP.2007.914731 (cit. on p. 1).

- [30] E. J. Candès, T. Strohmer, and V. Voroninski. “PhaseLift: Exact and Stable Signal Recovery from Magnitude Measurements via Convex Programming”. In: *Comm. Pure Appl. Math.* 66.8 (2013), pp. 1241–1274. DOI: 10 . 1002 / cpa . 21432 (cit. on pp. 9, 31).
- [31] G. Cardoso, Y. J. E. Idrissi, S. L. Corff, and E. Moulines. *Monte Carlo guided Diffusion for Bayesian linear inverse problems*. 2023. arXiv: 2308 . 07983 [stat.ML] (cit. on p. 11).
- [32] G. Cardoso, Y. J. E. Idrissi, S. L. Corff, and E. Moulines. *Monte Carlo guided Diffusion for Bayesian linear inverse problems*. 2023. arXiv: 2308 . 07983 [stat.ML] (cit. on p. 120).
- [33] P. Cascarano, A. Benfenati, U. S. Kamilov, and X. Xu. *Constrained Regularization by Denoising with Automatic Parameter Selection*. 2024. arXiv: 2311 . 03819 [math.OC] (cit. on p. 8).
- [34] S. H. Chan, X. Wang, and O. A. Elgendy. “Plug-and-Play ADMM for Image Restoration: Fixed-Point Convergence and Applications”. In: *IEEE Transactions on Computational Imaging* 3.1 (2017), pp. 84–98. DOI: 10 . 1109 / TCI . 2016 . 2629286 (cit. on pp. 7, 16).
- [35] H. Chang, Y. Lou, Y. Duan, and S. Marchesini. “Total variation–based phase retrieval for Poisson noise removal”. In: *SIAM J. Imag. Sci.* 11.1 (2018), pp. 24–55. DOI: 10 . 1137 / 16M1103270 (cit. on p. 7).
- [36] J. Chang, C. Duan, Y. Jiao, R. Li, J. Z. Yang, and C. Yuan. *Provable Diffusion Posterior Sampling for Bayesian Inversion*. 2025. arXiv: 2512 . 08022 [stat.ML] (cit. on p. 120).
- [37] J. H. Cho and J. A. Fessler. “Regularization designs for uniform spatial resolution and noise properties in statistical image reconstruction for 3D X-ray CT”. In: *IEEE Trans. Med. Imag.* 34.2 (Feb. 2015), pp. 678–689 (cit. on p. 9).
- [38] E. Chouzenoux, A. Jezierska, J.-C. Pesquet, and H. Talbot. “A Convex Approach for Image Restoration with Exact Poisson–Gaussian Likelihood”. In: *SIAM Journal on Imaging Sciences* 8.4 (2015), pp. 2662–2682. DOI: 10 . 1137 / 15M1014395 (cit. on pp. 18–20).
- [39] H. Chung, J. Kim, M. T. Mccann, M. L. Klasky, and J. C. Ye. “Diffusion Posterior Sampling for General Noisy Inverse Problems”. In: *The Eleventh International Conference on Learning Representations*. 2023 (cit. on pp. 11, 32, 40, 50, 60, 69, 82, 87, 101, 113–115, 120, 124).
- [40] H. Chung, J. Kim, and J. C. Ye. *Direct Diffusion Bridge using Data Consistency for Inverse Problems*. 2023. arXiv: 2305 . 19809 [cs.CV] (cit. on pp. 11, 119, 120).

- [41] H. Chung, S. Lee, and J. C. Ye. *Decomposed Diffusion Sampler for Accelerating Large-Scale Inverse Problems*. 2024. arXiv: 2303.05754 [cs.LG] (cit. on pp. 10, 12, 13, 55, 61, 63, 70, 74, 75, 78, 86, 101, 109, 120, 151).
- [42] H. Chung, D. Ryu, M. T. McCann, M. L. Klasky, and J. C. Ye. *Solving 3D Inverse Problems using Pre-trained 2D Diffusion Models*. 2022. arXiv: 2211.10655 [cs.CV] (cit. on pp. 13, 34, 55, 59, 60, 63, 68, 70, 77, 101, 142, 150).
- [43] H. Chung, B. Sim, D. Ryu, and J. C. Ye. *Improving Diffusion Models for Inverse Problems using Manifold Constraints*. 2022. arXiv: 2206.00941 [cs.LG] (cit. on pp. 11, 82).
- [44] H. Chung, B. Sim, D. Ryu, and J. C. Ye. “Improving Diffusion Models for Inverse Problems using Manifold Constraints”. In: *Advances in Neural Information Processing Systems*. Vol. 35. 2022, pp. 25683–25696 (cit. on pp. 13, 63).
- [45] H. Chung, B. Sim, and J. C. Ye. *Come-Closer-Diffuse-Faster: Accelerating Conditional Diffusion Models for Inverse Problems through Stochastic Contraction*. 2022. arXiv: 2112.05146 [eess.IV] (cit. on pp. 69, 107, 124).
- [46] H. Chung and J. C. Ye. “Score-based diffusion models for accelerated MRI”. In: *Medical Image Analysis* 80 (2022), p. 102479. ISSN: 1361-8415. DOI: <https://doi.org/10.1016/j.media.2022.102479> (cit. on pp. 10, 17, 40, 50, 51, 124, 132).
- [47] H. Chung and J. C. Ye. *Deep Diffusion Image Prior for Efficient OOD Adaptation in 3D Inverse Problems*. 2024. arXiv: 2407.10641 [cs.CV] (cit. on pp. 14, 15, 75, 83, 86, 93, 101, 102, 105, 107, 118, 163).
- [48] Z.-X. Cui, C. Cao, S. Liu, Q. Zhu, J. Cheng, H. Wang, Y. Zhu, and D. Liang. *Self-Score: Self-Supervised Learning on Score-Based Models for MRI Reconstruction*. 2022. arXiv: 2209.00835 [eess.IV] (cit. on p. 17).
- [49] K. Dabov, A. Foi, V. Katkovnik, and K. Egiazarian. “Image denoising with block-matching and 3D filtering”. In: *Proc. SPIE 6064, Image Processing: Algorithms and Systems, Neural Networks, and Machine Learning*. Vol. 6064. Feb. 2006, p. 606414. DOI: 10.1117/12.643267 (cit. on pp. 14, 92, 114).
- [50] J. Dainty and J. Fienup. “Phase retrieval and image reconstruction for astronomy”. In: *Imag. Recov. Theory Appl.* 13 (Jan. 1987), pp. 231–275 (cit. on p. 8).
- [51] G. Daras, H. Chung, C.-H. Lai, Y. Mitsufuji, J. C. Ye, P. Milanfar, A. G. Dimakis, and M. Delbracio. *A survey on diffusion models for inverse problems*. 2024 (cit. on p. 101).

- [52] I. Daubechies. *Ten Lectures on Wavelets*. Society for Industrial and Applied Mathematics, 1992. doi: 10 . 1137 / 1 . 9781611970104. eprint: <https://epubs.siam.org/doi/pdf/10.1137/1.9781611970104> (cit. on pp. 7, 14, 16).
- [53] M. Delbracio and P. Milanfar. *Inversion by Direct Iteration: An Alternative to Denoising Diffusion for Image Restoration*. 2024. arXiv: 2303 . 11435 [eess . IV] (cit. on p. 11).
- [54] P. Dhariwal and A. Nichol. “Diffusion Models Beat GANs on Image Synthesis”. In: *Advances in Neural Information Processing Systems*. Vol. 34. 2021, pp. 8780–8794 (cit. on pp. 16, 90).
- [55] Z. Ding, M. Zhang, J. Wu, and Z. Tu. *Patched Denoising Diffusion Models For High-Resolution Image Synthesis*. 2023. arXiv: 2308 . 01316 [cs . CV] (cit. on pp. 12, 38, 83).
- [56] B. Efron. “Tweedie’s formula and selection bias”. In: *Journal of the American Statistical Association* 106.496 (2011), pp. 1602–1614 (cit. on pp. 38, 39, 61, 104).
- [57] P. Esser, R. Rombach, and B. Ommer. *Taming Transformers for High-Resolution Image Synthesis*. 2021. arXiv: 2012 . 09841 [cs . CV] (cit. on p. 12).
- [58] Z. Fabian, J. Haldar, R. Leahy, and M. Soltanolkotabi. “3D Phase Retrieval at Nano-Scale via Accelerated Wirtinger Flow”. In: *2020 28th European Signal Processing Conference (EUSIPCO)*. 2021, pp. 2080–2084. doi: 10 . 23919/Eusipco47968 . 2020 . 9287703 (cit. on p. 9).
- [59] L. A. Feldkamp, L. C. Davis, and J. W. Kress. “Practical cone beam algorithm”. In: *J. Opt. Soc. Am. A* 1.6 (June 1984), pp. 612–619 (cit. on pp. 9, 68).
- [60] B. T. Feng, R. Baptista, and K. L. Bouman. *Neural Approximate Mirror Maps for Constrained Diffusion Models*. 2024. arXiv: 2406 . 12816 [cs . LG] (cit. on pp. 13, 82, 120).
- [61] B. T. Feng, J. Smith, M. Rubinstein, H. Chang, K. L. Bouman, and W. T. Freeman. *Score-Based Diffusion Models as Principled Priors for Inverse Imaging*. 2023. arXiv: 2304 . 11751 [cs . CV] (cit. on pp. 13, 82, 83).
- [62] J. Fienup. “Phase Retrieval Algorithms: A Comparison”. In: *Appl. Opt.* 21 (1982) (cit. on pp. 9, 30).
- [63] Y. Gao, F. Yang, and L. Cao. “Pixel Super-Resolution Phase Retrieval for Lensless On-Chip Microscopy via Accelerated Wirtinger Flow”. In: *Cells* 11.13 (2022). doi: 10 . 3390/cells11131999 (cit. on p. 9).

- [64] R. W. Gerchberg and W. O. Saxton. “Practical Algorithm for Determination of Phase from Image and Diffraction Plane Pictures”. In: *OPTIK* 35.2 (1972), 237–246. ISSN: 0030-4026 (cit. on pp. 9, 30).
- [65] T. Goldstein and S. Osher. “The split Bregman method for L1-regularized problems”. In: *SIAM journal on imaging sciences* 2.2 (2009), pp. 323–343 (cit. on p. 63).
- [66] A. Graikos, N. Malkin, N. Jojic, and D. Samaras. “Diffusion Models as Plug-and-Play Priors”. In: *Advances in Neural Information Processing Systems*. 2022 (cit. on p. 17).
- [67] M. A. Griswold, P. M. Jakob, R. M. Heidemann, M. Nittka, V. Jellus, J. Wang, B. Kiefer, and A. Haase. “Generalized autocalibrating partially parallel acquisitions (GRAPPA)”. In: *Magnetic Resonance in Medicine* 47.6 (2002), pp. 1202–1210 (cit. on p. 10).
- [68] J. Gu, S. Zhai, Y. Zhang, J. M. Susskind, and N. Jaitly. “Matryoshka diffusion models”. In: *Proc. Intl. Conf. on Learning Representations*. 2024 (cit. on p. 54).
- [69] J. Gu, S. Zhai, Y. Zhang, J. Susskind, and N. Jaitly. *Matryoshka Diffusion Models*. 2024. arXiv: 2310.15111 [cs.CV] (cit. on p. 12).
- [70] M. Guizar-Sicairos and J. R. Fienup. “Holography with extended reference by autocorrelation linear differential operation”. In: *Opt. Express* 15.26 (2007), pp. 17592–17612 (cit. on p. 8).
- [71] T. Gureyev, A. Pogany, D. Paganin, and S. Wilkins. “Linear algorithms for phase retrieval in the Fresnel region”. In: *Optics Communications* 231.1 (2004), pp. 53–70. ISSN: 0030-4018. DOI: <https://doi.org/10.1016/j.optcom.2003.12.020> (cit. on pp. 13, 82).
- [72] Y. He, S. Yang, H. Chen, X. Cun, M. Xia, Y. Zhang, X. Wang, R. He, Q. Chen, and Y. Shan. “ScaleCrafter: tuning-free higher-resolution visual generation with diffusion models”. In: *Proc. Intl. Conf. on Learning Representations*. 2024 (cit. on pp. 13, 54).
- [73] J. Ho, A. Jain, and P. Abbeel. “Denoising Diffusion Probabilistic Models”. In: 33 (2020). Ed. by H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, pp. 6840–6851 (cit. on pp. 13, 16, 22, 34, 36, 39, 55, 63, 82, 88, 100, 101, 105, 119, 135).
- [74] T. Hong, L. Hernandez-Garcia, and J. A. Fessler. “A Complex Quasi-Newton Proximal Method for Image Reconstruction in Compressed Sensing MRI”. In: *IEEE Transactions on Computational Imaging* 10 (2024), pp. 372–384. ISSN: 2573-0436. DOI: 10.1109/tci.2024.3369404 (cit. on pp. 50, 78, 92, 114).

- [75] T. Hong, X. Xu, J. Hu, and J. A. Fessler. “Provable Preconditioned Plug-and-Play Approach for Compressed Sensing MRI Reconstruction”. In: *arXiv preprint arXiv:2405.03854* (2024) (cit. on pp. 10, 14, 113).
- [76] T. Hong, Z. Xu, J. Hu, and J. A. Fessler. “Using Randomized Nyström Preconditioners to Accelerate Variational Image Reconstruction”. In: *IEEE Transactions on Computational Imaging* 11 (2025), pp. 1630–1643. DOI: 10.1109/TCI.2025.3622903 (cit. on p. 14).
- [77] T. Hong, I. Yavneh, and M. Zibulevsky. “Solving RED With Weighted Proximal Methods”. In: *IEEE Signal Processing Letters* 27 (2020), pp. 501–505. DOI: 10.1109/LSP.2020.2979062 (cit. on p. 14).
- [78] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. *LoRA: Low-Rank Adaptation of Large Language Models*. 2021. arXiv: 2106.09685 [cs.CL] (cit. on pp. 15, 86).
- [79] J. Hu, Z. Li, B. Song, L. Shen, and J. A. Fessler. “SPAR: Refine a single pretrained diffusion model to solve inverse problems in many modalities”. In: *IEEE Trans. Computational Imaging* (2025). Submitted (cit. on pp. 3, 100).
- [80] J. Hu, B. T.-W. Lin, J. H. Vega, and N. R.-L. Tsiang. “Predictive Models of Driver Deceleration and Acceleration Responses to Lead Vehicle Cutting In and Out”. In: *Transportation Research Record* 2677.5 (2023), pp. 92–102. DOI: 10.1177/03611981221128277. eprint: <https://doi.org/10.1177/03611981221128277> (cit. on pp. 7, 9).
- [81] J. Hu, B. Song, J. A. Fessler, and L. Shen. “Test-Time Adaptation Improves Inverse Problem Solving With Patch-Based Diffusion Models”. In: *IEEE Transactions on Computational Imaging* 11 (2025), pp. 980–991. DOI: 10.1109/TCI.2025.3587407 (cit. on pp. 3, 4, 41, 82, 102, 103, 113).
- [82] J. Hu, B. Song, X. Xu, L. Shen, and J. A. Fessler. *Learning Image Priors through Patch-based Diffusion Models for Solving Inverse Problems*. 2024. arXiv: 2406.02462 [cs.CV] (cit. on pp. 3, 4, 34, 41, 46, 71, 75, 83, 84, 86, 87, 91, 93, 94, 103, 105, 110, 111, 113, 152, 165, 176).
- [83] Y. Hu, J. Liu, X. Xu, and U. S. Kamilov. *Monotonically Convergent Regularization by Denoising*. 2022. arXiv: 2202.04961 [eess.IV] (cit. on pp. 8, 50, 51, 92, 93).
- [84] P. J. Huber. *Robust statistics*. New York: Wiley, 1981 (cit. on p. 23).
- [85] K. Jaganathan, Y. C. Eldar, and B. Hassibi. *Phase retrieval: an overview of recent developments*. <https://arxiv.org/abs/1510.07713>. 2015 (cit. on p. 8).

- [86] A. Jalal, M. Arvinte, G. Daras, E. Price, A. G. Dimakis, and J. Tamir. “Robust Compressed Sensing MRI with Deep Generative Priors”. In: *Advances in Neural Information Processing Systems*. Vol. 34. Curran Associates, Inc., 2021, pp. 14938–14954 (cit. on pp. 13, 17, 82, 83).
- [87] Y. Janati, A. Durmus, E. Moulines, and J. Olsson. *Divide-and-Conquer Posterior Sampling for Denoising Diffusion Priors*. 2024. arXiv: 2403.11407 [stat.ML] (cit. on p. 12).
- [88] X. Jiang, S. Rajan, and X. Liu. “Wirtinger Flow Method With Optimal Stepsize for Phase Retrieval”. In: *IEEE Sig. Proc. Letters* 23.11 (2016), pp. 1627–1631. DOI: 10.1109/LSP.2016.2611940 (cit. on p. 9).
- [89] K. H. Jin, M. T. McCann, E. Froustey, and M. Unser. “Deep convolutional neural network for inverse problems in imaging”. In: *IEEE Transactions on Image Processing* 26.9 (2017), pp. 4509–4522 (cit. on pp. 8, 10, 63, 78, 150).
- [90] Y. Jo, S. Y. Chun, and J. Choi. *Rethinking Deep Image Prior for Denoising*. 2021. arXiv: 2108.12841 [eess.IV] (cit. on p. 14).
- [91] U. S. Kamilov, H. Mansour, and B. Wohlberg. “A Plug-and-Play Priors Approach for Solving Nonlinear Imaging Inverse Problems”. In: *IEEE Signal. Proc. Let.* 24.12 (Dec. 2017), pp. 1872–1876 (cit. on pp. 23, 28).
- [92] U. S. Kamilov, C. B. Bouman, G. T. Buzzard, and B. Wohlberg. “Plug-and-Play Methods for Integrating Physical and Learned Models in Computational Imaging”. In: *IEEE Signal Process. Mag.* 40.1 (Jan. 2023), pp. 85–97 (cit. on p. 16).
- [93] T. Karras, M. Aittala, T. Aila, and S. Laine. *Elucidating the Design Space of Diffusion-Based Generative Models*. 2022. arXiv: 2206.00364 [cs.CV] (cit. on pp. 12, 38, 39, 48, 74, 92, 105, 110, 129, 135, 136, 167).
- [94] T. Karras, S. Laine, and T. Aila. “A Style-Based Generator Architecture for Generative Adversarial Networks”. In: *IEEE Transactions on Pattern Analysis & Machine Intelligence* 43.12 (Dec. 2021), pp. 4217–4228. ISSN: 1939-3539 (cit. on pp. 94, 112).
- [95] B. Kawar, M. Elad, S. Ermon, and J. Song. *Denoising Diffusion Restoration Models*. 2022. arXiv: 2201.11793 [eess.IV] (cit. on pp. 11, 12, 40, 69, 75, 101, 111, 113, 176).
- [96] B. Kawar, G. Vaksman, and M. Elad. *SNIPS: Solving Noisy Inverse Problems Stochastically*. 2021. arXiv: 2105.14951 [eess.IV] (cit. on p. 82).
- [97] D. Kim and J. A. Fessler. “Optimized first-order methods for smooth convex minimization”. In: *Math Program* 159.1 (2016). DOI: 10.1007/s10107-015-0949-3 (cit. on p. 9).

- [98] D. Kim, S. Shin, K. Song, W. Kang, and I.-C. Moon. “Soft Truncation: A Universal Training Technique of Score-based Diffusion Model for High Precision Score Estimation”. In: *International Conference on Machine Learning*. 2022 (cit. on p. 106).
- [99] D. P. Kingma and J. Ba. “Adam: A Method for Stochastic Optimization”. In: *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*. 2015 (cit. on p. 24).
- [100] F. Knoll, J. Zbontar, A. Sriram, M. J. Muckley, M. Bruno, A. Defazio, M. Parente, K. J. Geras, J. Katsnelson, H. Chandarana, Z. Zhang, M. Drozdal, A. Romero, M. Rabbat, P. Vincent, J. Pinkerton, D. Wang, N. Yakubova, E. Owens, C. L. Zitnick, and Y. W. Lui. “fastMRI: A Publicly Available Raw k-Space and DICOM Dataset of Knee Images for Accelerated MR Image Reconstruction Using Machine Learning”. In: *Radiology: Artificial Intelligence* 2.1 (2020), e190007. DOI: 10.1148/ryai.2020190007 (cit. on p. 112).
- [101] A. B. Konovalov. “Compressed-sensing-inspired reconstruction algorithms in low-dose computed tomography: A review”. In: *Physica Medica* 124 (2024), p. 104491. ISSN: 1120-1797. DOI: <https://doi.org/10.1016/j.ejmp.2024.104491> (cit. on p. 68).
- [102] T. Kwon and J. C. Ye. *Solving Video Inverse Problems Using Image Diffusion Models*. 2025. arXiv: 2409.02574 [cs.CV] (cit. on p. 121).
- [103] T. Kwon and J. C. Ye. *VISION-XL: High Definition Video Inverse Problem Solver using Latent Image Diffusion Models*. 2025. arXiv: 2412.00156 [cs.CV] (cit. on p. 121).
- [104] A. Lahiri, G. Maliakal, M. L. Klasky, J. A. Fessler, and S. Ravishankar. “Sparse-view cone beam CT reconstruction using data-consistent supervised and adversarial learning from scarce training data”. In: *IEEE Transactions on Computational Imaging* 9 (2023), pp. 13–28 (cit. on pp. 8, 10, 68).
- [105] K. Lange, D. R. Hunter, and I. Yang. “Optimization transfer using surrogate objective functions”. In: *J. Computational and Graphical Stat.* 9.1 (Mar. 2000), 1–20 (cit. on p. 169).
- [106] T. Lатычевская. “Iterative phase retrieval in coherent diffractive imaging: practical issues”. In: *Appl. Opt.* 57.25 (2018), pp. 7187–7197. DOI: 10.1364/AO.57.007187 (cit. on p. 8).
- [107] H. Lawrence, D. Barmherzig, H. Li, M. Eickenberg, and M. Gabriele. “Phase Retrieval with Holography and Untrained Priors: Tackling the Challenges of Low-Photon Nanoscale Imaging”. In: *Proceedings of the 2nd Mathematical and Scientific Machine Learning Conference*. Vol. 145. Aug. 2022, pp. 516–567 (cit. on pp. 18, 23).

- [108] S. Lee, H. Chung, J. Kim, and J. C. Ye. *Progressive Deblurring of Diffusion Models for Coarse-to-Fine Image Synthesis*. 2022. arXiv: 2207.11192 [cs.CV] (cit. on p. 17).
- [109] S. Lee, H. Chung, M. Park, J. Park, W.-S. Ryu, and J. C. Ye. *Improving 3D Imaging with Pre-Trained Perpendicular 2D Diffusion Models*. 2023. arXiv: 2303.08440 [eess.IV] (cit. on pp. 13, 40, 47, 55, 56, 60, 63, 70, 142, 150, 151).
- [110] J. Lehtinen, J. Munkberg, J. Hasselgren, S. Laine, T. Karras, M. Aittala, and T. Aila. “Noise2Noise: Learning Image Restoration without Clean Data”. In: *Proceedings of the 35th International Conference on Machine Learning*. PMLR, 2018, pp. 2971–2980 (cit. on p. 16).
- [111] H. Li and Z. Lin. “Accelerated Proximal Gradient Methods for Nonconvex Programming”. In: *Advances in Neural Information Processing Systems*. Vol. 28. 2015 (cit. on pp. 20–22, 31).
- [112] S. Li, D. Zeng, J. Peng, Z. Bian, H. Zhang, Q. Xie, Y. Wang, Y. Liao, S. Zhang, J. Huang, D. Meng, Z. Xu, and J. Ma. “An Efficient Iterative Cerebral Perfusion CT Reconstruction via Low-Rank Tensor Decomposition With Spatial–Temporal Total Variation Regularization”. In: *IEEE Transactions on Medical Imaging* 38.2 (2019), pp. 360–370. DOI: 10.1109/TMI.2018.2865198 (cit. on p. 14).
- [113] X. Li, Z. Zhang, X. Li, S. Chen, Z. Zhu, P. Wang, and Q. Qu. *Understanding Representation Dynamics of Diffusion Models via Low-Dimensional Modeling*. 2025. arXiv: 2502.05743 [cs.LG] (cit. on p. 107).
- [114] Z. Li, J. A. Fessler, J. K. Mikell, S. J. Wilderman, and Y. K. Dewaraja. “DblurDoseNet: A deep residual learning network for voxel radionuclide dosimetry compensating for single-Photon Emission Computerized Tomography Imaging resolution”. In: *Medical Physics* 49.2 (2022), pp. 1216–1230. DOI: 10.1002/mp.15397 (cit. on p. 22).
- [115] Z. Li, J. Hu, X. Xu, L. Shen, and J. A. Fessler. “Accelerated Wirtinger Flow With Score-Based Image Priors for Holographic Phase Retrieval in Poisson-Gaussian Noise Conditions”. In: *IEEE Transactions on Computational Imaging* 10 (2024), pp. 1384–1399. DOI: 10.1109/TCI.2024.3458418 (cit. on pp. 3, 4, 13, 16, 69, 82, 101).
- [116] Z. Li, K. Lange, and J. A. Fessler. “Poisson Phase Retrieval With Wirtinger Flow”. In: *2021 IEEE International Conference on Image Processing (ICIP)*. 2021, pp. 2828–2832. DOI: 10.1109/ICIP42928.2021.9506139 (cit. on p. 9).
- [117] Z. Li, K. Lange, and J. A. Fessler. “Poisson Phase Retrieval in Very Low-Count Regimes”. In: *IEEE Transactions on Computational Imaging* 8 (2022), pp. 838–850. DOI: 10.1109/TCI.2022.3209936 (cit. on p. 9).

- [118] Z. Li, X. Xu, J. Hu, J. Fessler, and Y. Dewaraja. “Reducing SPECT acquisition time by predicting missing projections with single-scan self-supervised coordinate-based learning”. In: *Journal of Nuclear Medicine* 64.supplement 1 (2023), P1014–P1014. eprint: <https://jnm.snmjournals.org/content> (cit. on pp. 7, 9, 13, 82).
- [119] J. Liang, P. Stoica, Y. Jing, and J. Li. “Phase Retrieval via the Alternating Direction Method of Multipliers”. In: *IEEE Sig. Proc. Letters* 25.1 (2018), pp. 5–9. DOI: 10.1109/LSP.2017.2767826 (cit. on pp. 9, 31).
- [120] G.-H. Liu, A. Vahdat, D.-A. Huang, E. A. Theodorou, W. Nie, and A. Anandkumar. *P²SB: Image-to-Image Schrödinger Bridge*. 2023. arXiv: 2302.05872 [cs.CV] (cit. on pp. 11, 119, 120).
- [121] J. Liu, Y. Sun, W. Gan, X. Xu, B. Wohlberg, and U. S. Kamilov. “Stochastic Deep Unfolding for Imaging Inverse Problems”. In: *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 2021, pp. 1395–1399. DOI: 10.1109/ICASSP39728.2021.9414332 (cit. on p. 8).
- [122] J. Liu, Y. Sun, X. Xu, and U. S. Kamilov. *Image Restoration using Total Variation Regularized Deep Image Prior*. 2018. arXiv: 1810.12864 [cs.CV] (cit. on pp. 7, 9, 14).
- [123] J. Liu, X. Xu, W. Gan, S. Shoushtari, and U. Kamilov. “Online Deep Equilibrium Learning for Regularization by Denoising”. In: *Advances in Neural Information Processing Systems*. Ed. by A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho. 2022 (cit. on pp. 8, 10).
- [124] Y. Liu, K. Zhang, Y. Li, Z. Yan, C. Gao, R. Chen, Z. Yuan, Y. Huang, H. Sun, J. Gao, L. He, and L. Sun. *Sora: A Review on Background, Technology, Limitations, and Opportunities of Large Vision Models*. 2024. arXiv: 2402.17177 [cs.CV] (cit. on p. 120).
- [125] Z. Liu, P. Luo, X. Wang, and X. Tang. “Deep Learning Face Attributes in the Wild”. In: *Proceedings of International Conference on Computer Vision (ICCV)*. Dec. 2015 (cit. on pp. 22, 23, 25–28, 32).
- [126] Z. Liu, P. Luo, X. Wang, and X. Tang. “Deep Learning Face Attributes in the Wild”. In: *Proceedings of International Conference on Computer Vision (ICCV)*. Dec. 2015 (cit. on pp. 48, 91, 110).
- [127] Y. Long, J. A. Fessler, and J. M. Balter. “A 3D forward and back-projection method for X-ray CT using separable footprint”. In: *Proc. Intl. Mtg. on Fully 3D Image Recon. in Rad. and Nuc. Med.* Winner of poster award. 2009, 146–9 (cit. on p. 17).

- [128] C. Lu, Y. Zhou, F. Bao, J. Chen, C. Li, and J. Zhu. *DPM-Solver: A Fast ODE Solver for Diffusion Probabilistic Model Sampling in Around 10 Steps*. 2022. arXiv: 2206.00927 [cs.LG] (cit. on pp. 12, 74, 135).
- [129] H. Lu, G. Yang, N. Fei, Y. Huo, Z. Lu, P. Luo, and M. Ding. *VDT: General-purpose Video Diffusion Transformers via Mask Modeling*. 2023. arXiv: 2305.13311 [cs.CV] (cit. on p. 120).
- [130] T. Luhman and E. Luhman. *Improving diffusion model efficiency through patching*. 2022 (cit. on p. 12).
- [131] W. Luo, W. Alghamdi, and Y. M. Lu. “Optimal Spectral Initialization for Signal Recovery With Applications to Phase Retrieval”. In: *IEEE Trans. Sig. Proc.* 67.9 (2019), pp. 2347–2356. DOI: 10.1109/TSP.2019.2904918 (cit. on p. 23).
- [132] M. Lustig, D. Donoho, and J. M. Pauly. “Sparse MRI: The application of compressed sensing for rapid MR imaging”. In: *Magnetic Resonance in Medicine* 58.6 (2007), pp. 1182–1195 (cit. on p. 10).
- [133] M. Lustig, D. L. Donoho, J. M. Santos, and J. M. Pauly. “Compressed sensing MRI”. In: *IEEE Signal Processing Magazine* 25.2 (2008), pp. 72–82 (cit. on p. 10).
- [134] X. Ma, Y. Wang, X. Chen, G. Jia, Z. Liu, Y.-F. Li, C. Chen, and Y. Qiao. *Latte: Latent Diffusion Transformer for Video Generation*. 2025. arXiv: 2401.03048 [cs.CV] (cit. on p. 121).
- [135] C. H. McCollough, A. Bartley, R. Carter, B. Chen, T. Drees, P. Edwards, D. Holmes, A. Huang, F. Khan, S. Leng, K. McMillan, G. Michalak, K. Nunez, L. Yu, and J. Fletcher. “Results of the 2016 Low Dose CT Grand Challenge”. English (US). In: *Medical physics* 44.10 (Oct. 2017), e339–e352. ISSN: 0094-2405. DOI: 10.1002/mp.12345 (cit. on pp. 48, 61, 76, 91, 94, 110, 112).
- [136] C. A. Metzler, A. Maleki, and R. G. Baraniuk. “BM3D-PRGAMP: Compressive phase retrieval based on BM3D denoising”. In: *2016 IEEE International Conference on Image Processing (ICIP)*. 2016, pp. 2504–2508. DOI: 10.1109/ICIP.2016.7532810 (cit. on p. 7).
- [137] R. P. Millane. “Phase retrieval in crystallography and optics”. In: *J. Opt. Soc. Am. A* 7.3 (Mar. 1990), 394–411. DOI: 10.1364/JOSAA.7.000394 (cit. on p. 8).
- [138] T. Moon, M. Choi, G. Lee, J.-W. Ha, and J. Lee. “Fine-tuning Diffusion Models with Limited Data”. In: *NeurIPS 2022 Workshop on Score-Based Methods*. 2022 (cit. on pp. 15, 34, 88).

- [139] Y. Nesterov. “Smooth Minimization of Non-Smooth Functions”. In: *Math. Program.* 103.1 (May 2005), pp. 127–152. ISSN: 0025-5610. DOI: 10 . 1007 / s10107 - 004 - 0552 - 5 (cit. on p. 9).
- [140] A. Nichol and P. Dhariwal. “Improved Denoising Diffusion Probabilistic Models”. In: *International Conference on Machine Learning*. 2021 (cit. on pp. 68, 105).
- [141] G. Ongie, A. Jalal, C. A. Metzler, R. Baraniuk, A. G. Dimakis, and R. M. Willett. “Deep Learning Techniques for Inverse Problems in Imaging”. In: *IEEE Journal on Selected Areas in Information Theory* 1 (2020), pp. 39–56 (cit. on p. 16).
- [142] O. Ozdenizci and R. Legenstein. “Restoring Vision in Adverse Weather Conditions With Patch-Based Denoising Diffusion Models”. In: *IEEE Transactions on Pattern Analysis & Machine Intelligence* 45.08 (Aug. 2023), pp. 10346–10357. ISSN: 1939-3539. DOI: 10 . 1109 / TPAMI . 2023 . 3238179 (cit. on pp. 12, 36, 41, 46, 50, 51, 137, 138).
- [143] C. Y. Park, M. T. McCann, C. Garcia-Cardona, B. Wohlberg, and U. S. Kamilov. “Random walks with Tweedie: A unified view of score-based diffusion models”. In: *IEEE Sig. Proc. Mag.* 42.3 (2025), 40–51. DOI: 10 . 1109 / MSP . 2025 . 3590608 (cit. on p. 101).
- [144] W. Peebles and S. Xie. “Scalable Diffusion Models with Transformers”. In: *arXiv.org* (2023). DOI: 10 . 48550 / arxiv . 2212 . 09748 (cit. on pp. 12, 13, 68).
- [145] K. P. Pruessmann, M. Weiger, M. B. Scheidegger, and P. Boesiger. “SENSE: sensitivity encoding for fast MRI”. In: *Magnetic Resonance in Medicine* 42.5 (1999), pp. 952–962 (cit. on p. 10).
- [146] T. Qiu, P. Babu, and D. P. Palomar. “PRIME: phase retrieval via majorization-minimization”. In: *IEEE Trans. Sig. Proc.* 64.19 (Oct. 2016), 5174–86. DOI: 10 . 1109 / TSP . 2016 . 2585084 (cit. on pp. 9, 31).
- [147] A. W. Reed, H. Kim, R. Anirudh, K. A. Mohan, K. Champley, J. Kang, and S. Jayasuriya. *Dynamic CT Reconstruction from Limited Views with Implicit Neural Representations and Parametric Motion Fields*. 2021. arXiv: 2104 . 11745 [eess . IV] (cit. on pp. 13, 82).
- [148] W. Ren, X. Cao, J. Pan, X. Guo, W. Zuo, and M.-H. Yang. “Image Deblurring via Enhanced Low-Rank Prior”. In: *IEEE Transactions on Image Processing* 25.7 (2016), pp. 3426–3437. DOI: 10 . 1109 / TIP . 2016 . 2571062 (cit. on p. 100).
- [149] Y. Romano, M. Elad, and P. Milanfar. “The Little Engine That Could: Regularization by Denoising (RED)”. In: *SIAM Journal on Imaging Sciences* 10.4 (2017), pp. 1804–1844. DOI: 10 . 1137 / 16M1102884 (cit. on pp. 16, 23, 28).

- [150] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. *High-Resolution Image Synthesis with Latent Diffusion Models*. 2022. arXiv: 2112.10752 [cs.CV] (cit. on pp. 12, 68).
- [151] O. Ronneberger, P. Fischer, and T. Brox. “U-Net: Convolutional Networks for Biomedical Image Segmentation”. In: *ArXiv abs/1505.04597* (2015) (cit. on p. 50).
- [152] J. L. Rumberger, X. Yu, P. Hirsch, M. Dohmen, V. E. Guarino, A. Mokarian, L. Mais, J. Funke, and D. Kainmueller. *How Shift Equivariance Impacts Metric Learning for Instance Segmentation*. 2021. arXiv: 2101.05846 [cs.CV] (cit. on pp. 50, 51, 137, 138).
- [153] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. *ImageNet Large Scale Visual Recognition Challenge*. 2015. arXiv: 1409.0575 [cs.CV] (cit. on pp. 34, 55, 69).
- [154] E. K. Ryu, J. Liu, S. Wang, X. Chen, Z. Wang, and W. Yin. *Plug-and-Play Methods Provably Converge with Properly Trained Denoisers*. 2019. arXiv: 1905.05406 [cs.CV] (cit. on p. 14).
- [155] M. Saliba, T. Lатыchevskaia, J. Longchamp, and H. Fink. “Fourier Transform Holography: A Lensless Non-Destructive Imaging Technique”. In: *Microsc. Microanal.* 18.S2 (2012) (cit. on p. 8).
- [156] A. H. Al-Shabli, X. Xu, I. Selesnick, and U. S. Kamilov. *Bregman Plug-and-Play Priors*. 2022. arXiv: 2202.02388 [eess.IV] (cit. on p. 8).
- [157] Y. Shechtman, Y. C. Eldar, O. Cohen, H. N. Chapman, J. Miao, and M. Segev. “Phase Retrieval with Application to Optical Imaging: A contemporary overview”. In: *IEEE Sig. Proc. Mag.* 32.3 (2015), pp. 87–109. DOI: 10.1109/MSP.2014.2352673 (cit. on pp. 8, 9, 31).
- [158] S. Shoushtari, J. Liu, and U. S. Kamilov. “Diffusion models for phase retrieval in computational imaging”. In: *Proc., IEEE Asilomar Conf. on Signals, Systems, and Comp.* 2023, 779–83. DOI: 10.1109/IEEECONF59524.2023.10477083 (cit. on pp. 17, 22, 23, 28, 29).
- [159] E. Y. Sidky, C.-M. Kao, and X. Pan. “Accurate image reconstruction from few-views and limited-angle data in divergent-beam CT”. In: *Journal of X-Ray Science and Technology* 14.2 (2006), pp. 119–139 (cit. on p. 9).
- [160] I. Skorokhodov, W. Menapace, A. Siarohin, and S. Tulyakov. *Hierarchical Patch Diffusion Models for High-Resolution Video Generation*. 2024. arXiv: 2406.07792 [cs.CV] (cit. on p. 121).

- [161] M. Soltanolkotabi. “Structured Signal Recovery From Quadratic Measurements: Breaking Sample Complexity Barriers via Nonconvex Optimization”. In: *IEEE Trans. Info. Theory* 65.4 (2019), pp. 2374–2400. DOI: 10.1109/TIT.2019.2891653 (cit. on p. 9).
- [162] B. Song, J. Hu, Z. Luo, J. A. Fessler, and L. Shen. *DiffusionBlend: Learning 3D Image Prior through Position-aware Diffusion Score Blending for 3D Computed Tomography Reconstruction*. 2024. arXiv: 2406.10211 [cs.CV] (cit. on pp. 3, 4, 55, 68, 70, 75, 78, 82, 89).
- [163] B. Song, S. M. Kwon, Z. Zhang, X. Hu, Q. Qu, and L. Shen. *Solving Inverse Problems with Latent Diffusion Models via Hard Data Consistency*. 2023. arXiv: 2307.08123 [cs.CV] (cit. on pp. 12, 101).
- [164] B. Song, Z. Zhang, Z. Luo, J. Hu, W. Yuan, J. Jia, Z. Tang, G. Wang, and L. Shen. *CCS: Controllable and Constrained Sampling with Diffusion Models via Initial Noise Perturbation*. 2025. arXiv: 2502.04670 [cs.LG] (cit. on p. 118).
- [165] J. Song, C. Meng, and S. Ermon. *Denoising Diffusion Implicit Models*. 2022. arXiv: 2010.02502 [cs.LG] (cit. on pp. 11, 12, 39, 61, 74, 75, 135, 136, 152, 176).
- [166] Y. Song and S. Ermon. “Generative Modeling by Estimating Gradients of the Data Distribution”. In: *Advances in Neural Information Processing Systems*. Vol. 32. 2019 (cit. on pp. 11, 16, 39, 44, 50, 51, 87, 100, 114, 115, 119, 124, 139).
- [167] Y. Song, L. Shen, L. Xing, and S. Ermon. “Solving Inverse Problems in Medical Imaging with Score-Based Generative Models”. In: *International Conference on Learning Representations*. 2022 (cit. on pp. 13, 17, 82).
- [168] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. “Score-Based Generative Modeling through Stochastic Differential Equations”. In: *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. 2021 (cit. on pp. 10, 13, 16, 22, 34, 55, 68, 82, 101, 124, 135).
- [169] M. Sonogashira, M. Shonai, and M. Iiyama. “High-Resolution Bathymetry by Deep-Learning-Based Image Superresolution”. In: *PloS One* 15.7 (2020), e0235487–e0235487. DOI: 10.1371/journal.pone.0235487 (cit. on pp. 8, 10).
- [170] S. Sreehari, S. V. Venkatakrishnan, B. Wohlberg, G. T. Buzzard, L. F. Drummy, J. P. Simmons, and C. A. Bouman. “Plug-and-Play Priors for Bright Field Electron Tomography and Sparse Interpolation”. In: *IEEE Transactions on Computational Imaging* 2.4 (Dec. 2016), pp. 408–423 (cit. on p. 14).
- [171] J. Su, B. Xu, and H. Yin. “A survey of deep learning approaches to image restoration”. In: *Neurocomputing* 487 (2022), pp. 46–65. ISSN: 0925-2312. DOI: <https://doi.org/10.1016/j.neucom.2022.02.046> (cit. on p. 1).

- [172] Y. Sun, Z. Wu, X. Xu, B. E. Wohlberg, and U. Kamilov. “Scalable Plug-and-Play ADMM with Convergence Guarantees”. In: *IEEE Transactions on Computational Imaging* 7 (July 2021). ISSN: 2573-0436. DOI: 10.1109/TCI.2021.3094062 (cit. on p. 8).
- [173] O.D. Team. *ODL: Operator Discretization Library*. https://odlgroup.github.io/odl/guide/geometry_guide.html. Accessed: April 2024. 2022 (cit. on pp. 49, 92, 113).
- [174] P. Teng, X. Jiang, L. Cai, E. Lee, R. Zhang, J. Zhou, and J. W. Stayman. “3D diffusion posterior sampling for CT reconstruction”. In: *Proc. SPIE 13405 Medical Imaging: Phys. Med. Im.* 2025, p. 134050X. DOI: 10.1117/12.3047466 (cit. on p. 70).
- [175] J.-B. Thibault, K. Sauer, C. Bouman, and J. Hsieh. “A three-dimensional statistical approach to improved image quality for multi-slice helical CT”. In: *Med. Phys.* 34.11 (Nov. 2007), pp. 4526–4544 (cit. on p. 9).
- [176] X. Tian. “Fourier ptychographic reconstruction using mixed Gaussian-Poisson likelihood with total variation regularisation”. In: *Electronics Letters* 55.19 (Sept. 2019), 1041–3. DOI: 10.1049/e1.2019.1141 (cit. on p. 8).
- [177] D. Ulyanov, A. Vedaldi, and V. Lempitsky. “Deep Image Prior”. In: *International Journal of Computer Vision* 128.7 (Mar. 2020), pp. 1867–1888. ISSN: 1573-1405. DOI: 10.1007/s11263-020-01303-4 (cit. on p. 14).
- [178] S. V. Venkatakrisnan, C. A. Bouman, and B. Wohlberg. “Plug-and-Play Priors for Model Based Reconstruction”. In: *Proc. IEEE Global Conf. Signal Process. and Inf. Process.* Austin, TX, USA, Dec. 2013, pp. 945–948 (cit. on pp. 23, 28).
- [179] P. Vincent. “A Connection Between Score Matching and Denoising Autoencoders”. In: *Neural Computation* 23.7 (2011), pp. 1661–1674. DOI: 10.1162/NECO_a_00142 (cit. on p. 11).
- [180] P. Wang, H. Zhang, Z. Zhang, S. Chen, Y. Ma, and Q. Qu. *Diffusion Models Learn Low-Dimensional Distributions via Subspace Clustering*. 2024. arXiv: 2409.02426 [cs.LG] (cit. on p. 107).
- [181] Y. Wang, J. Yu, and J. Zhang. *Zero-Shot Image Restoration Using Denoising Diffusion Null-Space Model*. 2022. arXiv: 2212.00490 [cs.CV] (cit. on pp. 11, 12, 40, 50, 51, 69, 75, 82, 87, 101, 109, 113–115, 124, 176).
- [182] Z. Wang, J. Liu, G. Li, and H. Han. “Blind2Unblind: Self-Supervised Image Denoising with Visible Blind Spots”. In: *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2022, pp. 2017–2026. DOI: 10.1109/CVPR52688.2022.00207 (cit. on p. 16).

- [183] Z. Wang, Y. Jiang, H. Zheng, P. Wang, P. He, Z. Wang, W. Chen, and M. Zhou. *Patch Diffusion: Faster and More Data-Efficient Training of Diffusion Models*. 2023. arXiv: 2304.12526 [cs.CV] (cit. on pp. 12, 38, 39, 49, 83, 88, 105).
- [184] Z. Wang, A. Bovik, H. Sheikh, and E. Simoncelli. “Image quality assessment: from error visibility to structural similarity”. In: *IEEE Transactions on Image Processing* 13.4 (2004), pp. 600–612. DOI: 10.1109/TIP.2003.819861 (cit. on p. 25).
- [185] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli. “Image quality assessment: From error visibility to structural similarity”. In: *IEEE Transactions on Image Processing* 13.4 (Apr. 2004), pp. 600–612 (cit. on pp. 48, 91).
- [186] X. Wei, H. van Gorp, L. Gonzalez-Carabarin, D. Freedman, Y. C. Eldar, and R. J. G. van Sloun. “Deep Unfolding With Normalizing Flow Priors for Inverse Problems”. In: *IEEE Transactions on Signal Processing* 70 (2022), pp. 2962–2971. DOI: 10.1109/TSP.2022.3179807 (cit. on p. 16).
- [187] E. Whang, D. McAllister, A. Reddy, A. Kohli, and L. Waller. “SeidelNet: An Aberration-Informed Deep Learning Model for Spatially Varying Deblurring”. In: *SPIE*. Vol. 12438. 2023, 124380Y–124380Y–6. DOI: 10.1117/12.2650416 (cit. on p. 8).
- [188] Z. Wu, Y. Sun, J. Liu, and U. Kamilov. “Online Regularization by Denoising with Applications to Phase Retrieval”. In: *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*. 2019, pp. 3887–3895. DOI: 10.1109/ICCVW.2019.00482 (cit. on pp. 13, 82).
- [189] W. Xia, C. Niu, W. Cong, and G. Wang. *Sub-volume-based denoising diffusion probabilistic model for cone-beam CT reconstruction from incomplete data*. 2023 (cit. on p. 12).
- [190] W. Xia, W. Cong, and G. Wang. *Patch-Based Denoising Diffusion Probabilistic Model for Sparse-View CT Reconstruction*. 2022. arXiv: 2211.10388 [eess.IV] (cit. on p. 36).
- [191] W. Xia, H. W. Tseng, C. Niu, W. Cong, X. Zhang, S. Liu, R. Ning, S. Vedantham, and G. Wang. *Parallel Diffusion Model-based Sparse-view Cone-beam Breast CT*. 2024. arXiv: 2303.12861 [eess.IV] (cit. on p. 36).
- [192] J. Xu and B. M. W. Tsui. “Interior and sparse-view image reconstruction using a mixed region and voxel-based ML-EM algorithm”. In: *IEEE Trans. Nuc. Sci.* 59.5 (Oct. 2012), pp. 1997–2007 (cit. on p. 9).
- [193] R. Xu, M. Soltanolkotabi, J. P. Haldar, W. Unglaub, J. Zusman, A. F. J. Levi, and R. M. Leahy. *Accelerated Wirtinger Flow: A fast algorithm for ptychography*. <http://arxiv.org/abs/1806.05546>. 2018. arXiv: 1806.05546 [eess.IV] (cit. on pp. 8, 9).

- [194] X. Xu, W. Gan, S. V. V. N. Kothapalli, D. A. Yablonskiy, and U. S. Kamilov. *CoRRECT: A Deep Unfolding Framework for Motion-Corrected Quantitative R2* Mapping*. 2022. arXiv: 2210.06330 [eess.IV] (cit. on pp. 7, 9).
- [195] X. Xu, M. L. Klasky, M. T. McCann, J. Hu, and J. A. Fessler. “Swap-Net: A Memory-Efficient 2.5D Network for Sparse-View 3D Cone Beam CT Reconstruction to ICF Applications”. In: *IEEE Transactions on Computational Imaging* 11 (2025), pp. 872–887. DOI: 10.1109/TCI.2025.3572699 (cit. on p. 10).
- [196] X. Xu, J. Liu, Y. Sun, B. Wohlberg, and U. S. Kamilov. “Boosting the Performance of Plug-and-Play Priors via Denoiser Scaling”. In: *54th Asilomar Conf. on Signals, Systems, and Computers*. 2020, pp. 1305–1312. DOI: 10.1109/IEEECONF51394.2020.9443410 (cit. on pp. 7, 24, 50, 51, 92, 93, 111, 114).
- [197] X. Xu, Y. Sun, J. Liu, B. Wohlberg, and U. S. Kamilov. “Provable Convergence of Plug-and-Play Priors With MMSE Denoisers”. In: *IEEE Signal Processing Letters* 27 (2020), pp. 1280–1284. ISSN: 1558-2361. DOI: 10.1109/lsp.2020.3006390 (cit. on pp. 8, 10).
- [198] B. Yaman, S. A. H. Hosseini, S. Moeller, J. Ellermann, K. Ugurbil, and M. Akcakaya. “Self-supervised learning of physics-based reconstruction neural networks without fully-sampled reference data”. In: *Mag. Res. Med.* 84.6 (Dec. 2020), 3172–91. DOI: 10.1002/mrm.28378 (cit. on p. 98).
- [199] T. Yang, J. Hu, J. A. Fessler, and L. Shen. *Local Patches Meet Global Context: Scalable 3D Diffusion Priors for Computed Tomography Reconstruction*. 2025. arXiv: 2512.18161 [cs.CV] (cit. on pp. 3, 68).
- [200] F. Yu, A. Seff, Y. Zhang, S. Song, T. Funkhouser, and J. Xiao. *LSUN: Construction of a Large-scale Image Dataset using Deep Learning with Humans in the Loop*. 2016. arXiv: 1506.03365 [cs.CV] (cit. on pp. 34, 55, 69, 112).
- [201] Z. Yuan and H. Wang. “Phase retrieval via reweighted Wirtinger flow”. In: *Appl. Opt.* 56.9 (Mar. 2017), pp. 2418–2427. DOI: 10.1364/AO.56.002418 (cit. on p. 30).
- [202] B. Zhang, Z. Wu, B. T. Feng, Y. Song, Y. Yue, and K. L. Bouman. *STeP: A Framework for Solving Scientific Video Inverse Problems with Spatiotemporal Diffusion Priors*. 2025. arXiv: 2504.07549 [cs.CV] (cit. on p. 121).
- [203] H. Zhang, J. Zhou, Y. Lu, M. Guo, P. Wang, L. Shen, and Q. Qu. *The Emergence of Reproducibility and Consistency in Diffusion Models*. 2024. arXiv: 2310.05264 [cs.LG] (cit. on p. 69).
- [204] K. Zhang, Y. Li, W. Zuo, L. Zhang, L. Van Gool, and R. Timofte. “Plug-and-Play Image Restoration With Deep Denoiser Prior”. In: *IEEE Transactions on Pattern*

- Analysis and Machine Intelligence* 44.10 (2022), pp. 6360–6376. DOI: 10 . 1109 / TPAMI . 2021 . 3088914 (cit. on pp. 16, 111, 113).
- [205] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang. “Beyond a Gaussian Denoiser: Residual Learning of Deep CNN for Image Denoising”. In: *IEEE Transactions on Image Processing* 26.7 (2017), pp. 3142–3155. DOI: 10 . 1109 / TIP . 2017 . 2662206 (cit. on p. 23).
- [206] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang. *The Unreasonable Effectiveness of Deep Features as a Perceptual Metric*. 2018. arXiv: 1801 . 03924 [cs . CV] (cit. on p. 175).
- [207] X. Zhang, S. Wen, L. Han, F. Juefei-Xu, A. Srivastava, J. Huang, H. Wang, M. Tao, and D. N. Metaxas. *Spectrum-Aware Parameter Efficient Fine-Tuning for Diffusion Models*. 2024. arXiv: 2405 . 21050 [cs . CV] (cit. on pp. 15, 88).
- [208] Y. Zhang, P. Song, and Q. Dai. “Fourier ptychographic microscopy using a generalized Anscombe transform approximation of the mixed Poisson-Gaussian likelihood”. In: *Optics Express* 25.1 (Jan. 2017), 168–79. DOI: 10 . 1364 / OE . 25 . 000168 (cit. on pp. 8, 176).
- [209] J. Zhu, H. Ma, J. Chen, and J. Yuan. *DomainStudio: Fine-Tuning Diffusion Models for Domain-Driven Image Generation using Limited Data*. 2024. arXiv: 2306 . 14153 [cs . CV] (cit. on p. 15).
- [210] Y. Zhu, K. Zhang, J. Liang, J. Cao, B. Wen, R. Timofte, and L. V. Gool. *Denoising Diffusion Models for Plug-and-Play Image Restoration*. 2023. arXiv: 2305 . 08995 [cs . CV] (cit. on pp. 111, 113, 176).
- [211] Z. Zhuang, D. Yang, D. Barmherzig, and J. Sun. “Phase Retrieval Using Double Deep Image Priors”. In: *NeurIPS 2023 Workshop on Deep Learning and Inverse Problems*. 2023 (cit. on p. 32).