

VARIATIONAL DIFFUSION POSTERIOR SAMPLING WITH MIDPOINT GUIDANCE

Anonymous authors

Paper under double-blind review

ABSTRACT

Diffusion models have recently shown considerable potential in solving Bayesian inverse problems when used as priors. However, sampling from the resulting denoising posterior distributions remains a challenge as it involves intractable terms. To tackle this issue, state-of-the-art approaches formulate the problem as that of sampling from a surrogate diffusion model targeting the posterior and decompose its scores into two terms: the prior score and an intractable guidance term. While the former is replaced by the pre-trained score of the considered diffusion model, the guidance term has to be estimated. In this paper, we propose a novel approach that utilises a decomposition of the transitions which, in contrast to previous methods, allows a trade-off between the complexity of the intractable guidance term and that of the prior transitions. We validate the proposed approach through extensive experiments on linear and nonlinear inverse problems, including challenging cases with latent diffusion models as priors, and demonstrate its effectiveness in reconstructing electrocardiogram (ECG) from partial measurements for accurate cardiac diagnosis.

1 INTRODUCTION

Inverse problems aim to reconstruct signals from incomplete and noisy observations and are prevalent across various fields. In signal and image processing, common examples include signal deconvolution, image restoration, and tomographic image reconstruction (Stuart, 2010; Idier, 2013). Other applications extend to protein backbone motif scaffolding (Watson et al., 2023) and urban mobility modeling (Jiang et al., 2023). Due to the presence of noise and the inherent complexity of the measurement process, these problems are typically ill-posed, meaning they have an infinite number of possible solutions. While many of these solutions may fit the observed data, only a few align with the true underlying signal. Bayesian inverse problems provide a principled framework for addressing the challenges of signal reconstruction by incorporating prior knowledge, enabling more plausible and meaningful solutions. The *a priori* knowledge about the signal \mathbf{x} to be reconstructed is captured in the prior distribution $q(\mathbf{x})$, while the information about the observation \mathbf{y} is encoded in the likelihood function $p(\mathbf{y}|\mathbf{x})$. Given these components, solving the inverse problem boils down to sampling from the posterior distribution $\pi(\mathbf{x}) \propto p(\mathbf{y}|\mathbf{x})q(\mathbf{x})$, which integrates both prior knowledge and observational data.

The choice of the prior distribution is crucial for achieving accurate reconstructions. If the goal is to reconstruct high-resolution data, the prior must be able to model data with similar fine detail and complexity. Recently, denoising diffusion models (DDM) (Sohl-Dickstein et al., 2015; Song & Ermon, 2019; Ho et al., 2020; Song et al., 2021b) have emerged as a powerful approach in this context. These models can generate highly realistic and detailed reconstructions and are becoming increasingly popular as priors. Using a DDM, a sample \mathbf{X}_0 being approximately distributed according to a given data distribution q of interest can be generated by iteratively denoising an initial sample \mathbf{X}_n from a standard Gaussian distribution. The denoising process $(\mathbf{X}_k)_{k=n}^0$ consists of the iterative refinement of \mathbf{X}_k for $k \in \llbracket 1, n \rrbracket$ based on a parametric approximation of the score $\nabla \log q_k$, where $q_k(\mathbf{x}_k) = \int q(\mathbf{x}_0)q_{k|0}(\mathbf{x}_k|\mathbf{x}_0), d\mathbf{x}_0$, with $q_{k|0}(\cdot|\mathbf{x}_0)$ being a Gaussian noising transition.

Denoising Diffusion Models (DDMs) form a highly powerful class of prior distributions, but their use introduces significant challenges in posterior sampling. Specifically, since a DDM prior does not admit an explicit and tractable density, conventional Markov chain Monte Carlo (MCMC) methods,

054 such as the Metropolis–Hastings algorithm and its variants (Besag, 1994; Neal, 2011), cannot be
 055 applied in general. Furthermore, gradient-based MCMC methods often prove inefficient, as they tend
 056 to get trapped in local modes of the posterior.

057 The problem of sampling from the posterior in Bayesian inverse problems with DDM priors has
 058 recently been addressed in several papers (Kadkhodaie & Simoncelli, 2020; Song et al., 2021b; Kawar
 059 et al., 2021; 2022; Lugmayr et al., 2022; Ho et al., 2022; Chung et al., 2023; Song et al., 2023a). These
 060 methods typically modify the denoising process to account for the observation \mathbf{y} . One principled
 061 approach for accurate posterior sampling involves skewing the denoising process at each stage k
 062 using the score $\nabla \log \pi_k$, where π_k is defined analogously to q_k , with q replaced by the posterior
 063 π . This guides the denoising process in a manner that is consistent with the observation. Notably,
 064 the posterior score $\nabla \log \pi_k$ decomposes into two components: the prior score and an additional
 065 intractable term commonly referred to as *guidance*. Previous works leverage *pre-trained* DDMs for
 066 the prior score, while providing various *training-free* approximations for the guidance term (Ho et al.,
 067 2022; Chung et al., 2023; Song et al., 2023a; Finzi et al., 2023). This framework enables solving
 068 a wide range of inverse problems for a given prior. However, despite the notable success of these
 069 methods, the approximations often introduce significant errors, leading to inaccurate reconstructions
 070 when the posterior distribution is highly multimodal, the measurement process is strongly nonlinear,
 071 or the data is heavily contaminated by noise.

072 **Our contribution.** We begin with the observation that the posterior denoising step, which trans-
 073 forms a sample \mathbf{X}_{k+1} into a sample \mathbf{X}_k , does not necessarily require the conditional scores
 074 $\nabla \log \pi_{k+1}$ or $\nabla \log \pi_k$. Instead, we demonstrate that this step can be decomposed into two in-
 075 termediate phases: first, denoising \mathbf{X}_{k+1} into an intermediate *midpoint* state \mathbf{X}_{ℓ_k} , where $\ell_k < k$,
 076 and then noising back to obtain \mathbf{X}_k , unconditionally on the observation \mathbf{y} . This decomposition
 077 introduces an additional degree of freedom, as it only requires the estimation of the guidance term at
 078 the intermediate step ℓ_k rather than step $k + 1$. Building on this insight, we introduce MIDPOINT
 079 GUIDANCE POSTERIOR SAMPLING (MGPS), a novel diffusion posterior sampling scheme that ex-
 080 plicitly leverages this approach. Our algorithm develops a principled approximation of the denoising
 081 transition by utilizing a Gaussian variational approximation, combined with the guidance approxima-
 082 tion proposed by Chung et al. (2023) at the intermediate steps ℓ_k . The strong empirical performance
 083 of our method is demonstrated through an extensive set of experiments. In particular, we validate
 084 our approach on a Gaussian mixture toy example, as well as various linear and nonlinear image
 085 reconstruction tasks—inpainting, super-resolution, phase retrieval, deblurring, JPEG dequantization,
 086 high-dynamic range—using both pixel-space and latent diffusion models (LDM). We also apply our
 087 posterior sampling method to imputation tasks for 12-lead electrocardiograms (ECGs). Our results
 088 demonstrate that *training-free* posterior sampling performs on par with diffusion models specifically
 089 trained for ECG imputation (Alcaraz & Strodthoff, 2022). Moreover, our method demonstrates
 090 significantly improved performance compared to existing approaches.

091 2 POSTERIOR SAMPLING WITH DDM PRIOR

092 **Problem setup.** In this paper, we focus on the approximate sampling from a density of the form

$$093 \pi(\mathbf{x}) := p(\mathbf{y}|\mathbf{x})q(\mathbf{x})/\mathcal{Z}, \quad (2.1)$$

094 where $p(\mathbf{y}|\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$ is a *non-negative* and *differentiable* likelihood that *can be evaluated*
 095 *pointwise*, q is a prior distribution, and $\mathcal{Z} := \int p(\mathbf{y}|\mathbf{x})q(\mathbf{x}) d\mathbf{x}$ is the normalizing constant. We are
 096 interested in the case where a DDM p_0^θ for the prior has been pre-trained for the prior based on i.i.d.
 097 observations from q . As a by-product, we also have parametric approximations $(s_k^\theta)_{k=1}^n$ of the *scores*
 098 $(\nabla \log q_k)_{k=1}^n$, where the marginals $(q_k)_{k=1}^n$ are defined as $q_k(\mathbf{x}_k) := \int q(\mathbf{x}_0) q_{k|0}(\mathbf{x}_k|\mathbf{x}_0) d\mathbf{x}_0$ with
 099 $q_{k|0}(\mathbf{x}_k|\mathbf{x}_0) := \mathcal{N}(\mathbf{x}_k; \sqrt{\alpha_k}\mathbf{x}_0, v_k\mathbf{I}_d)$ and $v_k := 1 - \alpha_k$. Typically, the sequence $(\alpha_k)_{k=0}^n$ is a
 100 decreasing sequence with $\alpha_0 = 1$ and α_n approximately equals zero for n large enough.

101 These score approximations enable the generation of new samples from q according to one of
 102 the DDM sampling schemes (Song & Ermon, 2019; Ho et al., 2020; Song et al., 2021a;b; Karras
 103 et al., 2022). All these approaches boil down to simulating a Markov chain $(\mathbf{X}_k)_{k=n}^0$ backwards
 104 in time starting from a standard Gaussian $p_n := \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$ and following the Gaussian transitions
 105 $p_{k|k+1}^\theta(\mathbf{x}_k|\mathbf{x}_{k+1}) = \mathcal{N}(\mathbf{x}_k; \mathbf{m}_{k|k+1}^\theta(\mathbf{x}_{k+1}), v_{k|k+1}\mathbf{I}_d)$, where the mean functions $\mathbf{m}_{k|k+1}^\theta$ can be
 106 derived from the learned DDM, and $v_{k|k+1} > 0$ are fixed and pre-defined variances. The backward
 107

transitions are designed in such a way that the marginal law of \mathbf{X}_k approximates q_k ; see Appendix A.1 for more details on the form of the backward transitions.

In our problem setup, we assume we only have access to the approximated scores of the prior and no observation from neither posterior π nor the prior q . This setup encompasses of solving an inverse problem without training a conditional model from scratch based on a paired signal/observation dataset; a requirement typically imposed by conditional DDM frameworks (Song et al., 2021b; Batzolis et al., 2021; Tashiro et al., 2021; Saharia et al., 2022). On the other hand, we require access to $p(\mathbf{y}|\cdot)$. This setup includes many applications in Bayesian inverse problems, *e.g.*, image restoration, motif scaffolding in protein design (Trippe et al., 2023; Wu et al., 2023), and trajectory control in traffic-simulation frameworks (Jiang et al., 2023).

Conditional score. Since we have access to a DDM model for q , a natural approach to sampling from π is to define a DDM approximation based on the pre-trained score approximations for q_k . Following the DDM approach, the basic idea here is to sample sequentially from the marginals $\pi_k(\mathbf{x}_k) := \int \pi(\mathbf{x}_0)q_{k|0}(\mathbf{x}_k|\mathbf{x}_0) d\mathbf{x}_0$ by relying on approximations of the conditional scores $(\nabla \log \pi_k)_{k=1}^n$. The latter can be expressed in terms of the unconditional scores by using

$$\pi_k(\mathbf{x}_k) \propto \int p(\mathbf{y}|\mathbf{x}_0)q(\mathbf{x}_0)q_{k|0}(\mathbf{x}_k|\mathbf{x}_0) d\mathbf{x}_0, \quad (2.2)$$

where after defining the backward transition kernel $q_{0|k}(\mathbf{x}_0|\mathbf{x}_k) := q(\mathbf{x}_0)q_{k|0}(\mathbf{x}_k|\mathbf{x}_0)/q_k(\mathbf{x}_k)$ yields

$$\pi_k(\mathbf{x}_k) \propto p_k(\mathbf{y}|\mathbf{x}_k)q_k(\mathbf{x}_k), \quad \text{where } p_k(\mathbf{y}|\mathbf{x}_k) := \int p(\mathbf{y}|\mathbf{x}_0)q_{0|k}(\mathbf{x}_0|\mathbf{x}_k) d\mathbf{x}_0. \quad (2.3)$$

It then follows that

$$\nabla \log \pi_k(\mathbf{x}_k) = \nabla \log q_k(\mathbf{x}_k) + \nabla \log p_k(\mathbf{y}|\mathbf{x}_k).$$

While we have a parametric approximation of the first term on the r.h.s., the bottleneck is the second term, which is intractable due to the integration of $p(\mathbf{y}|\cdot)$ against the conditional distribution $q_{0|k}(\cdot|\mathbf{x}_k)$.

Existing approaches. To circumvent this computational bottleneck, previous works have proposed approximations that involve replacing $q_{0|k}(\cdot|\mathbf{x}_k)$ in (2.3) with either a Dirac delta mass (Chung et al., 2023) or a Gaussian approximation (Ho et al., 2022; Song et al., 2023a)

$$q_{0|k}(\mathbf{x}_0|\mathbf{x}_k) \approx \mathcal{N}(\mathbf{x}_0; \mathbf{m}_{0|k}^\theta(\mathbf{x}_k), v_{0|k}\mathbf{I}_d), \quad (2.4)$$

where $\mathbf{m}_{0|k}^\theta$ is a parametric approximation of $\mathbf{m}_{0|k}(\mathbf{x}_k) := \int \mathbf{x}_0 q_{0|k}(\mathbf{x}_0|\mathbf{x}_k) d\mathbf{x}_0$ and $v_{0|k}$ is a tuning parameter. Using Tweedie’s formula (Robbins, 1956), it holds that $\mathbf{m}_{0|k}(\mathbf{x}_k) = (-\mathbf{x}_k + \sqrt{\alpha_k} \nabla \log q_k(\mathbf{x}_k))/v_k$, which suggests the approximation $\mathbf{m}_{0|k}^\theta(\mathbf{x}_k) := (-\mathbf{x}_k + \sqrt{\alpha_k} \mathbf{s}_k^\theta(\mathbf{x}_k))/v_k$. As for the variance parameter $v_{0|k}$, Ho et al. (2022) suggest setting $v_{0|k} = v_k/\alpha_k$, while Song et al. (2023a) use $v_{0|k} = v_k$. When $p(\mathbf{y}|\cdot)$ is the likelihood of a linear model, *i.e.*, $p(\mathbf{y}|\mathbf{x}) := \mathcal{N}(\mathbf{y}; \mathbf{A}\mathbf{x}, \sigma_y^2 \mathbf{I}_{d_y})$, for $\mathbf{A} \in \mathbb{R}^{d_y \times d}$ and $\sigma_y > 0$, it can be exactly integrated against the Gaussian approximation (2.4), yielding $p_k(\mathbf{y}|\mathbf{x}_k) \approx \mathcal{N}(\mathbf{y}; \mathbf{A}\mathbf{m}_{0|k}^\theta(\mathbf{x}_k), \sigma_y^2 \mathbf{I}_{d_y} + v_{0|k} \mathbf{A}\mathbf{A}^\top)$. Chung et al. (2023), on the other hand, use the pointwise approximation, yielding $p_k(\mathbf{y}|\mathbf{x}_k) \approx p(\mathbf{y}|\mathbf{m}_{0|k}^\theta(\mathbf{x}_k))$. We denote by $\tilde{p}_k(\mathbf{y}|\cdot)$ the resulting approximation stemming from any of these methods. Approximate samples from the posterior distribution are then drawn by simulating a backward Markov chain $(\mathbf{X}_k)_{k=n}^0$, where $\mathbf{X}_n \sim p_n$ and then, given \mathbf{X}_{k+1} , \mathbf{X}_k is obtained via the update

$$\mathbf{X}_k := \tilde{\mathbf{X}}_k + w_{k+1}(\mathbf{X}_{k+1}) \nabla \log \tilde{p}_{k+1}(\mathbf{y}|\mathbf{X}_{k+1}), \quad \text{where } \tilde{\mathbf{X}}_k \sim p_{k|k+1}^\theta(\cdot|\mathbf{X}_{k+1}), \quad (2.5)$$

and w_{k+1} is a method-dependent weighting function. For instance, Chung et al. (2023) use $w_{k+1}(\mathbf{X}_{k+1}) = \zeta \sigma_y^2 / \|\mathbf{y} - \mathbf{A}\mathbf{m}_{0|k+1}^\theta(\mathbf{X}_{k+1})\|$, with $\zeta \in [0, 1]$. We note that the update (2.5) in general involves a vector-Jacobian product, *e.g.*, considering the approximation suggested by Chung et al. (2023), $\nabla_{\mathbf{x}_k} \log \tilde{p}_k(\mathbf{y}|\mathbf{x}_k) = \nabla_{\mathbf{x}_k} \mathbf{m}_{0|k}^\theta(\mathbf{x}_k)^\top \nabla_{\mathbf{x}_0} \log p(\mathbf{y}|\mathbf{x}_0)|_{\mathbf{x}_0 = \mathbf{m}_{0|k}^\theta(\mathbf{x}_k)}$, which incurs an additional computational cost compared to an unconditional diffusion step.

3 THE MGPS ALGORITHM

In this section we propose a novel scheme for the approximate inference of π . We start by presenting a midpoint decomposition of the backward transition that allows us to trade adherence to the prior backward dynamics for improved guidance approximation.

We preface our description of the MGPS algorithm with some additional notations. First, consider the DDPM forward process with transitions given by $q_{k|j}(\mathbf{x}_k|\mathbf{x}_j) := \mathcal{N}(\mathbf{x}_k; \sqrt{\alpha_k/\alpha_j}\mathbf{x}_j, (1-\alpha_k/\alpha_j)\mathbf{I}_d)$ for all $(j, k) \in \llbracket 0, n \rrbracket^2$ such that $j < k$. We also define the joint law

$$\pi_{0:n}(\mathbf{x}_{0:n}) := \pi(\mathbf{x}_0) \prod_{k=0}^{n-1} q_{k+1|k}(\mathbf{x}_{k+1}|\mathbf{x}_k) \quad (3.1)$$

of the forward process when initialized with the posterior (2.1) of interest. Note that with these definitions, π_k (given by (2.3)) is the marginal of $\pi_{0:n}(\mathbf{x}_{0:n})$ w.r.t. \mathbf{x}_k . The time-reversed decomposition of (3.1) writes

$$\pi_{0:n}(\mathbf{x}_{0:n}) := \pi_n(\mathbf{x}_n) \prod_{k=0}^{n-1} \pi_{k|k+1}(\mathbf{x}_k|\mathbf{x}_{k+1}), \quad (3.2)$$

where, more generally, for all $(i, j) \in \llbracket 0, n \rrbracket^2$ such that $i < j$,

$$\pi_{i|j}(\mathbf{x}_i|\mathbf{x}_j) := \pi_i(\mathbf{x}_i)q_{j|i}(\mathbf{x}_j|\mathbf{x}_i)/\pi_j(\mathbf{x}_j) \propto p_i(\mathbf{y}|\mathbf{x}_i)q_{i|j}(\mathbf{x}_i|\mathbf{x}_j), \quad (3.3)$$

where the marginals are given by (2.2), is the conditional density of \mathbf{X}_i given $\mathbf{X}_j = \mathbf{x}_j$ and \mathbf{y} . Using the decomposition (3.1), an exact draw \mathbf{X}_0 from (2.1) is obtained by first drawing \mathbf{X}_n from π_n and then simulating recursively \mathbf{X}_k from the transitions $\pi_{k|k+1}(\cdot|\mathbf{X}_{k+1})$ for all k . However, according to (3.3), each such transition involves the intractable likelihood function $p_k(\mathbf{y}|\mathbf{x}_k)$ which is hard to approximate accurately when k is large.

Midpoint decomposition. The transition (3.3) has two components: the conditional density $p_i(\mathbf{y}|\mathbf{x}_i)$ and the prior transition density $q_{i|j}(\mathbf{x}_i|\mathbf{x}_j)$. The approximation $p_i(\mathbf{y}|\mathbf{x}_i) \approx p(\mathbf{y}|\mathbf{m}_{0|i}^\theta(\mathbf{x}_i))$ of Chung et al. (2023) is accurate when i is small, whereas the Gaussian approximation of $q_{i|j}(\mathbf{x}_i|\mathbf{x}_j)$ proposed by Ho et al. (2020) is accurate when the difference $|i - j|$ is small. Thus, in order to combine the best of both worlds, we introduce a decomposition that transfers the problem of sampling from $\pi_{k|k+1}$ to that of sampling from $\pi_{\ell_k|k+1}$, i.e., drawing \mathbf{X}_{ℓ_k} given \mathbf{X}_{k+1} and \mathbf{y} , where $\ell_k < k$. The parameter ℓ_k balances the approximation errors of the two conditional densities in (3.3); see Figure 1. In order to make this construction, define, for each $\ell \in \llbracket 0, k \rrbracket$, the bridge kernel $q_{k|\ell, k+1}(\mathbf{x}_k|\mathbf{x}_\ell, \mathbf{x}_{k+1}) \propto q_{k|\ell}(\mathbf{x}_k|\mathbf{x}_\ell)q_{k+1|k}(\mathbf{x}_{k+1}|\mathbf{x}_k)$, which is Gaussian; see Appendix A.1. By convention, $q_{k|k, k+1}(\cdot|\mathbf{x}_k, \mathbf{x}_{k+1}) = \delta_{\mathbf{x}_k}$.

Lemma 3.1. For all $k \in \llbracket 1, n - 1 \rrbracket$ and $\ell_k \in \llbracket 0, k \rrbracket$ it holds that

$$\pi_{k|k+1}(\mathbf{x}_k|\mathbf{x}_{k+1}) = \int q_{k|\ell_k, k+1}(\mathbf{x}_k|\mathbf{x}_{\ell_k}, \mathbf{x}_{k+1})\pi_{\ell_k|k+1}(\mathbf{x}_{\ell_k}|\mathbf{x}_{k+1})d\mathbf{x}_{\ell_k}. \quad (3.4)$$

The proof is found in Appendix A.3. Lead by the decomposition provided by Lemma 3.1, we define

$$\hat{\pi}_{k|k+1}^\ell(\mathbf{x}_k|\mathbf{x}_{k+1}) := \int q_{k|\ell_k, k+1}(\mathbf{x}_k|\mathbf{x}_{\ell_k}, \mathbf{x}_{k+1})\hat{\pi}_{\ell_k|k+1}^\theta(\mathbf{x}_{\ell_k}|\mathbf{x}_{k+1})d\mathbf{x}_{\ell_k}, \quad (3.5)$$

where $\hat{\pi}_{\ell_k|k+1}^\theta(\mathbf{x}_{\ell_k}|\mathbf{x}_{k+1}) \propto \hat{p}_{\ell_k}^\theta(\mathbf{y}|\mathbf{x}_{\ell_k})p_{\ell_k|k+1}^\theta(\mathbf{x}_{\ell_k}|\mathbf{x}_{k+1})$ with $\hat{p}_{\ell_k}^\theta(\mathbf{y}|\mathbf{x}_{\ell_k}) := p(\mathbf{y}|\mathbf{m}_{0|\ell_k}^\theta(\mathbf{x}_{\ell_k}))$ and $p_{\ell_k|k+1}^\theta(\mathbf{x}_{\ell_k}|\mathbf{x}_{k+1}) := q_{\ell_k|0, k+1}(\mathbf{x}_{\ell_k}|\mathbf{m}_{0|k+1}^\theta(\mathbf{x}_{k+1}), \mathbf{x}_{k+1})$. Finally, for any sequence $\ell := (\ell_k)_{k=1}^n$ satisfying $\ell_k \leq k$, we define our surrogate model for $\pi_{0:n}$ and the posterior (2.1) as

$$\hat{\pi}_{0:n}^\ell(\mathbf{x}_{0:n}) := \hat{\pi}_n(\mathbf{x}_n) \prod_{k=0}^{n-1} \hat{\pi}_{k|k+1}^\ell(\mathbf{x}_k|\mathbf{x}_{k+1}), \quad \hat{\pi}_0^\ell(\mathbf{x}_0) := \int \hat{\pi}_{0:n}^\ell(\mathbf{x}_{0:n})d\mathbf{x}_{1:n}, \quad (3.6)$$

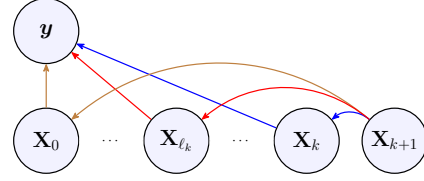


Figure 1: For each color, the different solid arrows indicate different conditional densities that need to be approximated for a given choice of the midpoint ℓ_k . The longer the arrow, the more difficult it is to approximate the corresponding conditional density. By placing the ℓ_k midway between zero and k , the shortest arrows are obtained.

where $\hat{\pi}_n(\mathbf{x}_n) = \mathcal{N}(\mathbf{x}_n; \mathbf{0}_d, \mathbf{I}_d)$ and $\hat{\pi}_{0|1}^\ell(\cdot|\mathbf{x}_1) := \delta_{\mathbf{m}_{0|1}^\theta(\mathbf{x}_1)}$ (similarly to Ho et al., 2020; Song et al., 2021a). Since the bridge kernel $q_{k|\ell_k, k+1}$ is Gaussian, it can be easily sampled. Therefore, we only need to focus on the approximate sampling from $\pi_{\ell_k|k+1}$. Moreover, since the decomposition (3.4) is valid for any ℓ_k , it can be chosen to better balance the approximation errors, as discussed above.

Choice of the sequence $(\ell_k)_{k=1}^{n-1}$. The accuracy of the surrogate model (3.6) ultimately depends on the design of the sequence $(\ell_k)_{k=1}^{n-1}$. In fact, for a given k , note that a decrease in ℓ_k ensures that the approximation $\hat{p}_{\ell_k}^\theta(\mathbf{y}|\mathbf{x}_{\ell_k})$ of $p_{\ell_k}(\mathbf{y}|\mathbf{x}_{\ell_k})$ becomes more accurate. For instance, setting $\ell_k = 0$ eliminates any error of the approximate likelihood (since $p_0(\mathbf{y}|\cdot) = p(\mathbf{y}|\cdot)$); however, decreasing ℓ_k can cause the distribution $\pi_{0|k+1}(\cdot|\mathbf{x}_{k+1})$ to become strongly multimodal, especially in the early stages of the diffusion process, making the DDPM Gaussian approximation less accurate. Conversely, setting ℓ_k closer to k , similar to the DPS approach of Chung et al. (2023), improves the accuracy of the approximation of $q_{\ell_k|k+1}$, but impairs the approximation $\hat{p}_{\ell_k}^\theta(\mathbf{y}|\mathbf{x}_{\ell_k})$ of $p_{\ell_k}(\mathbf{y}|\mathbf{x}_{\ell_k})$. The choice of ℓ_k is therefore subject to a particular trade-off, which we illustrate using the following Gaussian toy example. When considering the parameterization $\ell_k^n = \lfloor \eta k \rfloor$ with $\eta \in [0, 1]$, we show that the Wasserstein-2 distance between π and $\hat{\pi}_0^{\ell(\eta)}$ reaches its minimum for η around 0.5, which confirms that the additional flexibility obtained by introducing the midpoint can provide a surrogate model (3.6) with reduced approximation error.

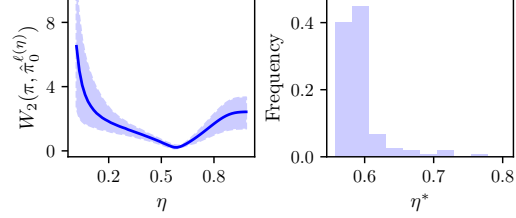


Figure 2: Right: average W_2 with 10%–90% quantile range. Left: distribution of the minimizing η^* .

Example 3.2. Consider a linear inverse problem $\mathbf{Y} = \mathbf{A}\mathbf{X} + \sigma_{\mathbf{y}}\mathbf{Z}$ with Gaussian noise and prior $q = \mathcal{N}(\mathbf{m}, \Sigma)$, where $\mathbf{m} \in \mathbb{R}^d$ and $\Sigma \in \mathbb{R}^{d \times d}$ is positive definite. In this setting, we have access to the exact denoiser $\mathbf{m}_{0|k}$ in a closed form, which allows us to take $\mathbf{m}_{0|k}^\theta = \mathbf{m}_{0|k}$. Moreover, since the transition (3.5) admits an explicit expression in this case, we can quantify the approximation error of (3.6) w.r.t. π for any sequence $(\ell_k)_{k=1}^{n-1}$. Note that in this case the true backward transitions $q_{\ell_k|k+1}$ are Gaussian and have the same mean as the Gaussian DDPM transitions, but differ in their covariance. All necessary derivations and computational details can be found in Appendix B. We consider sequences $\ell(\eta) := (\ell_k^\eta)_{k=1}^n$ with $\ell_k^\eta = \lfloor \eta k \rfloor$, $\eta \in [0, 1]$, and compute the Wasserstein-2 distance $W_2(\pi, \hat{\pi}_0^{\ell(\eta)})$ as a function of η on the basis of 500 samples $(\mathbf{A}, \mathbf{m}, \Sigma)$ for $d = 100$. The left panel of Figure 2 displays the average error $W_2(\pi, \hat{\pi}_0^{\ell(\eta)})$ as a function of η , while the right panel displays the associated distribution of $\eta^* := \operatorname{argmin}_{\eta \in [0, 1]} W_2(\pi, \hat{\pi}_0^{\ell(\eta)})$. Figure 2 illustrates that the smallest approximation error for the surrogate model (3.6) is reached at intermediate values close to $\eta = 0.5$ rather than $\eta = 1$, which corresponds to a DPS-like approximation.

Variational approximation. To sample from $\hat{\pi}_{0,n}^\theta$ in (3.6), we focus on simulating approximately and recursively (over k) the Markov chain with transition densities (3.5). Ideally, for $k = n$, we simply sample from $\hat{\pi}_n^\theta$; then, recursively, assuming that we have access to an approximate sample \mathbf{X}_{k+1} from $\hat{\pi}_{k+1}^\theta$, the next state \mathbf{X}_k is drawn from $\hat{\pi}_{k|k+1}^\theta(\cdot|\mathbf{X}_{k+1})$. However, as $\hat{\pi}_{k|k+1}^\theta(\cdot|\mathbf{X}_{k+1})$ is intractable, we propose using the Gaussian approximation that we specify next.

In the following, let $k \in \llbracket 1, n \rrbracket$ and $\ell_k \in \llbracket 0, k \rrbracket$ be fixed. For $\varphi = (\hat{\boldsymbol{\mu}}_{\ell_k}, \hat{\boldsymbol{\rho}}_{\ell_k}) \in (\mathbb{R}^d)^2$, consider the Gaussian variational distribution

$$\lambda_{k|k+1}^\varphi(\mathbf{x}_k|\mathbf{x}_{k+1}) := \int q_{k|\ell_k, k+1}(\mathbf{x}_k|\mathbf{x}_{\ell_k}, \mathbf{x}_{k+1}) \lambda_{\ell_k|k+1}^\varphi(\mathbf{x}_{\ell_k}) d\mathbf{x}_{\ell_k}, \quad (3.7)$$

$$\lambda_{\ell_k|k+1}^\varphi(\mathbf{x}_{\ell_k}) := \mathcal{N}(\mathbf{x}_{\ell_k}; \hat{\boldsymbol{\mu}}_{\ell_k}, \operatorname{diag}(e^{2\hat{\boldsymbol{\rho}}_{\ell_k}})). \quad (3.8)$$

In definition (3.8) the exponential function is applied element-wise to the vector $\hat{\boldsymbol{\rho}}_{\ell_k}$ and $\operatorname{diag}(e^{\hat{\boldsymbol{\rho}}_{\ell_k}})$ is the diagonal matrix with diagonal entries $e^{\hat{\rho}_{\ell_k}}$. Based on the family $\{\lambda_{k|k+1}^\varphi : \varphi \in (\mathbb{R}^d)^2\}$ and an approximate sample \mathbf{X}_{k+1} from $\hat{\pi}_{k+1}^\theta$, we then seek the best fitting parameter φ that minimizes the upper bound on the backward KL divergence obtained using the data-processing inequality

$$\operatorname{KL}(\lambda_{k|k+1}^\varphi(\cdot|\mathbf{X}_{k+1}) \parallel \hat{\pi}_{k|k+1}^\theta(\cdot|\mathbf{X}_{k+1})) \leq \operatorname{KL}(\lambda_{\ell_k|k+1}^\varphi \parallel \hat{\pi}_{\ell_k|k+1}^\theta(\cdot|\mathbf{X}_{k+1})) =: \mathcal{L}_k(\varphi; \mathbf{X}_{k+1}).$$

Here the gradient of the upper bound is given by

$$\nabla_\varphi \mathcal{L}_k(\varphi; \mathbf{X}_{k+1}) = -\mathbb{E}[\nabla_\varphi \log \hat{p}_{\ell_k}^\theta(\mathbf{y}|\hat{\boldsymbol{\mu}}_{\ell_k} + \operatorname{diag}(e^{\hat{\boldsymbol{\rho}}_{\ell_k}})\mathbf{Z})] + \nabla_\varphi \operatorname{KL}(\lambda_{\ell_k|k+1}^\varphi \parallel \hat{p}_{\ell_k|k+1}^\theta(\cdot|\mathbf{X}_{k+1})),$$

Algorithm 1 MIDPOINT GUIDANCE POSTERIOR SAMPLING

```

270 1: Input:  $(\ell_k)_{k=1}^n$  with  $\ell_n = n$  and  $\ell_1 = 1$ ; number  $M$  of gradient steps.
271
272 2:  $\mathbf{X}_n \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$ ,  $\hat{\mathbf{X}}_n \leftarrow \mathbf{X}_n$ 
273
274 3: for  $k = n - 1$  to 1 do
275 4:  $\hat{\boldsymbol{\mu}}_{\ell_k} \leftarrow \frac{\sqrt{\alpha_{\ell_k}(1-\alpha_{k+1}/\alpha_{\ell_k})}}{1-\alpha_{k+1}} \mathbf{m}_{0|\ell_{k+1}}^\theta(\hat{\mathbf{X}}_{\ell_{k+1}}) + \frac{\sqrt{\alpha_{k+1}/\alpha_{\ell_k}(1-\alpha_{\ell_k})}}{1-\alpha_{k+1}} \mathbf{X}_{k+1}$ 
276 5:  $\hat{\boldsymbol{\rho}}_{\ell_k} \leftarrow \frac{1}{2} \log \frac{(1-\alpha_{k+1}/\alpha_{\ell_k})(1-\alpha_{\ell_k})}{1-\alpha_{k+1}}$ 
277 6: for  $j = 1$  to  $M$  do
278 7:  $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$ 
279 8:  $(\hat{\boldsymbol{\mu}}_{\ell_k}, \hat{\boldsymbol{\rho}}_{\ell_k}) \leftarrow \text{OptimizerStep}(\nabla_{\boldsymbol{\varphi}} \tilde{\mathcal{L}}_k(\cdot, \mathbf{Z}; \mathbf{X}_{k+1}); \hat{\boldsymbol{\mu}}_{\ell_k}, \hat{\boldsymbol{\rho}}_{\ell_k})$ 
280 9: end for
281 10:  $\mathbf{Z}_{\ell_k}, \mathbf{Z}_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$ 
282 11:  $\hat{\mathbf{X}}_{\ell_k} \leftarrow \hat{\boldsymbol{\mu}}_{\ell_k} + \text{diag}(e^{\hat{\boldsymbol{\rho}}_{\ell_k}}) \mathbf{Z}_{\ell_k}$ 
283 12:  $\mathbf{X}_k \sim q_{k|\ell_k, k+1}(\cdot | \hat{\mathbf{X}}_{\ell_k}, \mathbf{X}_{k+1})$  (See (A.5))
284 13: end for
285 14:  $\mathbf{X}_0 \leftarrow \mathbf{m}_{0|1}^\theta(\mathbf{X}_1)$ 

```

where $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$ is independent of \mathbf{X}_{k+1} . The first term is dealt with using the reparameterization trick (Kingma & Welling, 2014) and can be approximated using a Monte Carlo estimator based on a single sample \mathbf{Z}_{k+1} ; we denote this estimate by $\nabla_{\boldsymbol{\varphi}} \tilde{\mathcal{L}}_k(\boldsymbol{\varphi}, \mathbf{Z}_{k+1}; \mathbf{X}_{k+1})$. The second term is the gradient of the KL divergence between two Gaussian distributions and can thus be computed in a closed form. Similarly to the previous approaches of Ho et al. (2022); Chung et al. (2023); Song et al. (2023a), our gradient estimator involves a vector-Jacobian product of the denoising network.

We now provide a summary of the MGPS algorithm, whose pseudocode is given in Algorithm 1. Given $(\ell_k)_{k=1}^{n-1}$, MGPS proceeds by simulating a Markov chain $(\mathbf{X}_k)_{k=n}^0$ starting from $\mathbf{X}_n \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$. Recursively, given the state \mathbf{X}_{k+1} , the state \mathbf{X}_k is obtained by

1. minimizing $\tilde{\mathcal{L}}_k(\cdot; \mathbf{X}_{k+1})$ by performing M stochastic gradient steps using the gradient estimator $\nabla_{\boldsymbol{\varphi}} \tilde{\mathcal{L}}_k(\cdot; \mathbf{X}_{k+1})$, yielding a parameter $\boldsymbol{\varphi}_k^*(\mathbf{X}_{k+1})$,
2. sampling $\hat{\mathbf{X}}_{\ell_k} \sim \lambda_{\ell_k|k+1}^{\boldsymbol{\varphi}_k^*}$, where we drop the dependence on \mathbf{X}_{k+1} in $\boldsymbol{\varphi}_k^*$, and then $\mathbf{X}_k \sim q_{k|\ell_k, k+1}(\cdot | \hat{\mathbf{X}}_{\ell_k}, \mathbf{X}_{k+1})$.

Remark 3.3. While conditionally on $\boldsymbol{\varphi}_k^*$ we have that $\hat{\mathbf{X}}_{\ell_k} \sim \lambda_{\ell_k|k+1}^{\boldsymbol{\varphi}_k^*}$, the actual law of $\hat{\mathbf{X}}_{\ell_k}$ given \mathbf{X}_{k+1} is not a Gaussian distribution as we need to marginalize over that of $\boldsymbol{\varphi}_k^*$ due to the randomness in the stochastic gradients.

A well-chosen initialization for the Gaussian variational approximation parameters is crucial for achieving accurate fitting with few optimization steps, ensuring the overall runtime of MGPS remains competitive with existing methods. Given $\hat{\mathbf{X}}_{\ell_{k+1}}$ sampled from $\lambda_{\ell_{k+1}|k+2}^{\boldsymbol{\varphi}_{k+1}^*}$ at the previous step, we choose the initial parameter $\boldsymbol{\varphi}_k$ so that

$$\lambda_{\ell_k|k+1}^{\boldsymbol{\varphi}_k}(\mathbf{x}_k) = q_{\ell_k|0, k+1}(\mathbf{x}_k | \mathbf{m}_{0|\ell_{k+1}}^\theta(\hat{\mathbf{X}}_{\ell_{k+1}}), \mathbf{X}_{k+1}). \quad (3.9)$$

This initialization is motivated by noting that $q_{\ell_k|0, k+1}(\cdot | \mathbf{m}_{0|\ell_{k+1}}^\pi(\mathbf{X}_{k+1}), \mathbf{X}_{k+1})$, the DDPM approximation of $\pi_{\ell_k|k+1}(\cdot | \mathbf{X}_{k+1})$, where $\mathbf{m}_{0|\ell_{k+1}}^\pi$ is a supposed denoiser for the posterior, is a reasonable candidate for the initialization. As it is intractable, we replace it with $\mathbf{m}_{0|\ell_{k+1}}^\theta(\hat{\mathbf{X}}_{\ell_{k+1}})$, our best current guess. Finally, in Appendix C.2 we devise a warm-start approach that improves the initialization during the early steps of the algorithm. This has proved advantageous in the most challenging problems.

Related works. We now discuss existing works that have similarities with our algorithm, MGPS. In essence, our method tries to reduce the approximation error incurred by DPS-like approximations. This has also been the focus of many works in the literature, which we review below.

When $p(\mathbf{y} | \cdot)$ is a likelihood associated to a linear inverse problem with Gaussian noise, Finzi et al. (2023); Stevens et al. (2023); Boys et al. (2023) leverage the fact that the Gaussian projection

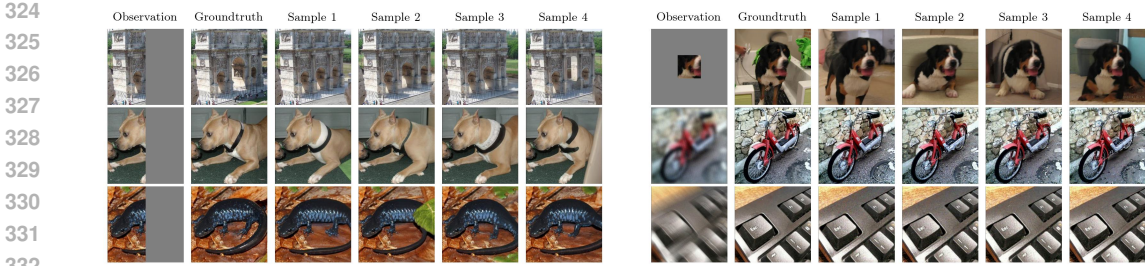


Figure 3: MGPS sample images for half mask (left), expand task, Gaussian blur and motion blur (right) on the ImageNet dataset.

of $q_{0|k}(\cdot|\mathbf{x}_k)$ minimizing the forward KL divergence can be approximated using the pre-trained denoisers; see [Meng et al. \(2021, Theorem 1\)](#). In the considered Gaussian setting, the likelihood $p(\mathbf{y}|\cdot)$ can be exactly integrated against the estimated Gaussian projection. However, this involves a matrix inversion that may be prohibitively expensive. [Boys et al. \(2023\)](#) circumvent the latter by using a diagonal approximation which involves the computation of d vector-Jacobian products. For general likelihoods $p(\mathbf{y}|\cdot)$, [Song et al. \(2023b\)](#) use the Gaussian approximations of [Ho et al. \(2022\)](#); [Song et al. \(2023a\)](#) to estimate $p_k(\mathbf{y}|\cdot)$ for non-linear likelihoods $p(\mathbf{y}|\cdot)$ using a vanilla Monte Carlo estimate. [Zhang et al. \(2023\)](#) also consider the surrogate transitions (3.5) with $\ell_k = k$, and, given an approximate sample \mathbf{X}_{k+1} from $\hat{\pi}_{k+1}^\ell$, estimate the next sample \mathbf{X}_k maximizing $\mathbf{x}_k \mapsto \hat{\pi}_{k|k+1}^\theta(\mathbf{x}_k|\mathbf{X}_{k+1})$ using gradient ascent. They also consider improved estimates of the likelihood $p_k(\mathbf{y}|\cdot)$ running a few diffusion steps. However, it is shown ([Zhang et al., 2023](#), last subtable in Table 2) that this brings minor improvements at the expense of a sharp increase in computational cost, which is mainly due to backpropagation over the denoisers. See also [Appendix C.3](#) for a discussion on how our method relates to the DPS ([Chung et al., 2023](#)) in the case where $\ell_k = k$.

Finally, a second line of work considers the distribution path $(\hat{\pi}_k)_{k=n}^0$, where $\hat{\pi}_k(\mathbf{x}_k) \propto p(\mathbf{y}|\mathbf{m}_{0|k}(\mathbf{x}_k))q_k(\mathbf{x}_k)$ for $k \in \llbracket 0, n-1 \rrbracket$ and $\hat{\pi}_n = \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$. This path bridges the Gaussian distribution and the posterior of interest. Furthermore, if one is able to accurately sample from $\hat{\pi}_{k+1}$, then these samples can be used to initialize a sampler targeting the next distribution $\hat{\pi}_k$. As these distributions are expected to be close, the sampling from $\hat{\pi}_k$ can also be expected to be accurate. Repeating this process yields approximate samples from π regardless of the approximation error in the likelihoods. This approach is pursued by [Rozet & Louppe \(2023\)](#), who combines the update (2.5) with a Langevin dynamics targeting $\hat{\pi}_k$. [Wu et al. \(2023\)](#) use instead sequential Monte Carlo ([Gordon et al., 1993](#)) to recursively build empirical approximations of each $\hat{\pi}_k$ by evolving a sample of N particles.

4 EXPERIMENTS

We now evaluate our algorithm on three different problems and compare it with several competitors. We begin in benchmarking our method on toy Gaussian-mixture targets and image experiments with both pixel space and latent diffusion models. We review the latter in [Appendix A.2](#). Finally, we perform inpainting experiments on ECG data. For experiments based on pixel-space diffusion models, MGPS is benchmarked against state-of-the-art methods in the literature: DPS ([Chung et al., 2023](#)), PGDM ([Song et al., 2023a](#)), DDNM ([Wang et al., 2023](#)), DIFFPIR ([Zhu et al., 2023](#)), and REDDIFF ([Mardani et al., 2024](#)). Regarding experiments using latent diffusion models, we compare against PSLD ([Rout et al., 2024](#)) and RESAMPLE ([Song et al., 2024](#)). Full details on the hyperparameters can be found in [Appendix D.1](#).

Gaussian mixture. We first evaluate the accuracy of our method and the impact of the hyperparameters on a toy linear inverse problem with Gaussian mixture (GM) as q and a Gaussian likelihood $p(\mathbf{y}|\mathbf{x}) := \mathcal{N}(\mathbf{y}; \mathbf{A}\mathbf{x}, \sigma_y^2 \mathbf{I}_d)$, where $\mathbf{A} \in \mathbb{R}^{d_y \times d}$ and $\sigma_y = 0.05$. We repeat the experiment of [Cardoso et al. \(2023, App B.3\)](#), where the prior is a GM with 25 well-separated strata. For this specific example, both the posterior and the denoiser $\mathbf{m}_{0|k}$ are available in a closed form. In particular, the posterior can be shown to be a GM. The full details are provided in [Appendix D.3](#). We consider the settings (20, 1) and (200, 1) for the dimensions (d, d_y) . For each method, we reconstruct 1000

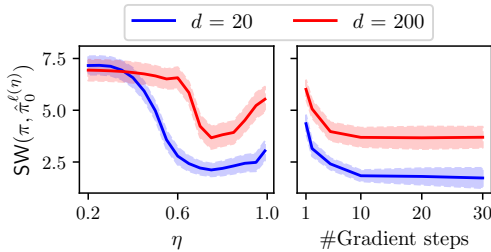


Figure 4: Left: SW as a function of η with $\ell_k = \lfloor \eta k \rfloor$. Right: SW as a function of the number of gradient steps, for a specific choice of $(\ell_k)_k$.

Table 1: 95 % confidence interval for the SW on the GM experiment.

	$d = 20, d_y = 1$	$d = 200, d_y = 1$
MGPS	2.11 ± 0.30	3.66 ± 0.53
DPS	8.93 ± 0.49	9.15 ± 0.44
PGDM	2.44 ± 0.36	5.23 ± 0.38
DDNM	4.24 ± 0.37	7.10 ± 0.50
DIFFPIR	4.14 ± 0.42	8.43 ± 0.92
REDDIFF	6.70 ± 0.45	8.35 ± 0.39

samples and then compute the sliced Wasserstein (SW) distance to exact samples from the posterior distribution. We repeat this procedure 100 times with randomly drawn matrices \mathbf{A} , and consider the resulting averaged SW distance. In every replication, the observation is generated as $\mathbf{Y} = \mathbf{A}\mathbf{X} + \sigma_y \mathbf{Z}$ where $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}_{d_y}, \mathbf{I}_{d_y})$ and $\mathbf{X} \sim q$. The results are reported in Table 1. It can be observed that MGPS with $\eta = 0.75$ outperforms all baselines in both settings. Although we have tuned the parameters of DPS it still exhibits considerable instability and often diverges. To account for these instabilities we set the SW to 10 when it diverges.

Ablations. Next, we perform an ablation study on the total number of gradient steps and the choice of the sequence $(\ell_k)_{k=1}^{n-1}$. The right-hand plot in Figure 4 shows the average SW distance as a function of the number of gradient steps per denoising step when using $\ell_k = \lfloor 3k/4 \rfloor$. For $d = 20$, we observe a monotonic decrease of the SW distance. For $d = 200$, the SW distance decreases and then stabilizes after 10 gradient steps, indicating that in this case, MGPS performs well without requiring many additional gradient steps. We also report the average SW as a function of $\eta \in [0, 1]$ in the left-hand plot of Figure 4, following the same sequence choices as in Example 3.2, *i.e.*, $\ell_k^\eta = \lfloor \eta k \rfloor$. In both settings, the best SW is achieved at $\eta = 0.75$. In dimension $d = 200$, the SW at $\eta = 0.75$ is nearly twice as good as at $\eta = 1$, which corresponds to a DPS-like approximation. This demonstrates that the introduced trade-off leads to non-negligible performance gains.

Images. We evaluate our algorithm on a wide range of linear and nonlinear inverse problems with noise level $\sigma_y = 0.05$. As linear problems we consider: image inpainting with a box mask of shape 150×150 covering the center of the image and half mask covering its right-hand side; image Super Resolution (SR) with factors $\times 4$ and $\times 16$; image deblurring with Gaussian blur; motion blur with 61×61 kernel size. We use the same experimental setup as Chung et al. (2023, Section 4) for the last two tasks. Regarding the nonlinear inverse problems on which we benchmark our algorithm, we consider: JPEG dequantization with quality factor 2%; phase retrieval with oversampling factor 2; non-uniform deblurring emulated via the forward model of Tran et al. (2021); high dynamic range (HDR) following Mardani et al. (2024, Section 5.2). Since the JPEG operator is not differentiable, we instead use the differentiable JPEG framework from Shin & Song (2017). We note that only DPS and REDDIFF apply to nonlinear problems.

Datasets and models. We test our algorithm on the 256×256 versions of the FFHQ (Karras et al., 2019) and ImageNet (Deng et al., 2009) datasets. We chose randomly 50 images from each to be used in the evaluation of the algorithms. We leverage pre-trained DDMs that are publicly available. For FFHQ, we use the pixel-space DDM of Choi et al. (2021) and the LDM of Rombach et al. (2022) with VQ4 autoencoder. We use the model by Dhariwal & Nichol (2021) for ImageNet.

Evaluation. We use LPIPS (Zhang et al., 2018) as the metric to evaluate the quality of the reconstructions. Given that the posterior distribution is expected to be highly multimodal for most of the tasks, pixel-wise comparisons with the ground truth, such as PSNR and SSIM, provide limited insight into reconstruction quality. We nonetheless provide them in the extended tables Table 7, 8 and 9 in the appendix. For all the considered algorithms, we use the hyperparameters proposed in their official implementations if available, otherwise we manually tune them to achieve the best reconstructions. We defer implementation details to Appendix D.1. For each inverse problem, we compute the LPIPS between the ground-truth image and a reconstruction and then average over the dataset. For phase retrieval, however, due to the task’s inherent instability, we follow the approach of Chung et al. (2023)

Table 2: Mean LPIPS for various linear and nonlinear imaging tasks on the FFHQ and ImageNet 256×256 datasets with $\sigma_y = 0.05$.

Task	FFHQ						ImageNet					
	MGPS	DPS	PGDM	DDNM	DIFFPIR	REDDIFF	MGPS	DPS	PGDM	DDNM	DIFFPIR	REDDIFF
SR ($\times 4$)	0.09	0.09	0.33	0.14	<u>0.13</u>	0.36	0.30	0.41	0.78	<u>0.34</u>	0.36	0.56
SR ($\times 16$)	<u>0.26</u>	0.24	0.44	0.30	0.28	0.51	<u>0.53</u>	0.50	0.60	0.70	0.63	0.83
Box inpainting	0.10	0.19	0.17	<u>0.12</u>	0.18	0.19	0.22	0.34	0.29	<u>0.28</u>	<u>0.28</u>	0.36
Half mask	0.20	0.24	0.26	<u>0.22</u>	0.23	0.28	0.29	0.44	0.38	<u>0.38</u>	<u>0.35</u>	0.44
Gaussian Deblur	<u>0.15</u>	0.16	0.87	0.19	0.12	0.23	<u>0.32</u>	0.35	1.00	0.45	0.29	0.54
Motion Deblur	0.13	0.16	—	—	—	0.21	0.22	0.39	—	—	—	0.40
JPEG (QF = 2)	0.16	0.39	1.10	—	—	<u>0.32</u>	0.42	0.63	1.31	—	—	0.51
Phase retrieval	0.11	0.46	—	—	—	<u>0.25</u>	0.47	0.62	—	—	—	<u>0.60</u>
Nonlinear deblur	0.23	<u>0.52</u>	—	—	—	0.66	0.44	0.88	—	—	—	<u>0.67</u>
High dynamic range	0.07	0.49	—	—	—	<u>0.20</u>	0.10	0.85	—	—	—	<u>0.21</u>

by selecting the best reconstruction out of four; further discussion on this can be found in Chung et al. (2023, Appendix C.6).

Results. We report the results in Table 2 for FFHQ and ImageNet with DDM prior, and Table 3 for FFHQ with LDM prior. These reveal that MGPS consistently outperforms the competing algorithms across both linear and nonlinear problems. Notably, on nonlinear tasks, MGPS achieves up to a twofold reduction in LPIPS compared to the other methods with DDM prior. It effectively manages the additional nonlinearities introduced by using LDMs and surpasses the state-of-the-art, as demonstrated in Table 3. Reconstructions obtained with MGPS are displayed in Figure 3 and Figure 5. Further examples and comparisons with other competitors are provided in Appendix D.6. It is seen that MGPS provides high quality reconstructions even for the most challenging examples. For the DDM prior, the results are obtained using the warm-start approach, see Algorithm 3, and setting $\ell_k = \lfloor k/2 \rfloor$ on all the tasks based on the FFHQ dataset. As for the ImageNet dataset, we use the same configuration on all the tasks, except for Gaussian deblur and motion deblur. For these tasks we found that using $\ell_k = \lfloor k/2 \rfloor \mathbb{1}_{k > \lfloor n/2 \rfloor} + k \mathbb{1}_{k < \lfloor n/2 \rfloor}$ improves performance. We use a similar strategy for the LDM prior. In Appendix D.2.1, we measure the runtime and GPU memory requirements for each algorithm. The memory requirement of our algorithm is the same as that of DPS and PGDM, but the runtime of MGPS with $n = 300$ is slightly larger. With fewer diffusion steps, the runtime is halved compared to DPS and PGDM and comparable to that of DIFFPIR, DDNM and REDDIFF. Still, MGPS remains competitive and consistently outperforms the baselines, particularly on nonlinear tasks. See Table 7 and 8 for detailed results using $n \in \{50, 100\}$ diffusion steps.

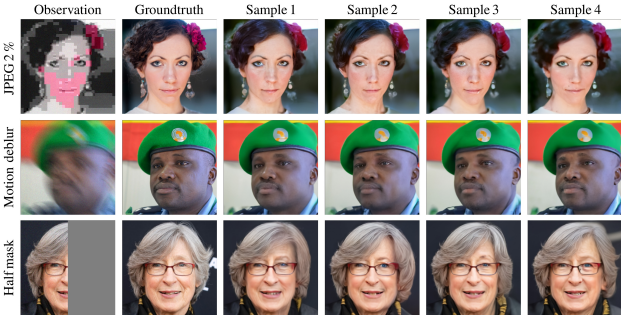


Figure 5: MGPS samples with LDM on FFHQ dataset.

Table 3: Mean LPIPS with LDM on FFHQ.

Task	MGPS	RESAMPLE	PSLD
SR ($\times 4$)	0.11	<u>0.20</u>	0.22
SR ($\times 16$)	0.30	0.36	<u>0.35</u>
Box inpainting	0.16	<u>0.22</u>	0.26
Half mask	0.25	<u>0.30</u>	0.31
Gaussian Deblur	<u>0.16</u>	0.15	0.35
Motion Deblur	0.18	<u>0.19</u>	0.41
JPEG (QF = 2)	0.20	0.26	—
Phase retrieval	0.34	0.41	—
Nonlinear deblur	0.26	0.30	—
High dynamic range	0.15	0.15	—

ECG. An electrocardiogram (ECG) is an electrical recording of the heart’s activity in which the signals generated by the heartbeats are recorded in order to diagnose various cardiac conditions such as cardiac arrhythmias and conduction abnormalities. Unlike static images, ECGs constitute complex time-series data consisting of 12 electrical signals acquired using 10 electrodes, 4 of which are attached to the limbs and record the ‘limb leads’, while the others record the ‘precordial leads’ around the heart. We study two conditional generation problems in ECGs. The first is a forecasting or missing-block (MB) reconstruction problem, where one half of the ECG is reconstructed from the other (see fig. 8). This task evaluates the algorithm’s ability to capture temporal information to predict a coherent signal. The second problem is an inpainting or missing-leads (ML) reconstruction, where the entire ECG is reconstructed from the lead I (see fig. 6). The question is whether we can capture

the subtle information contained in lead I and reconstruct a coherent ECG with the same diagnosis as the real ECG. This task is challenging because lead I, being acquired with limb electrodes far from the heart, may contain very subtle features related to specific cardiac conditions. We train a state-space diffusion model (Goel et al., 2022) to generate ECGs using the 20k training ECGs from the PTB-XL dataset (Wagner et al., 2020), and benchmark the posterior sampling algorithm on the 2k test ECGs (see appendix D.5). We report the Mean Absolute Error (MAE) and the Root Mean Squared Error (RMSE) between ground-truth and reconstructions in Table 4. We demonstrate that a diffusion model trained to generate ECGs can serve as a prior to solve imputation tasks without additional training or fine-tuning, yielding superior results to a diffusion model trained conditionally on the MB task (Alcaraz & Strodthoff, 2022). The rationale for this result is discussed in Appendix D.5.2. We report in Table 5 the balanced accuracy of diagnosing Left Bundle Branch Block (LBBB), Right Bundle Branch Block (RBBB), Atrial Fibrillation (AF), and Sinus Bradycardia (SB) using the downstream classifier proposed in (Strodthoff et al., 2020) (see appendix D.5.1), applied to both ground-truth and to reconstructed samples from lead I. See appendix D.5.2 for sample figures. MGPS outperforms all other posterior sampling algorithms with just 50 diffusion steps.

Table 4: MAE and RMSE for missing block task on the PTB-XL dataset.

Metric	MGPS	DPS	PGDM	DDNM	DIFFPIR	REDDIFF	TRAINEDDIFF
MAE	$0.111 \pm 2e-3$	$0.117 \pm 4e-3$	$0.118 \pm 2e-3$	$0.103 \pm 2e-3$	$0.115 \pm 2e-3$	$0.171 \pm 3e-3$	$0.116 \pm 2e-3$
RMSE	$0.225 \pm 4e-3$	$0.232 \pm 4e-3$	$0.233 \pm 4e-3$	$0.224 \pm 4e-3$	$0.233 \pm 4e-3$	$0.287 \pm 5e-3$	$0.266 \pm 3e-3$

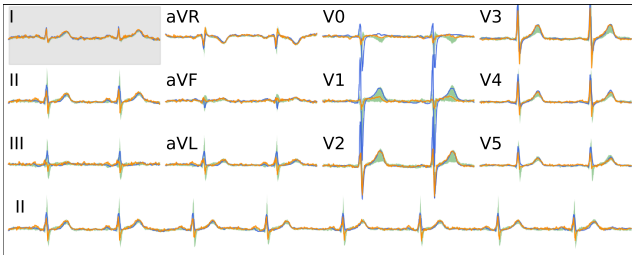


Figure 6: 10s ECG reconstruction from lead I. Ground-truth in blue, 10%-90% quantile range in green, random sample in orange.

Table 5: Balanced acc. downstream diagnosis from ECG reconstructed from lead I.

Method	RBBB	LBBB	AF	SB
MGPS ₅₀	<u>0.81</u>	<u>0.92</u>	0.94	<u>0.66</u>
MGPS ₃₀₀	0.90	0.93	<u>0.92</u>	0.66
DPS	0.54	0.84	0.79	0.50
PGDM	0.65	0.87	0.88	0.55
DDNM	0.71	0.83	0.86	0.59
DIFFPIR	0.57	0.80	0.77	0.53
REDDIFF	0.73	0.86	0.88	0.60
Ground-truth	0.99	0.98	0.94	0.70

5 CONCLUSION

We have introduced MGPS, a novel posterior sampling algorithm designed to solve general inverse problems using both diffusion models and latent diffusion models as a prior. Our approach is based on trading off the approximation error in the prior backward dynamics for a more accurate approximation of the guidance term. This strategy has proven to be effective in a variety of numerical experiments across various tasks. The results show that MGPS consistently performs competitively and often even superior to state-of-the-art methods.

Limitations and future directions. Although the proposed method is promising, it has certain limitations that invite further exploration. First, a detailed analysis of how the algorithm’s performance depends on the choice of the intermediate time steps $(\ell_k)_{k=1}^{n-1}$ remains a challenging but critical subject for future research. Although our image and ECG experiments suggest that setting $\ell_k = \lfloor k/2 \rfloor$ leads to significant performance gains on most tasks, we have observed that using an adaptive sequence $\ell_k = \lfloor \eta_k k \rfloor$, where η_k increases as k decreases, further enhances results on certain tasks, such as Gaussian and motion deblurring, as well as when using latent diffusion models. A second direction for improvement lies in refining the approximation of $p_k(\mathbf{y}|\cdot)$ which we believe could lead to overall algorithmic improvements. Specifically, devising an approximation $\hat{p}_k^\theta(\mathbf{y}|\cdot)$ such that $\nabla \log \hat{p}_k^\theta(\mathbf{y}|\cdot)$ does not require a vector-Jacobian product could significantly reduce the algorithm’s runtime. Lastly, we see potential in using the two-stage approximation introduced in our warm-start strategy (see Algorithm 3) at every diffusion step. Although this technique is promising, it currently leads to instability as k decreases. Understanding and resolving this instability will be a key focus of our future work.

REFERENCES

- 540
541
542 Juan Lopez Alcaraz and Nils Strodthoff. Diffusion-based time series imputation and forecasting
543 with structured state space models. *Transactions on Machine Learning Research*, 2022. ISSN
544 2835-8856. URL <https://openreview.net/forum?id=hHiIbk7Apw>.
545
- 546 Juan Miguel Lopez Alcaraz and Nils Strodthoff. Diffusion-based conditional ecg generation with
547 structured state space models. *Computers in Biology and Medicine*, 163:107115, 2023. ISSN
548 0010-4825. doi: <https://doi.org/10.1016/j.compbio.2023.107115>. URL <https://www.sciencedirect.com/science/article/pii/S0010482523005802>.
549
- 550 Georgios Batzolis, Jan Stanczuk, Carola-Bibiane Schönlieb, and Christian Etmann. Conditional
551 image generation with score-based diffusion models. *arXiv preprint arXiv:2111.13606*, 2021.
552
- 553 J.E. Besag. Comments on “Representations of knowledge in complex systems” by U. Grenander and
554 M. Miller. *J. Roy. Statist. Soc. Ser. B*, 56:591–592, 1994.
- 555 Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and*
556 *Statistics)*. Springer-Verlag, Berlin, Heidelberg, 2006. ISBN 0387310738.
557
- 558 Benjamin Boys, Mark Girolami, Jakiw Pidstrigach, Sebastian Reich, Alan Mosca, and O Deniz Aky-
559 ildiz. Tweedie moment projected diffusions for inverse problems. *arXiv preprint arXiv:2310.06721*,
560 2023.
- 561 Gabriel Cardoso, Yazid Janati El Idrissi, Sylvain Le Corff, and Eric Moulines. Monte Carlo guided
562 diffusion for Bayesian linear inverse problems. *arXiv preprint arXiv:2308.07983*, 2023.
563
- 564 Jooyoung Choi, Sungwon Kim, Yonghyun Jeong, Youngjune Gwon, and Sungroh Yoon. Ilvr: Condi-
565 tioning method for denoising diffusion probabilistic models. in 2021 *ieee*. In *CVF international*
566 *conference on computer vision (ICCV)*, volume 1, pp. 2, 2021.
- 567 Hyungjin Chung, Jeongsol Kim, Michael Thompson Mccann, Marc Louis Klasky, and Jong Chul Ye.
568 Diffusion posterior sampling for general noisy inverse problems. In *The Eleventh International*
569 *Conference on Learning Representations*, 2023. URL [https://openreview.net/forum?](https://openreview.net/forum?id=OnD9zGAGT0k)
570 [id=OnD9zGAGT0k](https://openreview.net/forum?id=OnD9zGAGT0k).
- 571 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
572 hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,
573 pp. 248–255. Ieee, 2009.
574
- 575 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances*
576 *in neural information processing systems*, 34:8780–8794, 2021.
- 577 Marc Anton Finzi, Anudhyan Boral, Andrew Gordon Wilson, Fei Sha, and Leonardo Zepeda-Núñez.
578 User-defined event sampling and uncertainty quantification in diffusion models for physical
579 dynamical systems. In *International Conference on Machine Learning*, pp. 10136–10152. PMLR,
580 2023.
- 581 Karan Goel, Albert Gu, Chris Donahue, and Christopher Ré. It’s raw! audio generation with
582 state-space models. *International Conference on Machine Learning (ICML)*, 2022.
583
- 584 N. J. Gordon, D. J. Salmond, and A. F. M. Smith. Novel approach to nonlinear/ non-Gaussian
585 Bayesian state estimation. *IEE Proceedings-F*, 140(2):107–113, April 1993.
586
- 587 Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured
588 state spaces. In *The International Conference on Learning Representations (ICLR)*, 2022.
- 589 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
590 *Neural Information Processing Systems*, 33:6840–6851, 2020.
591
- 592 Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J
593 Fleet. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633–8646,
2022.

- 594 Jérôme Idier. *Bayesian approach to inverse problems*. John Wiley & Sons, 2013.
- 595
- 596 Chiyu Jiang, Andre Cornman, Cheolho Park, Benjamin Sapp, Yin Zhou, Dragomir Anguelov, et al.
- 597 MotiDiffuser: Controllable multi-agent motion prediction using diffusion. In *Proceedings of the*
- 598 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9644–9653, 2023.
- 599
- 600 Zahra Kadkhodaie and Eero P Simoncelli. Solving linear inverse problems using the prior implicit in
- 601 a denoiser. *arXiv preprint arXiv:2007.13640*, 2020.
- 602
- 603 Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative
- 604 adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*
- 605 *recognition*, pp. 4401–4410, 2019.
- 606
- 607 Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-
- 608 based generative models. *Advances in Neural Information Processing Systems*, 35:26565–26577,
- 609 2022.
- 610
- 611 Bahjat Kawar, Gregory Vaksman, and Michael Elad. Snips: Solving noisy inverse problems stochas-
- 612 tically. *Advances in Neural Information Processing Systems*, 34:21757–21769, 2021.
- 613
- 614 Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising diffusion restoration
- 615 models. *Advances in Neural Information Processing Systems*, 35:23593–23606, 2022.
- 616
- 617 Diederik Kingma and Max Welling. Auto-encoding variational Bayes. In *International Conference*
- 618 *on Learning Representations*, 2014.
- 619
- 620 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *The 3rd*
- 621 *International Conference for Learning Representations, San Diego, 2015*, 2015.
- 622
- 623 Anji Liu, Mathias Niepert, and Guy Van den Broeck. Image inpainting via tractable steering of
- 624 diffusion models. *arXiv preprint arXiv:2401.03349*, 2023.
- 625
- 626 Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool.
- 627 Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the*
- 628 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11461–11471, 2022.
- 629
- 630 Morteza Mardani, Jiaming Song, Jan Kautz, and Arash Vahdat. A variational perspective on solving
- 631 inverse problems with diffusion models. In *The Twelfth International Conference on Learning*
- 632 *Representations*, 2024. URL <https://openreview.net/forum?id=1Y04EE3SPB>.
- 633
- 634 Chenlin Meng, Yang Song, Wenzhe Li, and Stefano Ermon. Estimating high order gradients of
- 635 the data distribution by denoising. *Advances in Neural Information Processing Systems*, 34:
- 636 25359–25369, 2021.
- 637
- 638 R.M. Neal. MCMC using Hamiltonian dynamics. *Handbook of Markov chain Monte Carlo*, 2011.
- 639
- 640 Ingram Olkin and Friedrich Pukelsheim. The distance between two random vectors with given
- 641 dispersion matrices. *Linear Algebra and its Applications*, 48:257–263, 1982.
- 642
- 643 Herbert E. Robbins. An empirical bayes approach to statistics. In *Proceedings of the Third Berkeley*
- 644 *Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory of*
- 645 *Statistics*, 1956. URL <https://api.semanticscholar.org/CorpusID:26161481>.
- 646
- 647 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
- resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Confer-*
- ence on Computer Vision and Pattern Recognition*, pp. 10684–10695, 2022.
- Litu Rout, Negin Raoof, Giannis Daras, Constantine Caramanis, Alex Dimakis, and Sanjay Shakkottai.
- Solving linear inverse problems provably via posterior sampling with latent diffusion models.
- Advances in Neural Information Processing Systems*, 36, 2024.
- François Rozet and Gilles Louppe. Score-based data assimilation. *Advances in Neural Information*
- Processing Systems*, 36:40521–40541, 2023.

- 648 Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi.
649 Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis and*
650 *Machine Intelligence*, 45(4):4713–4726, 2022.
- 651 Richard Shin and Dawn Song. Jpeg-resistant adversarial images. In *NIPS 2017 workshop on machine*
652 *learning and computer security*, volume 1, pp. 8, 2017.
- 654 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised
655 learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*,
656 pp. 2256–2265. PMLR, 2015.
- 657 Bowen Song, Soo Min Kwon, Zecheng Zhang, Xinyu Hu, Qing Qu, and Liyue Shen. Solving inverse
658 problems with latent diffusion models via hard data consistency. In *The Twelfth International*
659 *Conference on Learning Representations*, 2024. URL [https://openreview.net/forum?](https://openreview.net/forum?id=j8hdRqOUhN)
660 [id=j8hdRqOUhN](https://openreview.net/forum?id=j8hdRqOUhN).
- 662 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International*
663 *Conference on Learning Representations*, 2021a. URL [https://openreview.net/](https://openreview.net/forum?id=StlgjarCHLP)
664 [forum?id=StlgjarCHLP](https://openreview.net/forum?id=StlgjarCHLP).
- 665 Jiaming Song, Arash Vahdat, Morteza Mardani, and Jan Kautz. Pseudoinverse-guided diffusion
666 models for inverse problems. In *International Conference on Learning Representations*, 2023a.
667 URL https://openreview.net/forum?id=9_gsMA8MRKQ.
- 669 Jiaming Song, Qinsheng Zhang, Hongxu Yin, Morteza Mardani, Ming-Yu Liu, Jan Kautz, Yongxin
670 Chen, and Arash Vahdat. Loss-guided diffusion models for plug-and-play controllable generation.
671 In *International Conference on Machine Learning*, pp. 32483–32498. PMLR, 2023b.
- 672 Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution.
673 *Advances in neural information processing systems*, 32, 2019.
- 674 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
675 Poole. Score-based generative modeling through stochastic differential equations. In *International*
676 *Conference on Learning Representations*, 2021b.
- 678 Tristan SW Stevens, Hans van Gorp, Faik C Meral, Junseob Shin, Jason Yu, Jean-Luc Robert,
679 and Ruud JG van Sloun. Removing structured noise with diffusion models. *arXiv preprint*
680 *arXiv:2302.05290*, 2023.
- 681 Nils Strodthoff, Patrick Wagner, Tobias Schaeffter, and Wojciech Samek. Deep learning for ecg
682 analysis: Benchmarks and insights from ptb-xl. *arXiv preprint 2004.13701*, 2020.
- 684 Andrew M Stuart. Inverse problems: a Bayesian perspective. *Acta numerica*, 19:451–559, 2010.
- 685 Yusuke Tashiro, Jiaming Song, Yang Song, and Stefano Ermon. Csd: Conditional score-based diffu-
686 sion models for probabilistic time series imputation. *Advances in Neural Information Processing*
687 *Systems*, 34:24804–24816, 2021.
- 689 Phong Tran, Anh Tuan Tran, Quynh Phung, and Minh Hoai. Explore image deblurring via encoded
690 blur kernel space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*
691 *recognition*, pp. 11956–11965, 2021.
- 693 Brian L. Trippe, Jason Yim, Doug Tischer, David Baker, Tamara Broderick, Regina Barzilay, and
694 Tommi S. Jaakkola. Diffusion probabilistic modeling of protein backbones in 3d for the motif-
695 scaffolding problem. In *The Eleventh International Conference on Learning Representations*, 2023.
696 URL <https://openreview.net/forum?id=6TxBxqNME1Y>.
- 697 Tim Van Erven and Peter Harremo. Rényi divergence and kullback-leibler divergence. *IEEE*
698 *Transactions on Information Theory*, 60(7):3797–3820, 2014.
- 699 Patrick Wagner, Nils Strodthoff, Ralf-Dieter Boussejot, Dieter Kreisler, Fatima Lunze, Wojciech
700 Samek, and Tobias Schaeffter. Ptb-xl, a large publicly available electrocardiography dataset.
701 *Scientific Data*, 7:154, 05 2020. doi: 10.1038/s41597-020-0495-6.

702 Yinhuai Wang, Jiwen Yu, and Jian Zhang. Zero-shot image restoration using denoising diffusion
703 null-space model. In *The Eleventh International Conference on Learning Representations, 2023*.
704 URL <https://openreview.net/forum?id=mRieQgMtNTQ>.
705

706 Joseph L Watson, David Juergens, Nathaniel R Bennett, Brian L Trippe, Jason Yim, Helen E Eisenach,
707 Woody Ahern, Andrew J Borst, Robert J Ragotte, Lukas F Milles, et al. De novo design of protein
708 structure and function with rfdiffusion. *Nature*, 620(7976):1089–1100, 2023.

709 Luhuan Wu, Brian L. Trippe, Christian A Naesseth, John Patrick Cunningham, and David Blei.
710 Practical and asymptotically exact conditional sampling in diffusion models. In *Thirty-seventh*
711 *Conference on Neural Information Processing Systems, 2023*. URL <https://openreview.net/forum?id=eWKqrlzCRv>.
712

713 Guanhua Zhang, Jiabao Ji, Yang Zhang, Mo Yu, Tommi Jaakkola, and Shiyu Chang. Towards
714 coherent image inpainting using denoising diffusion implicit models. In *International Conference*
715 *on Machine Learning*, pp. 41164–41193. PMLR, 2023.
716

717 Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable
718 effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on*
719 *computer vision and pattern recognition*, pp. 586–595, 2018.

720 Yuanzhi Zhu, Kai Zhang, Jingyun Liang, Jiezhang Cao, Bihan Wen, Radu Timofte, and Luc Van Gool.
721 Denoising diffusion models for plug-and-play image restoration. In *Proceedings of the IEEE/CVF*
722 *Conference on Computer Vision and Pattern Recognition*, pp. 1219–1229, 2023.
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

A BACKGROUND ON DENOISING DIFFUSION MODELS

A.1 DENOISING DIFFUSION PROBABILISTIC MODELS

In this section, we provide further background on DDMs based on the DDPM framework (Ho et al., 2020; Dhariwal & Nichol, 2021; Song et al., 2021a). We rely on definitions provided in the main text.

DDPMs define generative models for q relying only on parametric approximations $(\mathbf{m}_{0|t}^\theta)_{t=1}^T$ of the denoisers $(\mathbf{m}_{0|t})_{t=1}^T$. These approximate denoisers are usually defined through the parameterization

$$\mathbf{m}_{0|t}^\theta(\mathbf{x}_t) = (\mathbf{x}_t - \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_t^\theta(\mathbf{x}_t)) / \sqrt{\alpha_t} \quad (\text{A.1})$$

and trained by minimizing the denoising loss

$$\sum_{t=1}^T w_t \mathbb{E} [\|\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_t^\theta(\sqrt{\alpha_t} \mathbf{X}_0 + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_t)\|^2], \quad (\text{A.2})$$

w.r.t. the neural network parameter θ , where $(\boldsymbol{\epsilon}_t)_{t=1}^T$ are i.i.d. standard normal vectors, $\mathbf{X}_0 \sim q$, and $(w_t)_{t=1}^T$ are some nonnegative weights. Having trained the denoisers, the generative model for q is defined as follows. Let $(t_k)_{k=0}^n$ be an increasing sequence of time steps in $\llbracket 0, T \rrbracket$ with $t_0 = 0$. We assume that t_n is large enough so that $q_{t_n}(\mathbf{x}_{t_n}) = \int q(\mathbf{x}_0) q_{t_n|0}(\mathbf{x}_{t_n} | \mathbf{x}_0) d\mathbf{x}_0$ is approximately the density of a multivariate standard normal distribution. For convenience we assign the index k to any quantity depending on t_k ; e.g., we denote q_{t_k} by q_k . Consider the backward decomposition

$$q_{0:n}(\mathbf{x}_{0:n}) = q(\mathbf{x}_0) \prod_{k=0}^{n-1} q_{k+1|k}(\mathbf{x}_{k+1} | \mathbf{x}_k) = q_n(\mathbf{x}_n) \prod_{k=0}^{n-1} q_{k|k+1}(\mathbf{x}_k | \mathbf{x}_{k+1})$$

of the forward process initialized at q , where $q_{k|k+1}(\mathbf{x}_k | \mathbf{x}_{k+1}) \propto q_k(\mathbf{x}_k) q_{k+1|k}(\mathbf{x}_{k+1} | \mathbf{x}_k)$. Next, for $(j, \ell, k) \in \llbracket 0, n \rrbracket^3$ such that $j < \ell < k$, define

$$\mathbf{m}_{\ell|j,k}(\mathbf{x}_j, \mathbf{x}_k) := \frac{\sqrt{\alpha_\ell/\alpha_j}(1 - \alpha_k/\alpha_\ell)}{1 - \alpha_k/\alpha_j} \mathbf{x}_j + \frac{\sqrt{\alpha_k/\alpha_\ell}(1 - \alpha_\ell/\alpha_j)}{1 - \alpha_k/\alpha_j} \mathbf{x}_k, \quad (\text{A.3})$$

$$v_{\ell|j,k} := \frac{(1 - \alpha_\ell/\alpha_j)(1 - \alpha_k/\alpha_\ell)}{1 - \alpha_k/\alpha_j}. \quad (\text{A.4})$$

Then the bridge kernel is

$$\begin{aligned} q_{\ell|j,k}(\mathbf{x}_\ell | \mathbf{x}_j, \mathbf{x}_k) &= q_{\ell|j}(\mathbf{x}_\ell | \mathbf{x}_j) q_{k|\ell}(\mathbf{x}_k | \mathbf{x}_\ell) / q_{k|j}(\mathbf{x}_k | \mathbf{x}_j) \\ &= \mathcal{N}(\mathbf{x}_\ell; \mathbf{m}_{\ell|j,k}(\mathbf{x}_j, \mathbf{x}_k), v_{\ell|j,k} \mathbf{I}_d). \end{aligned} \quad (\text{A.5})$$

Using the bridge kernel, $q_{0:n}$ is approximated using the variational approximation

$$p_{0:n}^\theta(\mathbf{x}_{0:n}) = p_n^\theta(\mathbf{x}_n) \prod_{k=0}^{n-1} p_{k|k+1}^\theta(\mathbf{x}_k | \mathbf{x}_{k+1}),$$

where $p_{0|1}^\theta(\mathbf{x}_0 | \mathbf{x}_1) = \mathcal{N}(\mathbf{x}_0; \mathbf{m}_{0|1}^\theta(\mathbf{x}_1), v_{0|1} \mathbf{I}_d)$, $v_{0|1}$ being a tunable parameter, and

$$p_{k|k+1}^\theta(\mathbf{x}_k | \mathbf{x}_{k+1}) := q_{k|0,k+1}(\mathbf{x}_k | \mathbf{m}_{0|k+1}^\theta(\mathbf{x}_{k+1}), \mathbf{x}_{k+1}) \quad k \in \llbracket 1, n-1 \rrbracket. \quad (\text{A.6})$$

When $n = T$, the denoising objective (A.2) corresponds to the KL divergence $\text{KL}(q_{0:n} \| p_{0:n}^\theta)$ for a specific choice of weights $(w_t)_t$. In practice, a DDPM is trained using the objective (A.2) with large T but at inference n is usually much smaller.

A.2 LATENT DIFFUSION MODELS

Latent diffusion models (LDM) (Rombach et al., 2022) define a DDM in a latent space. Let $\mathcal{E} : \mathbb{R}^d \rightarrow \mathbb{R}^p$ be an encoder function and $\mathcal{D} : \mathbb{R}^p \rightarrow \mathbb{R}^d$ a decoder function. We assume that these functions satisfy $\mathcal{D}_\# \mathcal{E}_\# q \approx q$, where, for instance, $\mathcal{E}_\# q$ denotes the law of the random variable $\mathcal{E}(\mathbf{X})$ where $\mathbf{X} \sim q$. A LDM approximating q is given by $\mathcal{D}_\# p_0^\theta$, where p_0^θ is a diffusion model trained on

810 samples from $\mathcal{E}_{\#}q$. Finally, when solving inverse problems with LDMs, we assume that the target
 811 distribution is instead

$$812 \quad \pi(\mathbf{x}) \propto p(\mathbf{y}|\mathbf{x})\mathcal{D}_{\#}\mathcal{E}_{\#}q(\mathbf{x}).$$

813 Let $\mathbf{X} \sim \pi$, $D \sim \mathcal{D}_{\#}\mathcal{E}_{\#}q$, $E \sim \mathcal{E}_{\#}q$, and $\mathbf{Z} \sim \bar{\pi}$, where $\bar{\pi}(\mathbf{z}) \propto p(\mathbf{y}|\mathcal{D}_{\#}(\mathbf{z}))\mathcal{E}_{\#}q(\mathbf{z})$. For any bounded
 814 function f on \mathbb{R}^d , we have, following the definition of π , that

$$815 \quad \mathbb{E}[f(\mathbf{X})] = \frac{\mathbb{E}[p(\mathbf{y}|D)f(D)]}{\mathbb{E}[p(\mathbf{y}|D)]} = \frac{\mathbb{E}[p(\mathbf{y}|\mathcal{D}(E))f(\mathcal{D}(E))]}{\mathbb{E}[p(\mathbf{y}|\mathcal{D}(E))]} = \mathbb{E}[f(\mathcal{D}(\mathbf{Z}))].$$

816 Hence $\text{Law}(\mathbf{X}) = \text{Law}(\mathcal{D}(\mathbf{Z}))$. As a result, to sample approximately from π , we first sample
 817 approximately from $\bar{\pi}$ using any diffusion posterior sampling algorithm with pre-trained DDM for
 818 $\mathcal{E}_{\#}q$, then decode the obtained samples using \mathcal{D} .

822 A.3 MIDPOINT DECOMPOSITION

823 Before we proceed with the midpoint decomposition, we first recall that under the joint distribution
 824 obtained by initializing the posterior π with the forward process (see (3.1)), it holds that for all $i < j$,

$$825 \quad \pi_i(\mathbf{x}_i)q_{j|i}(\mathbf{x}_j|\mathbf{x}_i) = \pi_j(\mathbf{x}_j)\pi_{i|j}(\mathbf{x}_i|\mathbf{x}_j), \quad (\text{A.7})$$

826 where $\pi_{i|j}(\mathbf{x}_i|\mathbf{x}_j) := \pi_i(\mathbf{x}_i)q_{j|i}(\mathbf{x}_j|\mathbf{x}_i)/\pi_j(\mathbf{x}_j)$ and integrates to one.

827 *Proof of Lemma 3.1.* Let $(\ell, k) \in \llbracket 0, n \rrbracket$ be such that $\ell < k$. Applying repeatedly the definition of
 828 the bridge kernel (A.5) and the identity (A.7),

$$829 \quad \begin{aligned} 830 \quad q_{k|\ell, k+1}(\mathbf{x}_k|\mathbf{x}_\ell, \mathbf{x}_{k+1}) &= \frac{q_{k|\ell}(\mathbf{x}_k|\mathbf{x}_\ell)q_{k+1|k}(\mathbf{x}_{k+1}|\mathbf{x}_k)}{q_{k+1|\ell}(\mathbf{x}_{k+1}|\mathbf{x}_\ell)} \\ 831 &= \frac{\pi_\ell(\mathbf{x}_\ell)q_{k|\ell}(\mathbf{x}_k|\mathbf{x}_\ell)q_{k+1|k}(\mathbf{x}_{k+1}|\mathbf{x}_k)}{\pi_\ell(\mathbf{x}_\ell)q_{k+1|\ell}(\mathbf{x}_{k+1}|\mathbf{x}_\ell)} \\ 832 &= \frac{\pi_{\ell|k}(\mathbf{x}_\ell|\mathbf{x}_k)\pi_{k|k+1}(\mathbf{x}_k|\mathbf{x}_{k+1})\pi_{k+1}(\mathbf{x}_{k+1})}{\pi_{\ell|k+1}(\mathbf{x}_\ell|\mathbf{x}_{k+1})\pi_{k+1}(\mathbf{x}_{k+1})} \\ 833 &= \frac{\pi_{\ell|k}(\mathbf{x}_\ell|\mathbf{x}_k)\pi_{k|k+1}(\mathbf{x}_k|\mathbf{x}_{k+1})}{\pi_{\ell|k+1}(\mathbf{x}_\ell|\mathbf{x}_{k+1})}. \end{aligned}$$

834 It then follows that

$$835 \quad \begin{aligned} 836 \quad \pi_{k|k+1}(\mathbf{x}_k|\mathbf{x}_{k+1}) &= \int \pi_{\ell|k}(\mathbf{x}_\ell|\mathbf{x}_k)\pi_{k|k+1}(\mathbf{x}_k|\mathbf{x}_{k+1}) d\mathbf{x}_\ell \\ 837 &= \int q_{k|\ell, k+1}(\mathbf{x}_k|\mathbf{x}_\ell, \mathbf{x}_{k+1})\pi_{\ell|k+1}(\mathbf{x}_\ell|\mathbf{x}_{k+1}) d\mathbf{x}_\ell. \end{aligned}$$

838 \square

851 B THE GAUSSIAN CASE

853 B.1 DERIVATION

854 In this section we derive the recursions verified by the first and second moments of the marginal
 855 distribution $\hat{\pi}_0^\ell$ of the surrogate model (3.6) in the simplified setting of Example 3.2. We recall
 856 that in this specific example we assume that $q = \mathcal{N}(\mathbf{m}, \Sigma)$ where $(\mathbf{m}, \Sigma) \in \mathbb{R}^d \times \mathcal{S}_d^{++}$ and
 857 $p(\mathbf{y}|\cdot) : \mathbf{x} \mapsto \mathcal{N}(\mathbf{y}; \mathbf{A}\mathbf{x}, \sigma_y^2 \mathbf{I}_{d_y})$.

858 **Denoiser and DDPM transitions.** Since we are dealing with a Gaussian prior, the denoiser $\mathbf{m}_{0|k}$
 859 can be computed in closed form for any $k \in \llbracket 1, n \rrbracket$. Using Bishop (2006, Eqn. 2.116), we have that

$$860 \quad \begin{aligned} 861 \quad q_{0|k}(\mathbf{x}_0|\mathbf{x}_k) &\propto q(\mathbf{x}_0)q_{k|0}(\mathbf{x}_k|\mathbf{x}_0) \\ 862 &= \mathcal{N}(\mathbf{x}_0; \Sigma_{0|k}((\sqrt{\alpha_k}/v_k)\mathbf{x}_k + \Sigma^{-1}\mathbf{m}), \Sigma_{0|k}), \end{aligned}$$

where $\Sigma_{0|k} := ((\alpha_k/v_k)\mathbf{I} + \Sigma^{-1})^{-1}$. Hence,

$$\mathbf{m}_{0|k}(\mathbf{x}_k) = \Sigma_{0|k}((\sqrt{\alpha_k}/v_k)\mathbf{x}_k + \Sigma^{-1}\mathbf{m}),$$

and we assume in the remainder of this section that $\mathbf{m}_{0|k}^\theta = \mathbf{m}_{0|k}$. From the expression of $q_{0|k}(\cdot|\mathbf{x}_k)$ we can immediately derive the more general backward transitions $q_{\ell|k}(\cdot|\mathbf{x}_k)$ for $\ell \in \llbracket 1, k-1 \rrbracket$ by noting that

$$q_{\ell|k}(\mathbf{x}_\ell|\mathbf{x}_k) = \int q_{\ell|0,k}(\mathbf{x}_\ell|\mathbf{x}_0, \mathbf{x}_k)q_{0|k}(\mathbf{x}_0|\mathbf{x}_k) d\mathbf{x}_0.$$

From this, (A.3), (A.4), and the law of total expectation and covariance, it follows that $q_{\ell|k}(\cdot|\mathbf{x}_k) = \mathbf{N}(\mathbf{m}_{\ell|k}(\mathbf{x}_k), \Sigma_{\ell|k}(\mathbf{x}_k))$, where

$$\mathbf{m}_{\ell|k}(\mathbf{x}_k) = \mathbf{m}_{\ell|0,k}(\mathbf{m}_{0|k}(\mathbf{x}_k), \mathbf{x}_k), \quad \Sigma_{\ell|k} = \frac{\alpha_\ell(1 - \alpha_k/\alpha_\ell)^2}{(1 - \alpha_k)^2} \Sigma_{0|k} + v_{\ell|0,k} \mathbf{I}_d.$$

On the other hand, the DDPM transitions are

$$p_{\ell|k}^\theta(\cdot|\mathbf{x}_k) = q_{\ell|0,k}(\cdot|\mathbf{m}_{0|k}(\mathbf{x}_k), \mathbf{x}_k) = \mathbf{N}(\mathbf{m}_{\ell|0,k}(\mathbf{m}_{0|k}(\mathbf{x}_k), \mathbf{x}_k), v_{\ell|0,k} \mathbf{I}_d), \quad (\text{B.1})$$

which shows that in this case, the true transitions and approximate ones differ only by their covariance.

Moments recursion. For $k \in \llbracket 0, n \rrbracket$ we let $\hat{\pi}_k^\ell$ denote the \mathbf{x}_k marginal of the surrogate model (3.6). We remind the reader that $\hat{\pi}_n^\ell = \mathbf{N}(\mathbf{0}_d, \mathbf{I}_d)$. The marginals satisfy the recursion

$$\hat{\pi}_k^\ell(\mathbf{x}_k) = \int \hat{\pi}_{k|k+1}^\ell(\mathbf{x}_k|\mathbf{x}_{k+1})\hat{\pi}_{k+1}^\ell(\mathbf{x}_{k+1}) d\mathbf{x}_{k+1}, \quad k \in \llbracket 0, n-1 \rrbracket.$$

Since $\hat{p}_k^\theta(\mathbf{y}|\mathbf{x}_k) = \mathcal{N}(\mathbf{y}; \mathbf{A}\mathbf{m}_{0|k}^\theta(\mathbf{x}_k), \sigma_y^2 \mathbf{I}_d)$ and $\mathbf{m}_{0|k}^\theta$ is linear in \mathbf{x}_k , it is easily seen that $\hat{\pi}_{k|k+1}^\theta(\cdot|\mathbf{x}_{k+1})$ is the density of a Gaussian distribution. Consequently, by definition (3.5) and the definition (A.5) of the bridge kernel, this is also the case for $\hat{\pi}_{k|k+1}^\ell(\cdot|\mathbf{x}_{k+1})$. Now assume that

$$\hat{\pi}_{k+1}^\ell(\mathbf{x}_{k+1}) = \mathcal{N}(\mathbf{x}_{k+1}; \hat{\boldsymbol{\mu}}_{k+1}^\ell, \hat{\boldsymbol{\Sigma}}_{k+1}^\ell), \quad (\text{B.2})$$

$$\hat{\pi}_{k|k+1}^\ell(\mathbf{x}_k|\mathbf{x}_{k+1}) = \mathcal{N}(\mathbf{x}_k; \mathbf{M}_{k|k+1}^\ell \mathbf{x}_{k+1} + \mathbf{c}_{k|k+1}^\ell, \hat{\boldsymbol{\Sigma}}_{k|k+1}^\ell), \quad (\text{B.3})$$

where $\mathbf{M}_{k|k+1}^\ell \in \mathbb{R}^{d \times d}$, $\hat{\boldsymbol{\Sigma}}_{k|k+1}^\ell \in \mathcal{S}_d^{++}$, and $\mathbf{c}_{k|k+1}^\ell \in \mathbb{R}^d$. Using the definition (A.5) of the bridge kernel we find that $\hat{\pi}_k^\ell = \mathbf{N}(\hat{\boldsymbol{\mu}}_k^\ell, \hat{\boldsymbol{\Sigma}}_k^\ell)$, where

$$\hat{\boldsymbol{\mu}}_k^\ell = \mathbf{M}_{k|k+1}^\ell \hat{\boldsymbol{\mu}}_{k+1}^\ell + \mathbf{c}_{k|k+1}^\ell,$$

$$\hat{\boldsymbol{\Sigma}}_k^\ell = \mathbf{M}_{k|k+1}^\ell \hat{\boldsymbol{\Sigma}}_{k+1}^\ell \mathbf{M}_{k|k+1}^{\top} + \hat{\boldsymbol{\Sigma}}_{k|k+1}^\ell.$$

Iterating these updates until reaching $k = 0$, starting from the initialization $\hat{\boldsymbol{\mu}}_n^\ell = \mathbf{0}_d$ and $\hat{\boldsymbol{\Sigma}}_n^\ell = \mathbf{I}_d$, yields the desired moments of the surrogate posterior π_0^ℓ . It now remains to show that the backward transition $\hat{\pi}_{k|k+1}^\ell(\cdot|\mathbf{x}_{k+1})$ writes in the form (B.3) and identify $\mathbf{M}_{k|k+1}^\ell$ and $\mathbf{c}_{k|k+1}^\ell$.

First, we write the approximate likelihood in the form

$$\hat{p}_k^\theta(\mathbf{y}|\mathbf{x}_k) = \mathcal{N}(\mathbf{y}; \hat{\mathbf{A}}_k \mathbf{x}_k + \mathbf{b}_k, \sigma_y^2 \mathbf{I}_d),$$

where

$$\hat{\mathbf{A}}_k = (\sqrt{\alpha_k}/v_k)\mathbf{A}\Sigma_{0|k}, \quad \mathbf{b}_k = \mathbf{A}\Sigma_{0|k}\Sigma^{-1}\mathbf{m}.$$

We also denote by $\mathbf{m}_{\ell|k}^\theta(\mathbf{x}_k)$ the mean of the Gaussian distribution with density given by the DDPM transition $p_{\ell|k}^\theta(\cdot|\mathbf{x}_k)$ in (B.1). We have that

$$\mathbf{m}_{\ell_k|k+1}^\theta(\mathbf{x}_{k+1}) = \mathbf{H}_{\ell_k|k+1} \mathbf{x}_{k+1} + \mathbf{h}_{\ell_k|k+1},$$

where

$$\mathbf{H}_{\ell_k|k+1} := \frac{\alpha_{\ell_k}(1 - \alpha_{k+1}/\alpha_{\ell_k})}{v_{\ell_k} v_{k+1}} \Sigma_{0|k} + \frac{\sqrt{\alpha_{k+1}}(1 - \alpha_{\ell_k})}{v_{k+1}} \mathbf{I}_d, \quad (\text{B.4})$$

$$\mathbf{h}_{\ell_k|k+1} := \frac{\sqrt{\alpha_{\ell_k}}(1 - \alpha_{k+1}/\alpha_{\ell_k})}{v_{k+1}} \Sigma_{0|k} \Sigma^{-1} \mathbf{m}. \quad (\text{B.5})$$

Then, applying (Bishop, 2006, Eqn. 2.116), we get

$$\hat{\pi}_{\ell_k|k+1}^\theta(\mathbf{x}_{\ell_k}|\mathbf{x}_{k+1}) = \mathcal{N}(\mathbf{x}_{\ell_k}; \widetilde{\mathbf{M}}_{\ell_k|k+1}^\theta \mathbf{x}_{k+1} + \tilde{\mathbf{c}}_{\ell_k|k+1}^\theta, \mathbf{\Gamma}_{\ell_k|k+1}),$$

where

$$\begin{aligned} \mathbf{\Gamma}_{\ell_k|k+1} &:= (v_{\ell_k|0,k+1}^{-1} \mathbf{I} + \sigma_{\mathbf{y}}^{-2} \hat{\mathbf{A}}_{\ell_k}^\top \hat{\mathbf{A}}_{\ell_k})^{-1}, \\ \widetilde{\mathbf{M}}_{\ell_k|k+1}^\theta &:= v_{\ell_k|0,k+1}^{-1} \mathbf{\Gamma}_{\ell_k|k+1} \mathbf{H}_{\ell_k|k+1}, \\ \tilde{\mathbf{c}}_{\ell_k|k+1}^\theta &:= \mathbf{\Gamma}_{\ell_k|k+1} [\sigma_{\mathbf{y}}^{-2} \hat{\mathbf{A}}_{\ell_k}^\top (\mathbf{y} - \mathbf{b}_{\ell_k}) + v_{\ell_k|0,k+1}^{-1} \mathbf{h}_{\ell_k|k+1}]. \end{aligned}$$

Finally, following (3.5) we find that

$$\begin{aligned} \mathbf{M}_{k|k+1}^\ell &= \frac{\sqrt{\alpha_k/\alpha_{\ell_k}(1-\alpha_{k+1}/\alpha_k)}}{1-\alpha_{k+1}/\alpha_{\ell_k}} \widetilde{\mathbf{M}}_{\ell_k|k+1}^\theta + \frac{\sqrt{\alpha_{k+1}/\alpha_k(1-\alpha_k/\alpha_{\ell_k})}}{1-\alpha_{k+1}/\alpha_{\ell_k}} \mathbf{I}_d, \\ \mathbf{c}_{k|k+1}^\ell &= \frac{\sqrt{\alpha_k/\alpha_{\ell_k}(1-\alpha_{k+1}/\alpha_k)}}{1-\alpha_{k+1}/\alpha_{\ell_k}} \tilde{\mathbf{c}}_{\ell_k|k+1}^\theta. \end{aligned}$$

B.2 EXPERIMENTAL SETUP

In the experiment described in Example 3.2, we set $d = 100$ and generate 500 instances of inverse problems $(\mathbf{A}, \mathbf{m}, \Sigma)$. For each instance, we compute the Wasserstein-2 distance between the resulting posterior distribution π and $\hat{\pi}_0^{\ell(\eta)}$.

Prior. The mean of the prior is sampled from a standard Gaussian distribution. To generate the covariance $\Sigma \in \mathcal{S}_d^{++}$, we first draw a matrix $\mathbf{G} \in \mathbb{R}^{d \times d}$ with i.i.d. entries sampled from $\mathcal{N}(0, 1)$. Then, for a better conditioning of Σ , we normalize the columns of \mathbf{G} and set $\Sigma = \bar{\lambda}^2 \mathbf{I} + \mathbf{G}\mathbf{G}^\top$, where $\bar{\lambda}^2$ is the mean of the squared singular values of \mathbf{G} .

Likelihood. To sample an ill-posed problem we generate a rank-deficient matrix $\mathbf{A} \in \mathbb{R}^{d_{\mathbf{y}} \times d}$ with $d_{\mathbf{y}} \leq d$. We sample uniformly $d_{\mathbf{y}}$ from the interval $\llbracket d/10, d \rrbracket$ and draw the entries of \mathbf{A} i.i.d. from $\mathcal{N}(0, 1)$. Regarding $\sigma_{\mathbf{y}}$, we sample it uniformly from the interval $[0.1, 0.5]$.

Finally, the resulting posterior is also Gaussian (Bishop, 2006, Eqn. 2.116)

$$\pi(\mathbf{x}) \propto p(\mathbf{y}|\mathbf{x})p(\mathbf{x}) \propto \mathcal{N}(\mathbf{x}; \mathbf{m}_{\mathbf{y}}, \Sigma_{\mathbf{y}}),$$

where $\Sigma_{\mathbf{y}} = (\Sigma^{-1} + (1/\sigma_{\mathbf{y}}^2)\mathbf{A}^\top \mathbf{A})^{-1}$ and $\mathbf{m}_{\mathbf{y}} = \Sigma_{\mathbf{y}}((1/\sigma_{\mathbf{y}}^2)\mathbf{A}^\top \mathbf{y} + \Sigma^{-1}\mathbf{m})$. The Wasserstein-2 distance between the true posterior and $\hat{\pi}_0^{\ell(\eta)}$ is thus (Olkin & Pukelsheim, 1982)

$$W_2(\pi, \hat{\pi}_0^{\ell(\eta)})^2 = \|\mathbf{m}_{\mathbf{y}} - \hat{\boldsymbol{\mu}}_0^{\ell(\eta)}\|^2 + \text{tr}(\Sigma_{\mathbf{y}} + \hat{\Sigma}_0^{\ell(\eta)} - 2(\Sigma_{\mathbf{y}}^{\frac{1}{2}} \hat{\Sigma}_0^{\ell(\eta)} \Sigma_{\mathbf{y}}^{\frac{1}{2}})^{\frac{1}{2}}).$$

C MORE DETAILS ON MGPS

C.1 DETAILS ON THE LOSS

First, by the Data Processing inequality (Van Erven & Harremos, 2014, Example 2),

$$\begin{aligned} &\text{KL}(\lambda_{\ell_k|k+1}^\varphi(\cdot|\mathbf{x}_{k+1}) \parallel \hat{\pi}_{\ell_k|k+1}^\theta(\cdot|\mathbf{x}_{k+1})) \\ &\leq \int \log \frac{q_{k|\ell_k,k+1}(\mathbf{x}_k|\mathbf{x}_{\ell_k}, \mathbf{x}_{k+1}) \lambda_{\ell_k|k+1}^\varphi(\mathbf{x}_{\ell_k})}{q_{k|\ell_k,k+1}(\mathbf{x}_k|\mathbf{x}_{\ell_k}, \mathbf{x}_{k+1}) \hat{\pi}_{\ell_k|k+1}^\theta(\mathbf{x}_{\ell_k}|\mathbf{x}_{k+1})} q_{k|\ell_k,k+1}(\mathbf{x}_k|\mathbf{x}_{\ell_k}, \mathbf{x}_{k+1}) \lambda_{\ell_k|k+1}^\varphi(\mathbf{x}_{\ell_k}) d\mathbf{x}_k d\mathbf{x}_{\ell_k} \\ &= \int \log \frac{\lambda_{\ell_k|k+1}^\varphi(\mathbf{x}_{\ell_k})}{\hat{\pi}_{\ell_k|k+1}^\theta(\mathbf{x}_{\ell_k}|\mathbf{x}_{k+1})} q_{k|\ell_k,k+1}(\mathbf{x}_k|\mathbf{x}_{\ell_k}, \mathbf{x}_{k+1}) \lambda_{\ell_k|k+1}^\varphi(\mathbf{x}_{\ell_k}) d\mathbf{x}_k d\mathbf{x}_{\ell_k} \\ &= \text{KL}(\lambda_{\ell_k|k+1}^\varphi \parallel \hat{\pi}_{\ell_k|k+1}^\theta(\cdot|\mathbf{x}_{k+1})) \end{aligned}$$

The gradient of $\varphi \mapsto \text{KL}(\lambda_{\ell_k|k+1}^\varphi \parallel \hat{\pi}_{\ell_k|k+1}^\theta(\cdot|\mathbf{x}_{k+1}))$ for a given \mathbf{x}_{k+1} writes

$$\begin{aligned} &\nabla_\varphi \mathcal{L}_k(\varphi; \mathbf{X}_{k+1}) \\ &= -\mathbb{E}[\nabla_\varphi \log \hat{p}_{\ell_k}^\theta(\mathbf{y}|\hat{\boldsymbol{\mu}}_{\ell_k} + \text{diag}(e^{\hat{\rho}_{\ell_k}})\mathbf{Z})] + \nabla_\varphi \text{KL}(\lambda_{\ell_k|k+1}^\varphi \parallel p_{\ell_k|k+1}^\theta(\cdot|\mathbf{X}_{k+1})) \quad (\text{C.1}) \end{aligned}$$

and the second term can be expressed in a closed form since both distributions are Gaussians, *i.e.*,

$$\begin{aligned} & \nabla_{\varphi} \text{KL}(\lambda_{\ell_k|k+1}^{\varphi} \parallel p_{\ell_k|k+1}^{\theta}(\cdot|\mathbf{x}_{k+1})) \\ &= \nabla_{\varphi} \left[-\sum_{j=1}^d \hat{\rho}_{\ell_k,j} + \frac{\|\hat{\boldsymbol{\mu}}_{\ell_k} - \mathbf{m}_{\ell_k|0,k+1}(\mathbf{m}_{0|k+1}^{\theta}(\mathbf{x}_{k+1}), \mathbf{x}_{k+1})\|^2 + \sum_{j=1}^d \frac{e^{2\hat{\rho}_{\ell_k,j}}}{v_{\ell_k|0,k+1}}}{2v_{\ell_k|0,k+1}} \right]. \end{aligned}$$

C.2 WARM START

In this section we describe the warm-start approach discussed in the main paper. The complete algorithm with the warm-start procedure is given in Algorithm 3. The original version presented in Algorithm 1 relies on first sampling approximately \mathbf{X}_{ℓ_k} , given \mathbf{X}_{k+1} , from the surrogate transition $\hat{\pi}_{\ell_k|k+1}^{\theta}(\cdot|\mathbf{X}_{k+1})$ and then sampling \mathbf{X}_k from $q_{k|\ell_k,k+1}(\cdot|\mathbf{X}_{\ell_k}, \mathbf{X}_{k+1})$. In our warm-start approach, which we apply only during the first iterations of the algorithm (see the hyperparameter settings in Table 6), we draw inspiration from the decomposition

$$\pi_{k|k+1}(\mathbf{x}_k|\mathbf{x}_{k+1}) = \int q_{k|1,k+1}(\mathbf{x}_k|\mathbf{x}_1, \mathbf{x}_{k+1}) \pi_{1|\ell_k}(\mathbf{x}_1|\mathbf{x}_{\ell_k}) \pi_{\ell_k|k+1}(\mathbf{x}_{\ell_k}|\mathbf{x}_{k+1}) d\mathbf{x}_1 d\mathbf{x}_{\ell_k}. \quad (\text{C.2})$$

It suggests introducing a second intermediary step that involves sampling approximately from the transition $\mathbf{X}_1 \sim \pi_{1|\ell_k}(\cdot|\mathbf{X}_{\ell_k})$. \mathbf{X}_k is then sampled from the bridge $q_{k|1,k+1}(\cdot|\mathbf{X}_1, \mathbf{X}_{k+1})$.

In order to draw approximate samples from $\pi_{1|\ell_k}(\cdot|\mathbf{X}_{\ell_k})$ we leverage again a Gaussian variational approximation $\lambda_{1|\ell_k}^{\psi} := N(\tilde{\boldsymbol{\mu}}_1, \text{diag}(e^{2\tilde{\rho}_1}))$, which we fit by minimizing a proxy of the KL divergence between $\lambda_{1|\ell_k}^{\psi}$ and $\pi_{1|\ell_k}(\cdot|\mathbf{X}_{\ell_k})$ of which the gradient w.r.t. $\psi := (\tilde{\boldsymbol{\mu}}_1, \tilde{\boldsymbol{\rho}}_1)$ writes

$$\begin{aligned} \nabla_{\psi} \mathcal{L}_{1|\ell_k}^{\text{ws}}(\psi; \mathbf{X}_{\ell_k})|_{\psi_0} &:= -\nabla_{\psi} \left[\sum_{j=1}^d \tilde{\rho}_{1,j} - \frac{\|\mathbf{X}_{\ell_k} - (\alpha_{\ell_k}/\alpha_1)^{1/2} \tilde{\boldsymbol{\mu}}_1\|^2 + (\alpha_{\ell_k}/\alpha_1) \sum_{j=1}^d e^{\tilde{\rho}_{1,j}^2}}{2(1 - \alpha_{\ell_k}/\alpha_1)} \right] |_{\psi_0} \\ &- \mathbb{E} \left[\nabla_{\psi} (\tilde{\boldsymbol{\mu}}_1 + \text{diag}(e^{\tilde{\rho}_1}) \mathbf{Z})^{\top} (\nabla_x \log p(\mathbf{y}|\cdot) + \mathbf{s}_1^{\theta}(\tilde{\boldsymbol{\mu}}_1 + \text{diag}(e^{\tilde{\rho}_1}) \mathbf{Z})|_{\psi_0}) \right]. \quad (\text{C.3}) \end{aligned}$$

This expression corresponds to the gradient w.r.t. ψ of $\text{KL}(\lambda_{1|\ell_k}^{\psi} \parallel \pi_{1|\ell_k}(\cdot|\mathbf{X}_{\ell_k}))$ combined with the score approximation \mathbf{s}_1^{θ} of $\nabla_x \log q_1$ and the mild likelihood approximation $\nabla_x \log p_1(\mathbf{y}|\cdot)$ of $\nabla_x \log p(\mathbf{y}|\cdot)$. Indeed, we have that

$$\begin{aligned} \text{KL}(\lambda_{1|\ell_k}^{\psi} \parallel \pi_{1|\ell_k}(\cdot|\mathbf{X}_{\ell_k})) &= -\mathbb{E}_{\lambda_{1|\ell_k}^{\psi}} [\log p_1(\mathbf{y}|\mathbf{X}_1) + \log q_1(\mathbf{X}_1)] \\ &+ \int \log \frac{\lambda_{1|\ell_k}^{\psi}(\mathbf{x}_1)}{q_{\ell_k|1}(\mathbf{X}_{\ell_k}|\mathbf{x}_1)} \lambda_{1|\ell_k}^{\psi}(\mathbf{x}_1) d\mathbf{x}_1. \end{aligned}$$

The second term can be computed in a closed form and its gradient corresponds to the first term in (C.3). As for the first term, under standard differentiability assumptions and applying the reparameterization trick, we find that

$$\begin{aligned} & \nabla_{\psi} \mathbb{E}_{\lambda_{1|\ell_k}^{\psi}} [\log p_1(\mathbf{y}|\mathbf{X}_1) + \log q_1(\mathbf{X}_1)] |_{\psi_0} \\ &= \nabla_{\psi} \mathbb{E} [(\log p_1(\mathbf{y}|\cdot) + \log q_1)(\tilde{\boldsymbol{\mu}}_1 + \text{diag}(e^{\tilde{\rho}_1}) \mathbf{Z})] |_{\psi_0} \\ &= \mathbb{E} \left[\nabla_{\psi} (\tilde{\boldsymbol{\mu}}_1 + \text{diag}(e^{\tilde{\rho}_1}) \mathbf{Z})^{\top} (\nabla_x \log p_1(\mathbf{y}|\cdot) + \nabla_x \log q_1)(\tilde{\boldsymbol{\mu}}_1 + \text{diag}(e^{\tilde{\rho}_1}) \mathbf{Z})|_{\psi_0} \right]. \end{aligned}$$

Plugging the previous approximations yields the second term in (C.3). The warm-start procedure is summarized in Algorithm 2. We also use a single sample Monte Carlo estimate to perform the optimization. Finally, at step k and given $\hat{\mathbf{X}}_{\ell_k}$, the initial parameter ψ_k of the variational approximation is chosen so that

$$\lambda_{1|\ell_k}^{\psi_k}(\mathbf{x}_1) = q_{1|0,\ell_k}(\mathbf{x}_1|\mathbf{m}_{0|\ell_k}^{\theta}(\hat{\mathbf{X}}_{\ell_k}), \hat{\mathbf{X}}_{\ell_k}).$$

Algorithm 2 WarmStart

```

1026 1: Input: step  $k$ , samples  $(\hat{\mathbf{X}}_{\ell_k}, \mathbf{X}_{k+1})$ , gradient steps  $M$ 
1027 2:  $\tilde{\boldsymbol{\mu}}_1 \leftarrow \mathbf{m}_{1|0,\ell_k}(\mathbf{m}_{0|\ell_k}^\theta(\hat{\mathbf{X}}_{\ell_k}, \mathbf{X}_{\ell_k}), \tilde{\boldsymbol{\rho}}_1 \leftarrow \frac{1}{2} \log v_{1|0,\ell_k}$ 
1028 3: for  $j = 1$  to  $M$  do
1029 4:    $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$ 
1030 5:    $(\tilde{\boldsymbol{\mu}}_1, \tilde{\boldsymbol{\rho}}_1) \leftarrow \text{OptimizerStep}(\nabla_{\psi} \tilde{\mathcal{L}}_{1|\ell_k}^{\text{ws}}(\cdot, \mathbf{Z}; \mathbf{X}_{\ell_k}); \tilde{\boldsymbol{\mu}}_1, \tilde{\boldsymbol{\rho}}_1)$ 
1031 6: end for
1032 7:  $\hat{\mathbf{X}}_1 \leftarrow \tilde{\boldsymbol{\mu}}_1 + \text{diag}(e^{\hat{\rho}_1}) \mathbf{Z}_1$  where  $\mathbf{Z}_1 \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$ 
1033 8:  $\mathbf{X}_k \leftarrow \mathbf{m}_{k|1,k+1}(\hat{\mathbf{X}}_1, \mathbf{X}_{k+1}) + v_{k|1,k+1}^{1/2} \mathbf{Z}_k$ 
1034 9: Output:  $\mathbf{X}_k, \hat{\mathbf{X}}_1$ 

```

Algorithm 3 MGPS with warm start strategy

```

1038 1: Input:  $(\ell_k)_{k=1}^n$  with  $\ell_n = n, \ell_1 = 1$ , gradient steps  $(M_k)_{k=1}^{n-1}$ , warm start threshold  $w$ .
1039 2:  $\mathbf{X}_n \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$ ,  $\hat{\mathbf{X}}_n \leftarrow \mathbf{X}_n$ 
1040 3:  $\hat{\mathbf{X}}_0 \leftarrow \mathbf{m}_{0|n}^\theta(\hat{\mathbf{X}}_n)$ 
1041 4: for  $k = n - 1$  to  $1$  do
1042 5:    $\hat{\boldsymbol{\mu}}_{\ell_k} \leftarrow \mathbf{m}_{\ell_k|0,k+1}(\hat{\mathbf{X}}_0, \mathbf{X}_{k+1}), \hat{\boldsymbol{\rho}}_{\ell_k} \leftarrow \frac{1}{2} \log v_{\ell_k|0,k+1}$ 
1043 6:   for  $j = 1$  to  $M_k$  do
1044 7:      $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$ 
1045 8:      $(\hat{\boldsymbol{\mu}}_{\ell_k}, \hat{\boldsymbol{\rho}}_{\ell_k}) \leftarrow \text{OptimizerStep}(\nabla_{\varphi} \tilde{\mathcal{L}}_{\ell_k}(\cdot, \mathbf{Z}; \mathbf{X}_{k+1}); \hat{\boldsymbol{\mu}}_{\ell_k}, \hat{\boldsymbol{\rho}}_{\ell_k})$ 
1046 9:   end for
1047 10:  $\mathbf{Z}_{\ell_k}, \mathbf{Z}_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$ 
1048 11:  $\hat{\mathbf{X}}_{\ell_k} \leftarrow \hat{\boldsymbol{\mu}}_{\ell_k} + \text{diag}(e^{\hat{\rho}_{\ell_k}}) \mathbf{Z}_{\ell_k}$ 
1049 12: if  $k \geq w$  then
1050 13:    $(\mathbf{X}_k, \hat{\mathbf{X}}_1) \leftarrow \text{WarmStart}(k, \hat{\mathbf{X}}_{\ell_k}, \mathbf{X}_{k+1}, M_k)$ 
1051 14:    $\hat{\mathbf{X}}_0 \leftarrow \hat{\mathbf{X}}_1$ 
1052 15: else
1053 16:    $\mathbf{X}_k \leftarrow \mathbf{m}_{k|\ell_k,k+1}(\hat{\mathbf{X}}_{\ell_k}, \mathbf{X}_{k+1}) + v_{k|\ell_k,k+1}^{1/2} \mathbf{Z}_k$ 
1054 17:    $\hat{\mathbf{X}}_0 \leftarrow \mathbf{m}_{0|\ell_k}^\theta(\mathbf{X}_{\ell_k})$ 
1055 18: end if
1056 19: end for
1057 20: Output:  $\hat{\mathbf{X}}_0$ 

```

After having sampled $\tilde{\mathbf{X}}_1 \sim \lambda_{1|\ell_k}^{\psi_k^*}$, we use it to first sample $\mathbf{X}_k \sim q_{k|1,k+1}(\cdot|\tilde{\mathbf{X}}_1, \mathbf{X}_{k+1})$ and then initialize the next variational approximation $\lambda_{\ell_k-1|k}^\varphi$. The initial parameter φ_{k-1} is set so that

$$\lambda_{\ell_k-1|k}^{\varphi_{k-1}}(\mathbf{x}_{\ell_k-1}) = q_{\ell_k-1|0,k}(\mathbf{x}_{\ell_k-1}|\tilde{\mathbf{X}}_1, \mathbf{X}_k),$$

see line 5 in Algorithm 3.

C.3 DIFFERENCE WITH DPS

In this section we detail the differences between MGPS and the DPS algorithm proposed in (Chung et al., 2023). While DPS cannot be seen as an instantiation of our methodology, we highlight the main differences by deriving a specific case of MGPS that closely resembles DPS.

We assume that $\ell_k = k$ for all $k \in \llbracket 1, n-1 \rrbracket$ and that $p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}; \mathcal{A}(\mathbf{x}), \sigma_{\mathbf{y}}^2 \mathbf{I}_{d_y})$ where $\mathcal{A}: \mathbb{R}^d \rightarrow \mathbb{R}^{d_y}$. We consider the same optimization procedure as in Algorithm 1 but we: (i) optimize only the mean parameter $\hat{\boldsymbol{\mu}}_k$ and fix the covariance of the variational approximation to $v_{k|k+1} \mathbf{I}_d$, (ii) perform a single stochastic optimization step.

Following that, the initialization strategy (3.9) boils down to setting

$$\lambda_{k|k+1}^{\varphi_k}(\mathbf{x}_k|\mathbf{x}_{k+1}) = p_{k|k+1}^\theta(\mathbf{x}_k|\mathbf{x}_{k+1}) \tag{C.4}$$

which means that the initial mean parameter is $\hat{\boldsymbol{\mu}}_k^0 := \mathbf{m}_{k|0,k+1}(\mathbf{m}_{0|k+1}^\theta(\mathbf{x}_{k+1}), \mathbf{x}_{k+1})$. With this initialization $\varphi_k \mapsto \text{KL}(\lambda_{k|k+1}^{\varphi_k}(\cdot|\mathbf{x}_{k+1}) \parallel p_{k|k+1}^\theta(\cdot|\mathbf{x}_{k+1}))$ has its global optimum attained at $\varphi_k = \hat{\boldsymbol{\mu}}_k^0$. Thus, the gradient (C.1) writes $\nabla_{\varphi_k} \mathcal{L}_k(\varphi_k; \mathbf{x}_{k+1})|_{\hat{\boldsymbol{\mu}}_k^0} = -\mathbb{E}[\nabla_{\varphi_k} \log \hat{p}_k^\theta(\mathbf{y}|\hat{\boldsymbol{\mu}}_k + \sqrt{v_{k|k+1}}\mathbf{Z})|_{\hat{\boldsymbol{\mu}}_k^0}]$. Next, updating the parameters of the variational approximation using a single stochastic gradient descent step with the gradient estimate $-\nabla_{\varphi_k} \log \hat{p}_k^\theta(\mathbf{y}|\hat{\boldsymbol{\mu}}_k + \sqrt{v_{k|k+1}}\mathbf{Z}_k)|_{\hat{\boldsymbol{\mu}}_k^0}$, where $\mathbf{Z}_k \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$ and step-size

$$\gamma_k = \frac{\zeta \sigma_y^2}{\|\mathbf{y} - \mathcal{A}(\mathbf{m}_{0|k}^\theta(\hat{\boldsymbol{\mu}}_k^0 + \sqrt{v_{k|k+1}}\mathbf{Z}_k))\|}$$

yields the variational approximation

$$\lambda_{k|k+1}^{\varphi_k^*}(\mathbf{x}_k|\mathbf{x}_{k+1}) = \mathcal{N}(\mathbf{x}_k; \hat{\boldsymbol{\mu}}_k^0 - \zeta \nabla_{\varphi_k} \|\mathbf{y} - \mathcal{A}(\mathbf{m}_{0|k}^\theta(\hat{\boldsymbol{\mu}}_k + \sqrt{v_{k|k+1}}\mathbf{Z}_k))\|_{|\hat{\boldsymbol{\mu}}_k^0}, v_{k|k+1}\mathbf{I}_d).$$

In order to simplify the expression we have used the fact that

$$\frac{-\nabla_{\varphi_k} \log \hat{p}_k^\theta(\mathbf{y}|\hat{\boldsymbol{\mu}}_k + \sqrt{v_{k|k+1}}\mathbf{Z}_k)|_{\hat{\boldsymbol{\mu}}_k^0}}{\|\mathbf{y} - \mathcal{A}(\mathbf{m}_{0|k}^\theta(\hat{\boldsymbol{\mu}}_k + \sqrt{v_{k|k+1}}\mathbf{Z}_k))\|_{|\hat{\boldsymbol{\mu}}_k^0}} = \frac{1}{\sigma_y^2} \nabla_{\varphi_k} \|\mathbf{y} - \mathcal{A}(\mathbf{m}_{0|k}^\theta(\hat{\boldsymbol{\mu}}_k + \sqrt{v_{k|k+1}}\mathbf{Z}_k))\|_{|\hat{\boldsymbol{\mu}}_k^0}.$$

Therefore, given \mathbf{X}_{k+1} , a draw from the variational approximation in MGPS is obtained following the update

$$\mathbf{X}_k = \tilde{\mathbf{X}}_k - \zeta \nabla_{\mathbf{x}_k} \|\mathbf{y} - \mathcal{A}(\mathbf{m}_{0|k}^\theta(\mathbf{x}_k))\|_{|\mathbf{x}_k = \tilde{\mathbf{X}}_k}, \quad (\tilde{\mathbf{X}}_k, \tilde{\mathbf{X}}_k') \stackrel{\text{i.i.d.}}{\sim} p_{k|k+1}^\theta(\cdot|\mathbf{X}_{k+1}).$$

On the other hand, the DPS transition is given by

$$\lambda_{k|k+1}^{\text{DPS}}(\mathbf{x}_k|\mathbf{x}_{k+1}) := \mathcal{N}(\mathbf{x}_k; \hat{\boldsymbol{\mu}}_k^0 - \zeta \nabla_{\mathbf{x}_{k+1}} \|\mathbf{y} - \mathcal{A}(\mathbf{m}_{0|k+1}^\theta(\mathbf{x}_{k+1}))\|, v_{k|k+1}\mathbf{I}_d).$$

and, hence a sample $\mathbf{X}_k^{\text{DPS}}$ is drawn following

$$\mathbf{X}_k^{\text{DPS}} = \tilde{\mathbf{X}}_k - \zeta \nabla_{\mathbf{x}_{k+1}} \|\mathbf{y} - \mathcal{A}(\mathbf{m}_{0|k+1}^\theta(\mathbf{x}_{k+1}))\|_{|\mathbf{x}_{k+1} = \mathbf{X}_{k+1}}, \quad \tilde{\mathbf{X}}_k \sim p_{k|k+1}^\theta(\cdot|\mathbf{X}_{k+1}).$$

The difference between MGPS and DPS is in: (i) the diffusion step used for the denoiser (k for MGPS, $k+1$ for DPS), (ii) the sample where the gradient is evaluated.

D EXPERIMENTS

In this section we provide the implementation details on our algorithm as well as the algorithms we benchmark against. We use the the hyperparameters recommended by the authors and tune them on each dataset if they are not provided.

D.1 IMPLEMENTATION DETAILS AND HYPERPARAMETERS FOR MGPS

We implement Algorithm 1 with $\ell_k = \lfloor k/2 \rfloor$ and use the Adam optimizer (Kingma & Ba, 2015) with a learning rate of 0.03 for optimization. The number of gradient steps is adjusted based on the complexity of the task: posterior sampling with the ImageNet DDM prior or FFHQ LDM prior is more challenging and therefore requires additional gradient steps. Detailed hyperparameters are provided in Table 6.

The warm-start strategy outlined in Appendix C.2 improved reconstruction plausibility and eliminated potential artifacts. A similar effect was observed when performing multiple gradient steps ($M = 20$) during the initial stages. For latent-space models, switching the intermediate step to $\ell_k = k$ for the second half of the diffusion process has been crucial and significantly enhanced reconstruction quality by mitigating the smoothing effect, which often removes important details. A similar strategy has been useful for the Gaussian deblurring and motion deblurring tasks on the ImageNet dataset.

D.2 IMPLEMENTATION DETAILS OF COMPETITORS

DPS. We implemented Chung et al. (2023, Algorithm 1) and refer to Chung et al. (2023, App. D) for the values of its hyperparameters. After tuning, we adopt $\gamma = 0.2$ for JPEG 2% and $\gamma = 0.07$ for High Dynamic Range tasks.

Table 6: The hyperparameters used in MGPS for the considered datasets.

	Warm start	Threshold (w)	Diffusion steps	ℓ_k	Learning rate	Gradient steps
FFHQ	✓	$\lfloor 3n/4 \rfloor$	$n \in \{50, 100, 300\}$	$\lfloor k/2 \rfloor$	0.03	$M_k = \begin{cases} 20 & \text{if } k \geq n - 5 \\ 20 & \text{if } k \bmod 10 = 0 \\ 2 & \text{otherwise} \end{cases}$
FFHQ (LDM)	✓	$\lfloor 3n/4 \rfloor$	$n \in \{50, 100, 300\}$	$\ell_k = \begin{cases} \lfloor k/2 \rfloor & \text{if } k > \lfloor n/2 \rfloor \\ k & \text{otherwise} \end{cases}$	0.03	$M_k = \begin{cases} 20 & \text{if } k \geq n - 5 \\ 10 & \text{if } k \bmod 10 = 0 \\ 5 & \text{otherwise} \end{cases}$
ImageNet	✓	$\lfloor 3n/4 \rfloor$	$n \in \{50, 100, 300\}$	$\lfloor k/2 \rfloor$	0.03	$M_k = \begin{cases} 20 & \text{if } k \geq n - 5 \\ 10 & \text{if } k \bmod 20 = 0 \\ 2 & \text{otherwise} \end{cases}$
Gaussian Mixture	✗	-	$n = 300$	$\lfloor k/2 \rfloor$	0.1	$M_k = \begin{cases} 20 & \text{if } k \geq n - 5 \\ 20 & \text{if } k \bmod 10 = 0 \\ 2 & \text{otherwise} \end{cases}$
PTB-XL (ECG)	✓	$\lfloor 3n/4 \rfloor$	$n \in \{50, 300\}$	$\lfloor k/2 \rfloor$	0.03	$M_k = \begin{cases} 20 & \text{if } k \geq n - 5 \\ 5 & \text{if } k \bmod 20 = 0 \\ 5 & \text{otherwise} \end{cases}$

DiffPIR. We implemented [Zhu et al. \(2023, Algorithm 1\)](#) to make it compatible with our existing code base. We used the hyperparameters recommended in the official, released version¹. Unfortunately, we did not manage to make the algorithm converge for nonlinear problems. While the authors give some guidelines to handle such problems ([Zhu et al., 2023, Eqn. \(13\)](#)), examples are missing in the paper and the released code. Similarly, we do not run it on the motion deblur task as the FFT-based solution provided in ([Zhu et al., 2023](#)) is only valid for circular convolution, yet we opted for the experimental setup of [Chung et al. \(2023\)](#) which uses convolution with reflect padding.

DDNM. We adapted the implementation in the released code² to our code base. Since DDNM utilizes the pseudo-inverse of the degradation operator, we noticed that it is unstable for operators whose SVD are prone to numerical errors, such as Gaussian Blur with wide convolution kernel.

RedDiff. We used the implementation of REDDIFF available in the released code³. On nonlinear problems, for which the pseudo-inverse of the observation is not available, we initialized the variational optimization with a sample from the standard Gaussian distribution.

PGDM. We relied on the implementation provided in REDDIFF’s code as some of its authors are co-authors of PGDM. The implementation features a subtle difference with [Song et al. \(2023a, Algorithm 1\)](#): in the last line of the algorithm, the guidance term g is multiplied by $\sqrt{\alpha_t}$ but in the implementation it is multiplied by $\sqrt{\alpha_{t-1}\alpha_t}$. This modification stabilizes the algorithm on most tasks. For JPEG 2% we found that it worsens the performance. In this case we simply multiply by $\sqrt{\alpha_t}$, as in the original algorithm.

PSLD We implemented the PSLD algorithm provided in [Rout et al. \(2024, Algorithm \(2\)\)](#) and used the hyperparameters provided by the authors in the publicly available implementation⁴.

ReSample We used the original code provided by the authors⁵ and modified it to expose several hyperparameters that were not directly accessible in the released version. Specifically, we exposed the tolerance ε and the maximum number of iterations N for solving the optimization problems related to hard data consistency, as well as the scaling factor for the variance of the stochastic resampling distribution γ . Our experiments revealed that the algorithm is highly sensitive to ε . We found that setting it equal to the noise level of the inverse problem gave the best reconstruction across tasks and noise levels. We set the maximum number of gradient iterations to $N = 200$ to make the algorithm less computationally prohibitive. Finally, we tuned γ for each task but found it has less impact on the quality of the reconstruction compared to ε .

¹<https://github.com/yuanzhi-zhu/DiffPIR>

²<https://github.com/wyhuai/DDNM>

³<https://github.com/NVLabs/RED-diff>

⁴<https://github.com/LituRout/PSLD>

⁵<https://github.com/soominkwon/resample>

D.2.1 RUNTIME AND MEMORY

To get the runtime and GPU memory consumption of an algorithm on a dataset, we average these two metrics over both samples of the dataset and the considered tasks in Section 4.

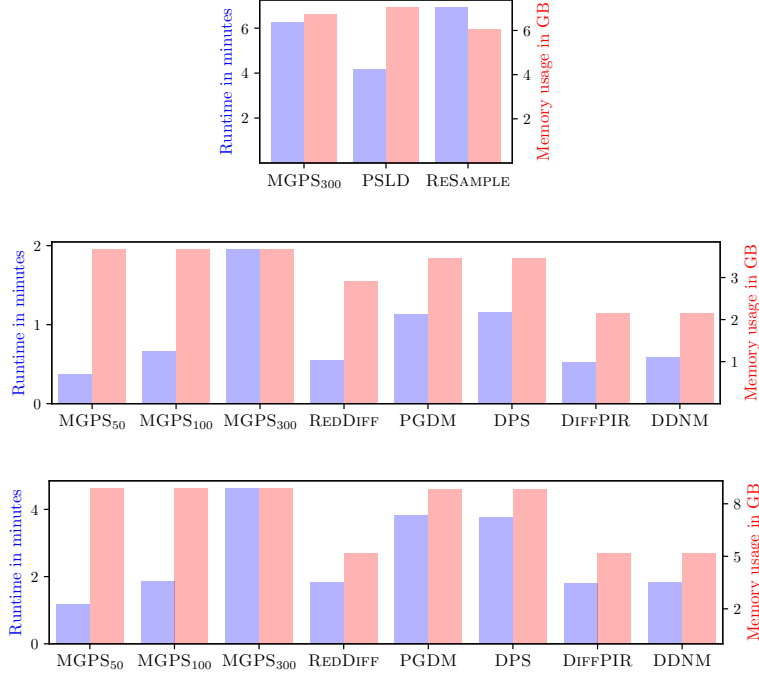


Figure 7: Runtime and memory requirement of the considered algorithms on datasets: FFHQ latent space (1st row), FFHQ pixel space (2nd row), and ImageNet (3rd row). The left axis displays the runtime in minutes, whereas the right axis the GPU memory requirement in Gigabytes (GB).

D.3 GAUSSIAN MIXTURES

In this section, we elaborate more on the Gaussian mixture experiment in Section 4, where we consider a linear inverse problem with a Gaussian mixture as prior. Therefore, the likelihood is $p(\mathbf{y}|\cdot) : \mathbf{x} \mapsto \mathcal{N}(\mathbf{y}; \mathbf{A}\mathbf{x}, \sigma_y^2 \mathbf{I}_{d_y})$. Recall that a Gaussian mixture with $C \in \mathbb{N}$ components, whose weights, means, and covariances are $w_i > 0$, $\mathbf{m}_i \in \mathbb{R}^d$, and $\sigma_i^2 \mathbf{I}_d$, respectively, has the density

$$q(\mathbf{x}) = \sum_{i=1}^C w_i \mathcal{N}(\mathbf{x}; \mathbf{m}_i, \sigma_i^2 \mathbf{I}_d).$$

Denoyer. The denoyer is obtained via Tweedie’s formula (Robbins, 1956),

$$\mathbf{m}_{0|k}(\mathbf{x}_k) = (\mathbf{x}_k + (1 - \alpha_k) \nabla \log q_k(\mathbf{x}_k)) / \sqrt{\alpha_k},$$

where the densities of the marginals $(q_k)_{k=1}^n$ are straightforward,

$$\begin{aligned} q_k(\mathbf{x}_k) &= \int q(\mathbf{x}_0) q_{k|0}(\mathbf{x}_k | \mathbf{x}_0) d\mathbf{x}_0 = \sum_{i=1}^C w_i \int \mathcal{N}(\mathbf{x}_0; \mathbf{m}_i, \sigma_i^2 \mathbf{I}_d) q_{k|0}(\mathbf{x}_k | \mathbf{x}_0) d\mathbf{x}_0 \\ &= \sum_{i=1}^C w_i \mathcal{N}(\mathbf{x}_k; \sqrt{\alpha_k} \mathbf{m}_i, (\alpha_k \sigma_i^2 + v_k) \mathbf{I}_d). \end{aligned}$$

Posterior. Furthermore, following Bishop (2006, Eqn. 2.116), the posterior can be shown to be a Gaussian mixture

$$\pi(\mathbf{x}) \propto g(\mathbf{x}) q(\mathbf{x}) \propto \sum_{i=1}^C \bar{w}_i \mathcal{N}(\mathbf{x}; \bar{\mathbf{m}}_i, \bar{\Sigma}_i),$$

with new weights, covariances, and means given by

$$\begin{aligned}\bar{w}_i &= w_i \mathcal{N}(\mathbf{y}; \mathbf{A} \mathbf{m}_i, \sigma_y^2 \mathbf{I}_{d_y} + \sigma_i^2 \mathbf{A} \mathbf{A}^\top), \\ \bar{\Sigma}_i &= ((1/\sigma_i^2) \mathbf{I}_d + (1/\sigma_y^2) \mathbf{A}^\top \mathbf{A})^{-1}, \\ \bar{\mathbf{m}}_i &= \bar{\Sigma}_i ((1/\sigma_y^2) \mathbf{A}^\top \mathbf{y} + (1/\sigma_i^2) \mathbf{m}_i),\end{aligned}$$

where the weights are un-normalized.

Experimental setup. We consider two setups where $d \in \{20, 200\}$ and generate Gaussian mixtures of $C = 25$ components with means $(\mathbf{m}_i)_{i=1}^C = \{(8i, 8j, \dots, 8i, 8j) \in \mathbb{R}^d : (i, j) \in \llbracket -2, 2 \rrbracket^2\}$, unit covariances, and weights w_i drawn from a uniform distribution on $[0, 1]$ then normalized to sum to 1. The linear inverse problem (\mathbf{y}, \mathbf{A}) is generated by first drawing the matrix $\mathbf{A} \in \mathbb{R}^{1 \times d}$ with entries i.i.d. from $\mathcal{N}(0, 1)$, then sampling \mathbf{x}^* from the prior q and finally computing $\mathbf{y} = \mathbf{A} \mathbf{x}^* + \sigma_y \varepsilon$ with $\varepsilon \sim \mathcal{N}(0, 1)$ to get the observation. The standard deviation of the inverse problem is fixed to $\sigma_y = 0.05$ across all problem instances.

To assess the performance of each algorithm, we draw 2000 samples and compare against 2000 samples from the true posterior distribution using the Sliced Wasserstein (SW) distance by averaging over 10^4 slices. In Table 1 we report the average SW and the 95% confidence interval over 100 replicates.

D.4 EXTENDED IMAGE EXPERIMENTS

We extend, in Tables 7 and 8, the results in Table 2 by including MGPS with $n = 50$ and $n = 100$. In the setting $n = 100$, the runtime of MGPS is twice lower than that of DPS, PGDM and comparable to that of DIFFPIR, DDNM and REDDIFF, see Figure 7. It is also outperforming them on most of the tasks, especially the nonlinear ones, as is seen in the table below. In the setting $n = 50$, MGPS has the lowest runtime among all the methods while maintaining competitive performance.

Finally, in Appendix D.6, we present sample reconstructions on the various tasks. For each algorithm, we generate five samples and select the four most visually appealing ones. Completely black or white images correspond to failure cases of DPS.

D.5 ECG EXPERIMENTS

D.5.1 IMPLEMENTATION DETAILS

PTB-XL dataset We use 12-lead ECGs at a sampling frequency of 100 Hz from the PTB-XL dataset (Wagner et al., 2020). For both training and generation, we do not use the augmented limb leads aVL, aVR and aVF as they can be obtained with the following relations: $aVL = (I - III)/2$, $aVF = (II + III)/2$ and $aVR = -(I + II)/2$. This leads to inputs of shape $T \times 9$ instead of $T \times 12$ (where T is the length of the signal). For evaluation metrics, we reconstruct the augmented limb leads. For the MB task, following (Alcaraz & Strodthoff, 2022), we use 2.56-second ECG random crops, hence the inputs are of shape 256×9 . For the ML task, we use the full 10-second ECGs and pad them into a tensor of shape 1024×9 . Table 10 summarizes the input shapes for models and algorithms for each task. Table 11 summarizes the distribution of train/val/test splits and the number of cardiac conditions (RBBB, LBBB, AF, SB) per split. Note that for posterior sampling, we only use the test set, and for both training and sampling we never use the cardiac condition.

Diffusion model Following (Alcaraz & Strodthoff, 2023; 2022) we used Structured State Space Diffusion (SSSD) models (Gu et al., 2022; Goel et al., 2022). SSSD generates a sequence $u(t)$ with the following differential equation depending on the input sequence $\mathbf{x}(t)$ and a hidden state $h(t)$:

$$h'(t) = Ah(t) + Bx(t) \tag{D.1}$$

$$u(t) = Ch(t) + Dx(t), \tag{D.2}$$

where A, B, C, D are transition matrices. We use the SSSD^{SA} (Alcaraz & Strodthoff, 2022) (also denoted as Sashimi (Goel et al., 2022)) architecture publicly available⁶⁷. We parametrize the matrix

⁶<https://github.com/AI4HealthUOL/SSSD>

⁷<https://github.com/albertfgu/diffwave-sashimi/tree/master>

Table 7: Mean LPIPS/PSNR/SSIM values for various linear and nonlinear imaging tasks on the FFHQ 256×256 dataset. Best is in **bold** and second best is underlined.

Task	MGPS ₅₀	MGPS ₁₀₀	MGPS ₃₀₀	DPS	PGDM	DDNM	DIFFPIR	REDDIFF
LPIPS ↓								
SR ($\times 4$)	0.13	<u>0.10</u>	0.09	0.09	0.33	0.14	0.13	0.36
SR ($\times 16$)	0.27	<u>0.26</u>	<u>0.26</u>	0.24	0.44	0.30	0.28	0.51
Box inpainting	0.16	<u>0.12</u>	0.10	0.19	0.17	<u>0.12</u>	0.18	0.19
Half mask	0.24	<u>0.22</u>	0.20	0.24	0.26	<u>0.22</u>	0.23	0.28
Gaussian Deblur	0.21	0.18	<u>0.15</u>	0.16	0.87	0.19	0.12	0.26
Motion Deblur	0.19	<u>0.15</u>	0.13	0.16	–	–	–	0.21
JPEG (QF = 2)	0.20	<u>0.17</u>	0.16	0.39	1.10	–	–	<u>0.32</u>
Phase retrieval	0.20	<u>0.14</u>	0.11	0.46	–	–	–	<u>0.25</u>
Nonlinear deblur	0.23	0.23	0.23	<u>0.52</u>	–	–	–	0.66
High dynamic range	0.13	<u>0.09</u>	0.07	0.49	–	–	–	<u>0.20</u>
PSNR ↑								
SR ($\times 4$)	27.83	27.79	27.79	<u>28.24</u>	23.34	29.52	27.17	27.25
SR ($\times 16$)	20.45	20.34	20.22	20.67	17.65	22.43	20.75	<u>21.91</u>
Box inpainting	21.55	22.22	22.68	18.39	21.13	<u>22.35</u>	21.96	21.79
Half mask	15.10	15.32	15.54	14.82	16.03	<u>16.16</u>	15.17	16.21
Gaussian Deblur	25.09	25.19	25.89	24.20	13.36	<u>26.69</u>	25.89	26.72
Motion Deblur	26.07	<u>26.64</u>	26.48	24.24	–	–	–	27.58
JPEG (QF = 2)	<u>25.00</u>	25.23	24.94	18.50	12.76	–	–	24.42
Phase retrieval	24.20	<u>26.60</u>	27.25	14.87	–	–	–	24.85
Nonlinear deblur	24.16	<u>24.21</u>	24.37	15.89	–	–	–	21.97
High dynamic range	24.77	<u>26.07</u>	27.74	16.83	–	–	–	21.25
SSIM ↑								
SR ($\times 4$)	0.81	0.80	0.79	<u>0.81</u>	0.50	0.85	0.77	0.70
SR ($\times 16$)	0.58	0.57	0.55	0.59	0.38	0.67	0.60	<u>0.62</u>
Box inpainting	0.80	0.81	<u>0.82</u>	0.76	0.70	0.83	0.80	0.70
Half mask	0.69	0.70	<u>0.71</u>	0.66	0.56	0.73	0.67	0.65
Gaussian Deblur	0.71	0.73	0.75	0.69	0.14	0.77	0.73	<u>0.76</u>
Motion Deblur	<u>0.77</u>	0.78	<u>0.77</u>	0.71	–	–	–	0.71
JPEG (QF = 2)	0.75	<u>0.74</u>	0.73	0.51	12.76	–	–	0.71
Phase retrieval	0.74	<u>0.78</u>	0.79	0.43	–	–	–	0.61
Nonlinear deblur	0.70	0.70	<u>0.69</u>	0.46	–	–	–	0.42
High dynamic range	0.79	<u>0.81</u>	0.86	0.48	–	–	–	0.71

$A = \Lambda - pp^*$ where p is a vector and Λ is a diagonal matrix. This parametrization with facilitates the control of the spectrum of A to enforce stability (negative real part of eigenvalues). We use a multi-scale architecture, with a stack of residual S4 blocks. The top tier processes the raw audio signal at its original sampling rate, while lower tiers process downsampled versions of the input signal. The outputs of the lower tiers are upsampled with up-pooling and down-pooling operations and combined with the input to the upper tier to provide a stronger conditioning signal. Hyper-parameters are described in table 12. The selected model is the model of the last epoch. Training time is 30 hours for 2.5 seconds ECGs, 42 hours for 10 seconds ECGs.

Algorithm parameters We use the same parameters as described in table 6 for MGPS and the implementation described in appendix D.2 for competitors. Since we had to run the experiments over 2k samples, we set the number of diffusion steps such that the runtime approximates 30 seconds for generating 10 samples. This leads to the parameters shown in table 13. For the ML experiment, we also tested MGPS with 300 diffusion steps, which resulted in a runtime of 5 minutes and 30 seconds per batch of 10 samples. The results improved, but we already outperform competitors with just 50 diffusion steps and a 30-second runtime.

Evaluation metrics For the MB task, for each observation, we generate 10 samples - instead of 100 as in in (Alcaraz & Strodthoff, 2022). We then compute the Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE) between each generated sample and the ground-truth. The final score is the average of these errors over the 10 generated samples. We report the confidence intervals over the test set in table 4 and table 15.

Table 8: Mean LPIPS/PSNR/SSIM values for various linear and nonlinear imaging tasks on the IMAGENET 256 × 256 dataset. Best is in **bold** and second best is underlined.

Task	MGPS ₅₀	MGPS ₁₀₀	MGPS ₃₀₀	DPS	PGDM	DDNM	DIFFPIR	REDDIFF
LPIPS ↓								
SR (×4)	0.36	<u>0.33</u>	0.30	0.41	0.78	0.34	0.36	0.56
SR (×16)	0.58	<u>0.55</u>	<u>0.53</u>	0.50	0.60	0.70	0.63	0.83
Box inpainting	0.31	<u>0.26</u>	0.22	0.34	0.29	0.28	0.28	0.36
Half mask	0.39	<u>0.34</u>	0.29	0.44	0.38	0.38	0.35	0.44
Gaussian Deblur	0.35	0.29	<u>0.32</u>	0.35	1.00	0.45	0.29	0.52
Motion Deblur	0.35	<u>0.25</u>	0.22	0.39	–	–	–	0.40
JPEG (QF = 2)	0.50	<u>0.46</u>	0.42	0.63	1.31	–	–	0.51
Phase retrieval	0.54	<u>0.52</u>	0.47	0.62	–	–	–	0.60
Nonlinear deblur	0.49	<u>0.47</u>	0.44	0.88	–	–	–	0.67
High dynamic range	0.22	<u>0.15</u>	0.10	0.85	–	–	–	0.21
PSNR ↑								
SR (×4)	24.68	24.70	<u>24.77</u>	23.52	15.67	25.55	24.26	24.24
SR (×16)	18.56	18.42	18.04	18.22	15.80	20.43	19.37	<u>19.95</u>
Box inpainting	17.52	17.95	18.25	14.34	17.35	20.08	<u>19.77</u>	18.90
Half mask	14.98	15.14	15.60	14.65	14.36	17.06	15.79	<u>16.96</u>
Gaussian Deblur	22.56	21.96	20.10	21.20	9.93	23.29	22.10	<u>23.27</u>
Motion Deblur	23.91	25.03	<u>24.50</u>	21.59	–	–	–	24.43
JPEG (QF = 2)	21.96	<u>22.17</u>	22.44	16.11	5.29	–	–	22.15
Phase retrieval	16.36	<u>16.94</u>	18.10	14.40	–	–	–	15.78
Nonlinear deblur	21.89	<u>22.23</u>	22.36	8.49	–	–	–	20.76
High dynamic range	23.93	<u>25.64</u>	27.04	9.32	–	–	–	22.88
SSIM ↑								
SR (×4)	0.65	<u>0.66</u>	<u>0.66</u>	0.60	0.23	0.70	0.63	0.61
SR (×16)	0.41	0.38	0.35	0.40	0.24	0.50	0.46	<u>0.47</u>
Box inpainting	0.71	0.74	<u>0.76</u>	0.70	0.62	0.78	0.74	0.67
Half mask	0.61	0.63	<u>0.65</u>	0.55	0.53	0.69	0.63	0.62
Gaussian Deblur	0.56	0.53	0.44	0.51	0.07	0.60	0.51	<u>0.57</u>
Motion Deblur	0.64	0.69	<u>0.65</u>	0.55	–	–	–	0.60
JPEG (QF = 2)	0.59	0.60	0.60	0.39	0.01	–	–	0.59
Phase retrieval	0.37	<u>0.40</u>	0.43	0.29	–	–	–	0.26
Nonlinear deblur	<u>0.57</u>	0.58	0.58	0.24	–	–	–	0.41
High dynamic range	0.75	<u>0.81</u>	0.84	0.25	–	–	–	0.73

Table 9: Mean LPIPS/PSNR/SSIM values for various linear and nonlinear imaging tasks on FFHQ 256 × 256 dataset with LDM prior. Best is in **bold** and second best is underlined.

Task	MGPS	RESAMPLE	PSLD	MGPS	RESAMPLE	PSLD	MGPS	RESAMPLE	PSLD
	LPIPS ↓			PSNR ↑			SSIM ↑		
SR (×4)	0.11	<u>0.20</u>	0.22	28.46	26.08	25.53	0.83	0.69	0.70
SR (×16)	0.30	0.36	<u>0.35</u>	20.64	<u>21.09</u>	21.42	<u>0.57</u>	0.56	0.63
Box inpainting	0.16	<u>0.22</u>	0.26	22.94	18.80	<u>20.39</u>	0.79	<u>0.75</u>	0.66
Half mask	0.25	<u>0.30</u>	0.31	15.11	14.59	<u>14.75</u>	0.69	<u>0.67</u>	0.61
Gaussian Deblur	<u>0.16</u>	0.15	0.35	27.57	<u>27.44</u>	19.95	0.79	<u>0.75</u>	0.47
Motion Deblur	0.18	<u>0.19</u>	0.41	<u>26.49</u>	26.85	18.14	0.77	<u>0.72</u>	0.39
JPEG (QF = 2)	0.20	0.26	–	24.75	24.33	–	0.72	0.67	–
Phase retrieval	0.34	0.41	–	22.21	19.05	–	0.62	0.47	–
Nonlinear deblur	0.26	0.30	–	23.79	24.44	–	0.70	0.68	–
High dynamic range	0.15	0.15	–	25.17	25.42	–	0.79	0.81	–

For the ML task, we use a classifier trained to detect four cardiac conditions: Right Bundle Branch Block (RBBB), Left Bundle Branch Block (LBBB), Atrial Fibrillation (AF), and Sinus Bradycardia (SB) on the PTB-XL dataset. We follow the XRESNET1D50 described in (Strodthoff et al., 2020) with hyper-parameters reported in table 14. We apply the classifier to the ground-truth ECG of the test set and to the samples generated from lead I. As in the MB task, for each observation, 10 samples are generated. The output of the classifier is averaged over these 10 samples. For each cardiac condition, we compute balanced accuracy to account for class imbalance (see table 11). The classification threshold is selected using PTB-XL validation-set.

Table 10: ECG size and input-shape per task.

Task	ECG size (seconds)	Total leads	Used leads	Input shape
MB	2,56	12	I, II, III, V1–6	256×9
ML	10	12	I, II, III, V1–6	1024×9

Table 11: PTB-XL dataset description.

Split	All	RBBB	LBBB	AF	SB
Train	17,403	432	428	1211	503
Val	2,183	55	54	151	64
Test	2,203	54	54	152	64

D.5.2 ADDITIONAL RESULTS

Discussion In table 4, we demonstrated that most posterior sampling algorithms outperform the trained diffusion model for the missing block reconstruction task. This result is particularly interesting as it suggests that training a diffusion model (which takes several days) is not necessary for this task. However, when (Alcaraz & Strodthoff, 2023) trained the model on the reconstruction of random missing blocks (RMB), where 50% of each lead is independently removed, the model outperformed all posterior sampling algorithms on the MB task. We report in table 15 the results of the top two algorithms, as well as the model trained on the MB task and the RMB task. The significant improvement between RMB and MB can be seen as an enhancement due to data augmentation.

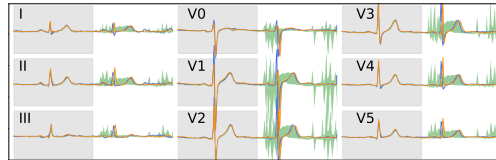


Figure 8: Missing block imputation with MGPS on 2.56s 12-lead ECG. Ground-truth in blue, 10%–90% quantile range in green, random sample in orange.

Generated samples

D.6 SAMPLE IMAGES

While it may appear that some of the methods underperform on some tasks/images compared to the original publications, for instance Figure 14 and Figure 16, they still produce competitive reconstructions on others; see for example Figure 15, Figure 18, and Figure 23. Similar patterns are also displayed in Zhang et al. (2023, Figure 10) and Liu et al. (2023, Figure 9). With that being said, we highlight that these discrepancies appear more frequently on the ImageNet dataset than the FFHQ one. This can be explained by the fact that ImageNet is notoriously challenging due to its diversity, encompassing 1000 classes. This also seems to happen on one of the most difficult tasks, namely the half mask one.

Table 12: ECG diffusion generative model hyper-parameters.

Hyper-parameter	Value
Residual layers	4
Pooling factor	[1, 2, 2]
Feature expansion	2
Diffusion embedding dim. 1	128
Diffusion embedding dim. 2	512
Diffusion embedding dim. 3	512
Diffusion steps	1000
Optimizer	Adam
Number of iterations	150k
Loss function	MSE
Learning rate	0.002
Batch size	128
Number of parameters	16 millions

Table 13: Number of diffusion steps used in posterior sampling algorithms for ECG tasks.

MGPS ₅₀	MGPS ₃₀₀	DPS	PGDM	DDNM	DIFFPIR	REDDIFF
50	300	400	200	200	500	400

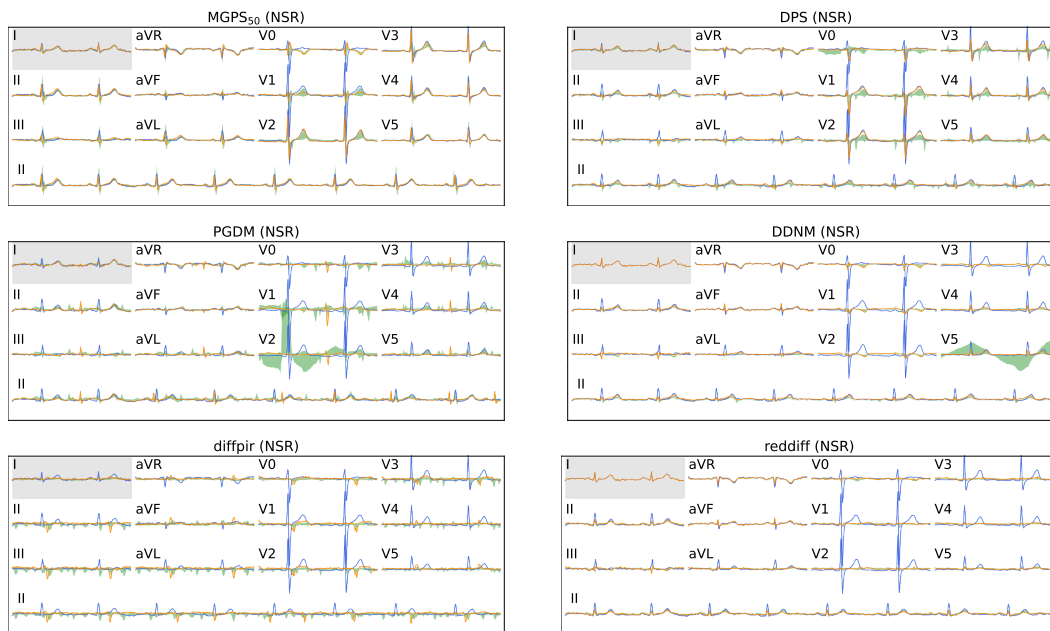


Figure 9: Missing lead reconstruction from lead I on 10s 12-lead Normal Sinus Rhythm (NSR) ECGs. Ground-truth in blue, 10%–90% quantile range in green, random sample in orange.

Table 14: XRESNET1D50 downstream classifier hyper-parameters.

Hyper-parameter	Value
Blocks	4
Layers per block	[3, 4, 6, 3]
Expansion	4
Stride	1
Optimizer	Adam
Learning rate	0.001
Batch size	0.001
Epochs	100

Table 15: MAE and RMSE for missing block task on the PTB-XL dataset.

Metric	MGPS	DDNM	TRAINEDDIFF-MB	TRAINEDDIFF-RMB
MAE	$0.111 \pm 2e-3$	$0.103 \pm 2e-3$	$0.116 \pm 2e-3$	$0.0879 \pm 2e-3$
RMSE	$0.225 \pm 4e-3$	$0.224 \pm 4e-3$	$0.266 \pm 3e-3$	$0.217 \pm 6e-3$

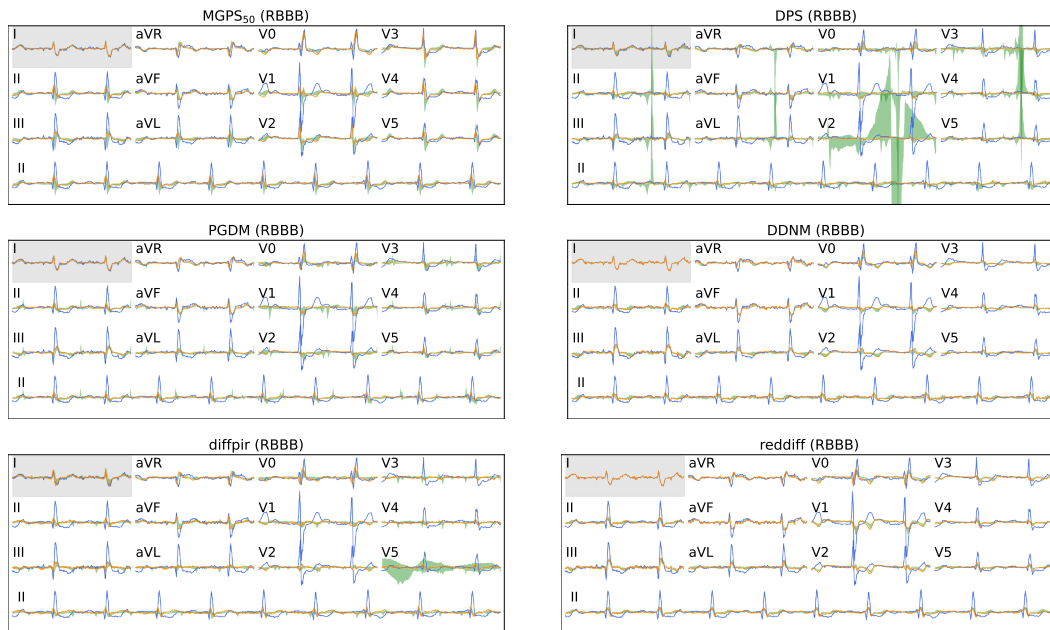


Figure 10: Missing lead reconstruction from lead I on 10s 12-lead Right Bundle Branch Block (RBBB) ECGs. Ground-truth in blue, 10%–90% quantile range in green, random sample in orange.

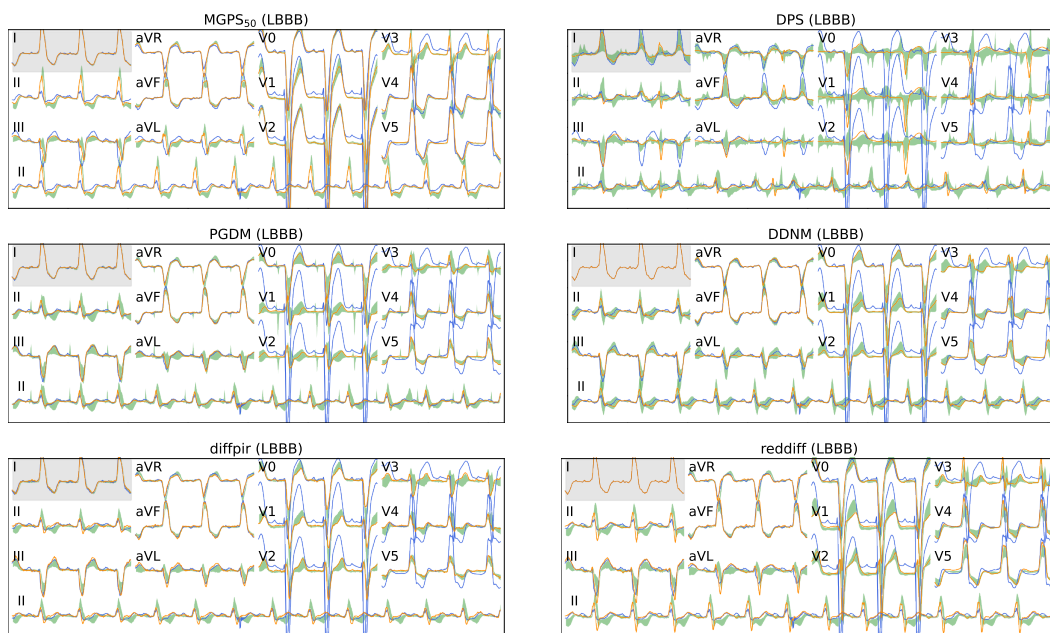
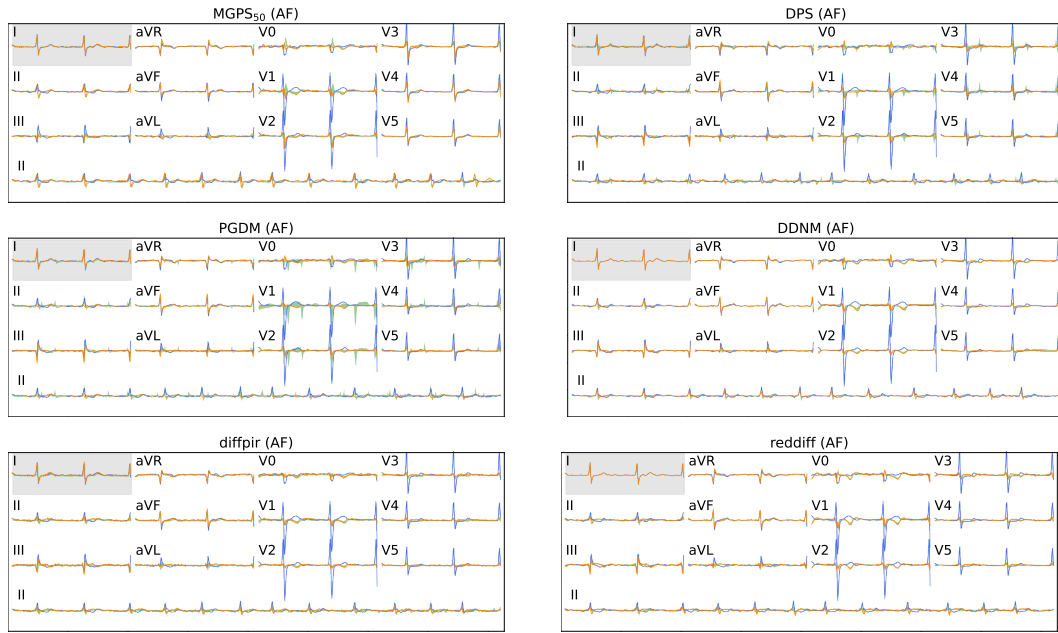


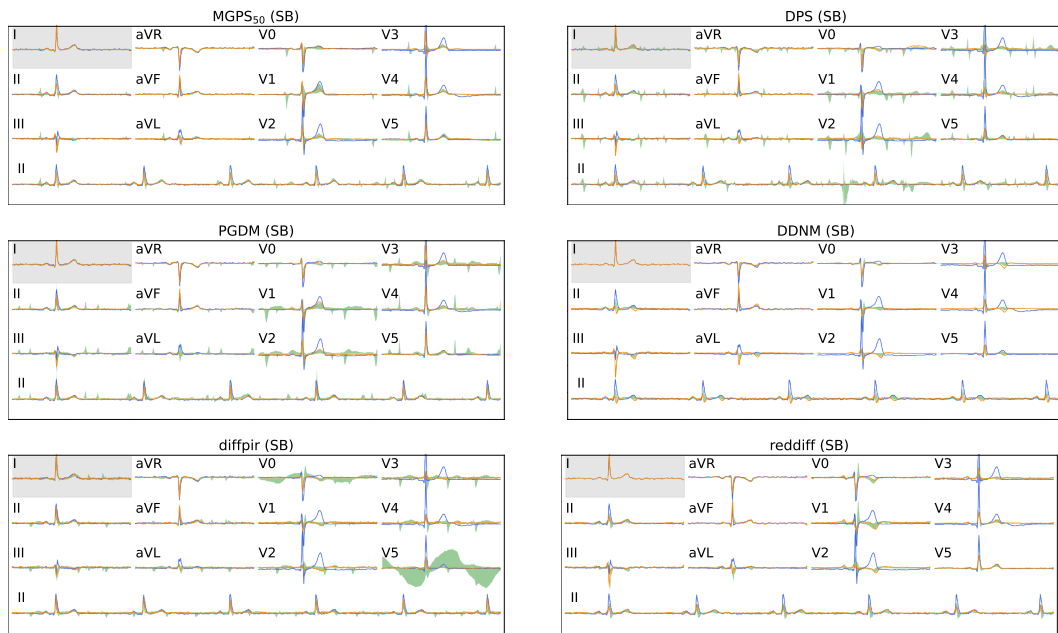
Figure 11: Missing lead reconstruction from lead I on 10s 12-lead Left Bundle Branch Block (LBBB) ECGs. Ground-truth in blue, 10%–90% quantile range in green, random sample in orange.

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588



1589 **Figure 12: Missing lead reconstruction from lead I on 10s 12-lead Atrial Fibrillation (AF) ECGs. Ground-truth in blue, 10%–90% quantile range in green, random sample in orange.**

1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615



1616 **Figure 13: Missing lead reconstruction from lead I on 10s 12-lead Sinus Bradycardia (SB) ECGs. Ground-truth in blue, 10%–90% quantile range in green, random sample in orange.**

1618
1619

1620
 1621
 1622
 1623
 1624
 1625
 1626
 1627
 1628
 1629
 1630
 1631
 1632
 1633
 1634
 1635
 1636
 1637
 1638
 1639
 1640
 1641
 1642
 1643
 1644
 1645
 1646
 1647
 1648
 1649
 1650
 1651
 1652
 1653
 1654
 1655
 1656
 1657
 1658
 1659
 1660
 1661
 1662
 1663
 1664
 1665
 1666
 1667
 1668
 1669
 1670
 1671
 1672
 1673

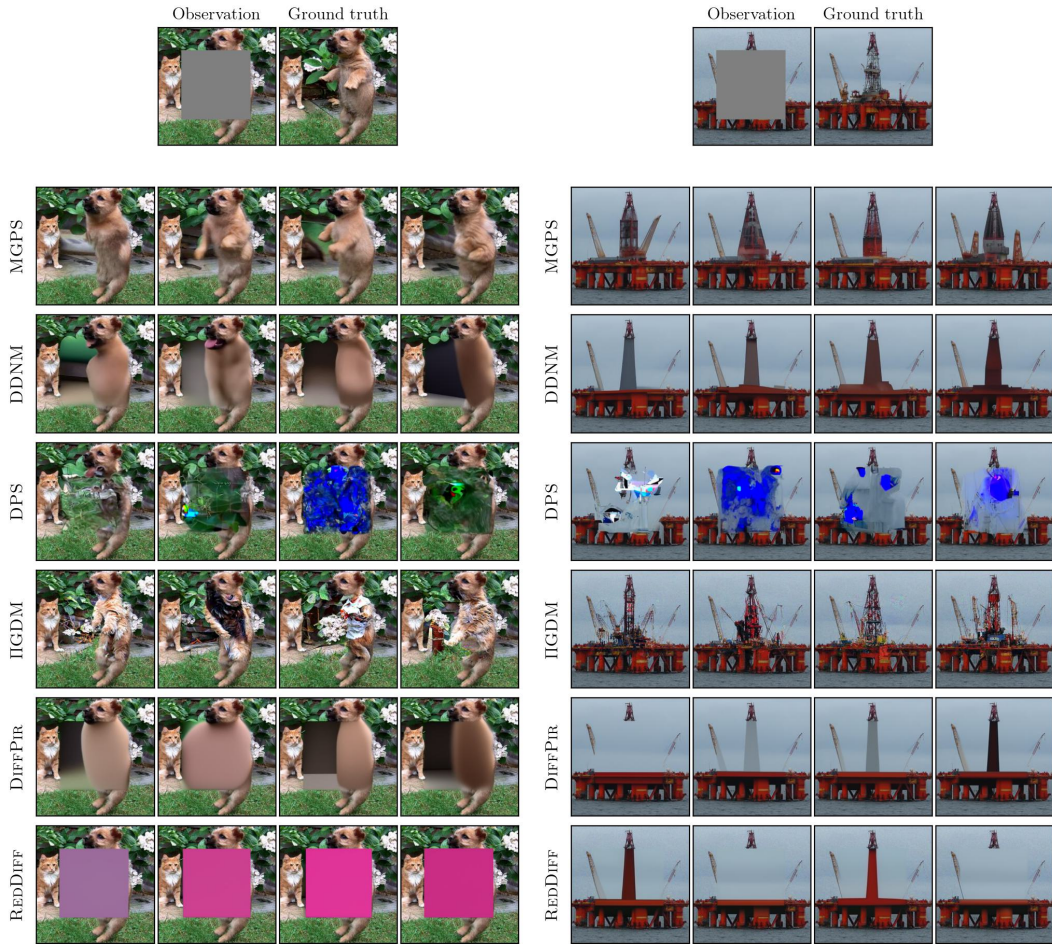


Figure 14: Sample reconstructions with box mask on ImageNet dataset.



Figure 15: Sample reconstructions with half mask on FFHQ dataset.

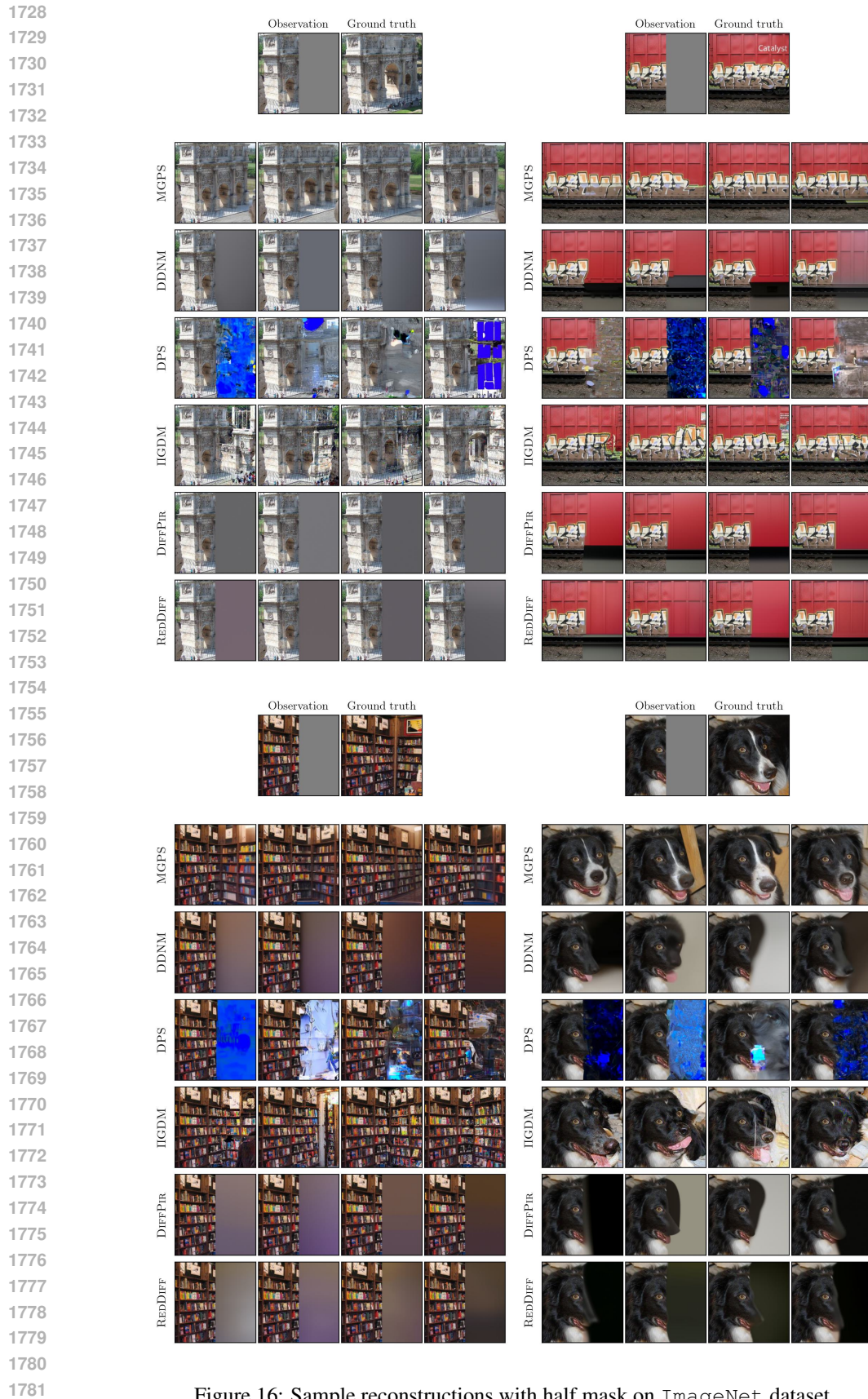


Figure 16: Sample reconstructions with half mask on ImageNet dataset.

1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835

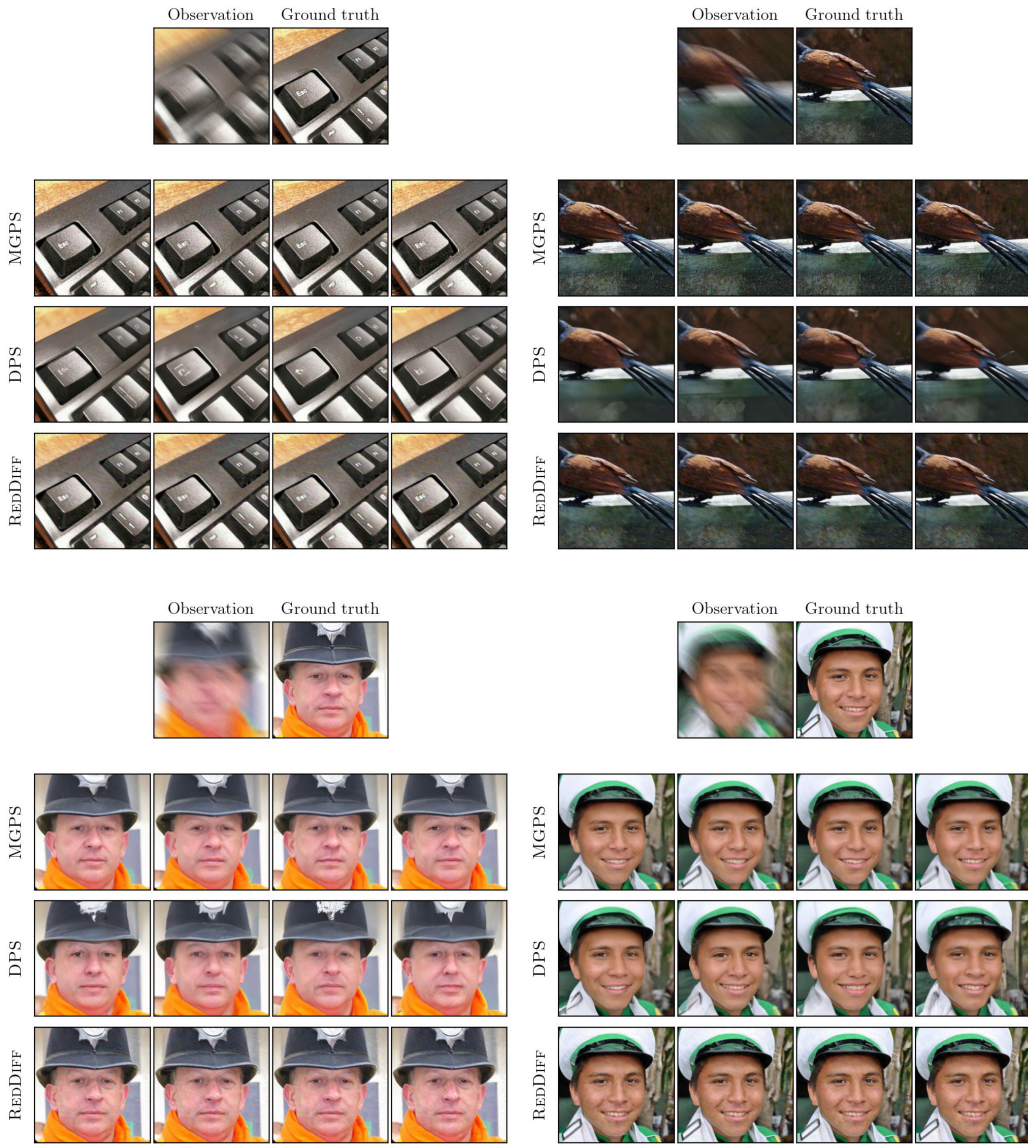


Figure 17: Sample reconstructions on motion deblurring task.

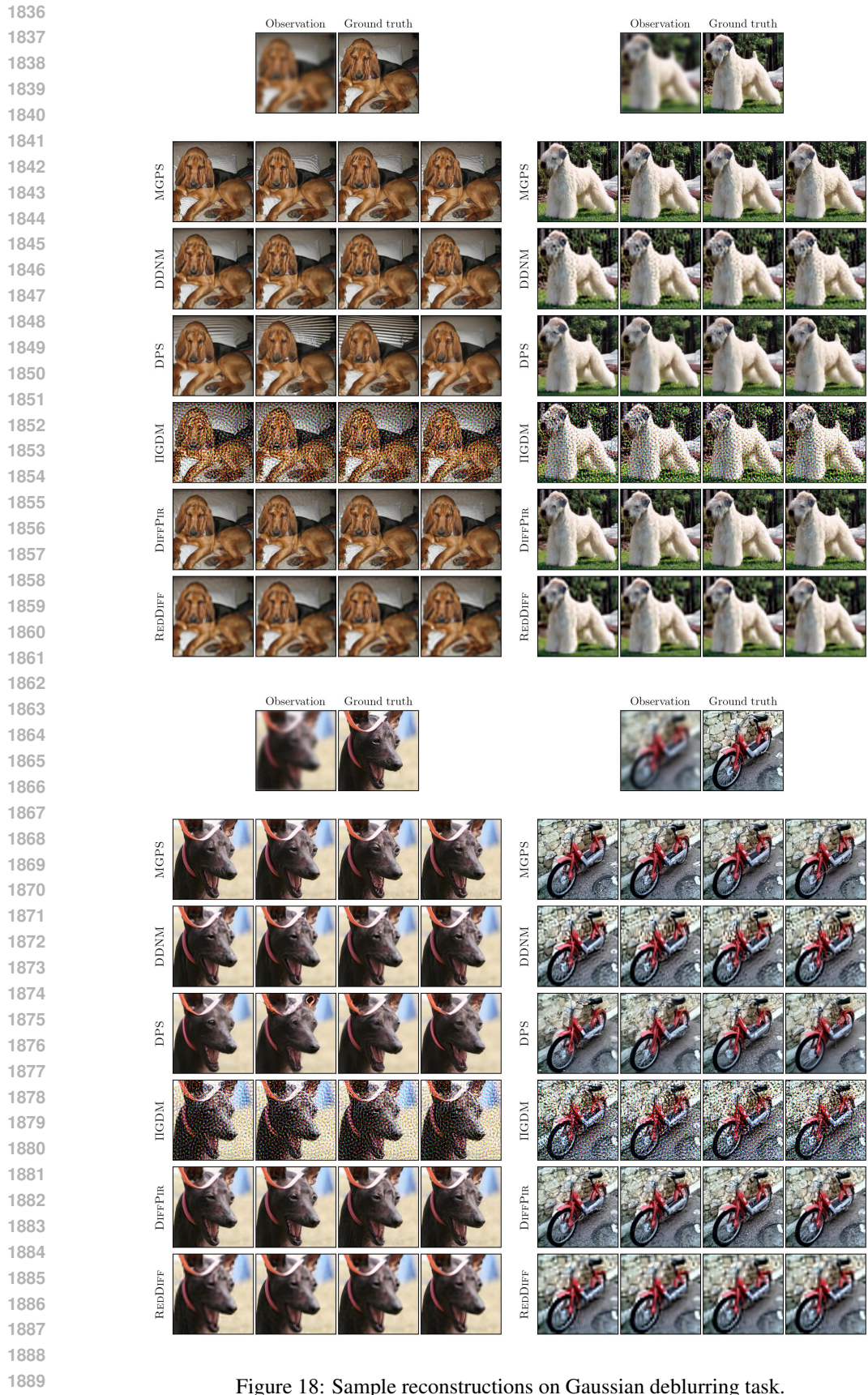


Figure 18: Sample reconstructions on Gaussian deblurring task.

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



Figure 19: Sample reconstructions on phase retrieval task.

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

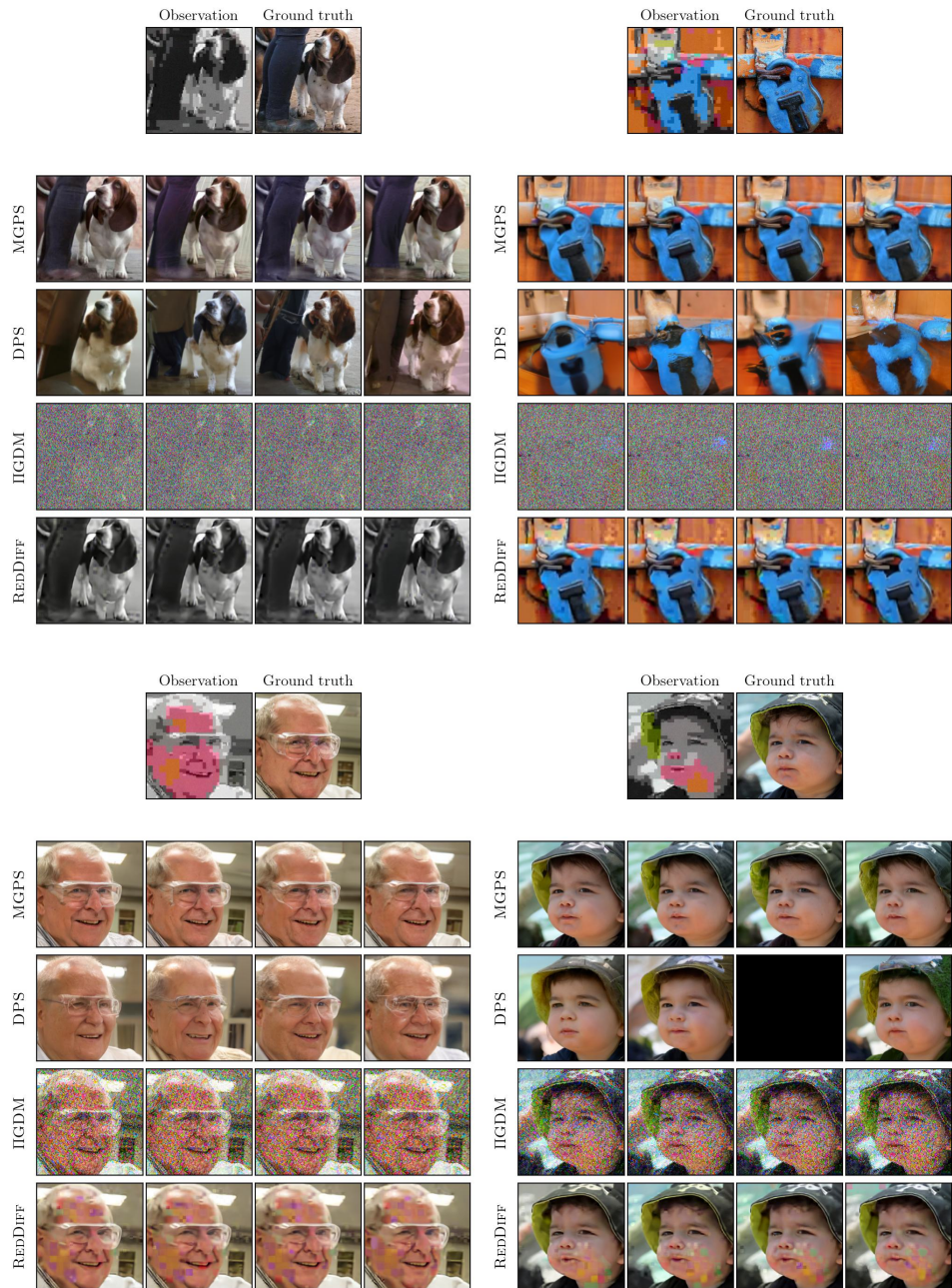


Figure 20: Sample reconstructions on JPEG 2 task.

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

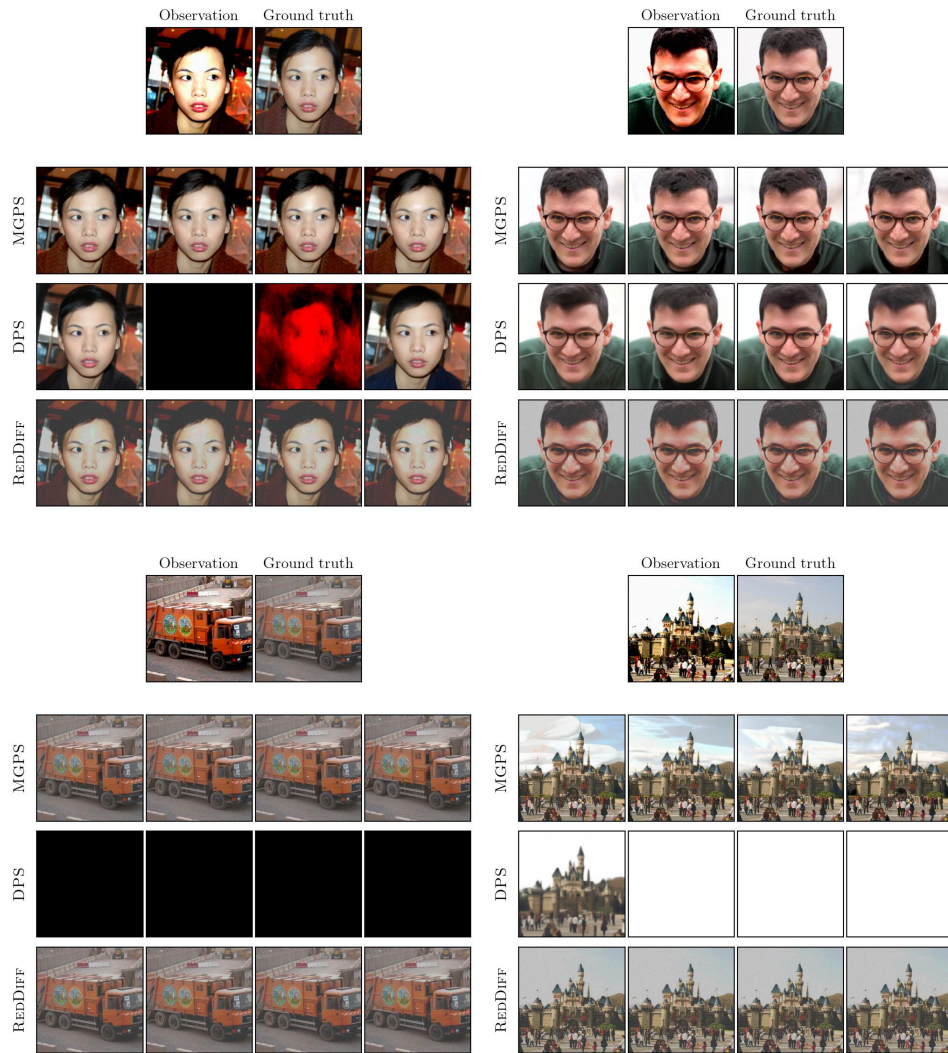
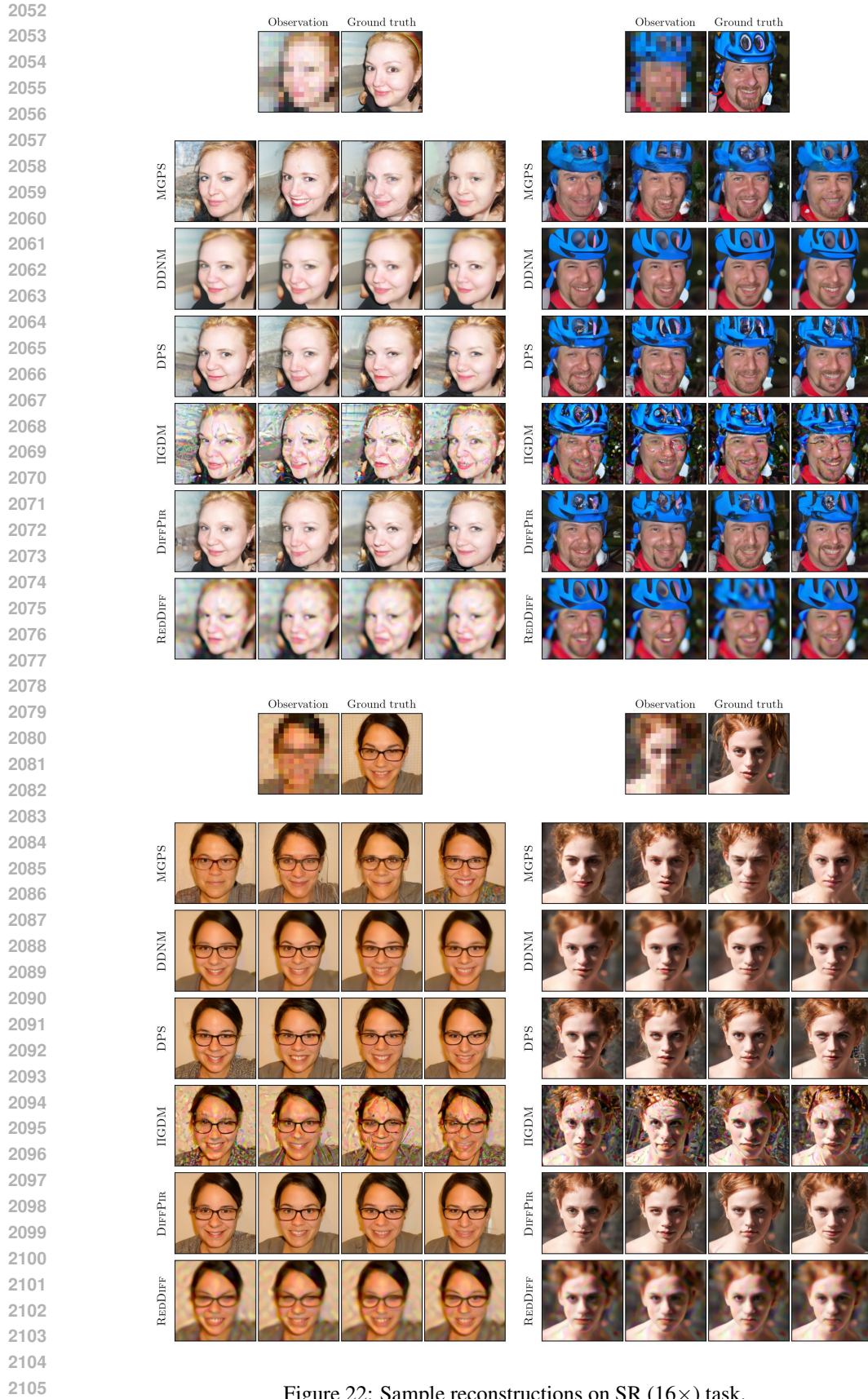


Figure 21: Sample reconstructions on high dynamic range task.



2106
 2107
 2108
 2109
 2110
 2111
 2112
 2113
 2114
 2115
 2116
 2117
 2118
 2119
 2120
 2121
 2122
 2123
 2124
 2125
 2126
 2127
 2128
 2129
 2130
 2131
 2132
 2133
 2134
 2135
 2136
 2137
 2138
 2139
 2140
 2141
 2142
 2143
 2144
 2145
 2146
 2147
 2148
 2149
 2150
 2151
 2152
 2153
 2154
 2155
 2156
 2157
 2158
 2159

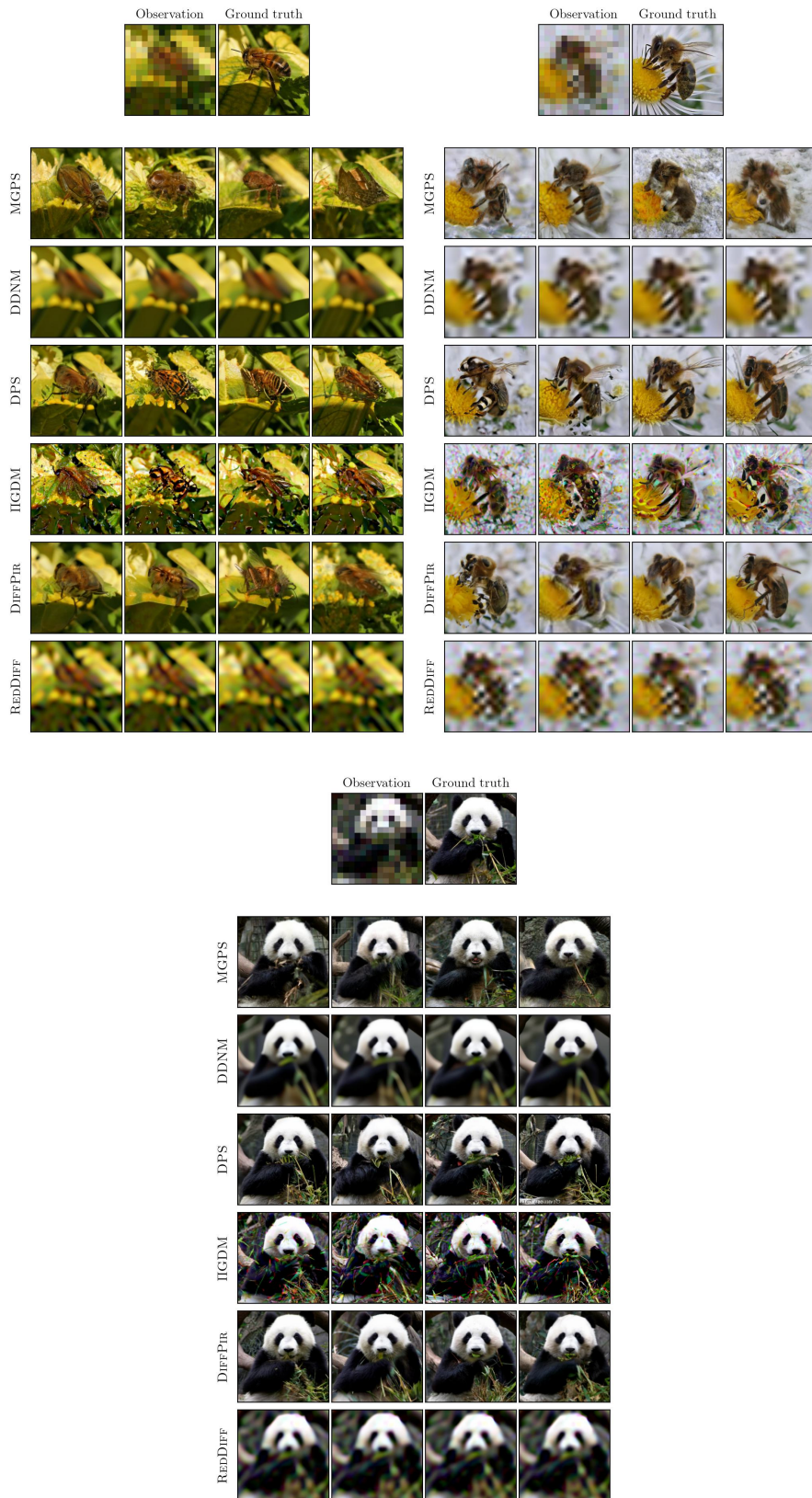
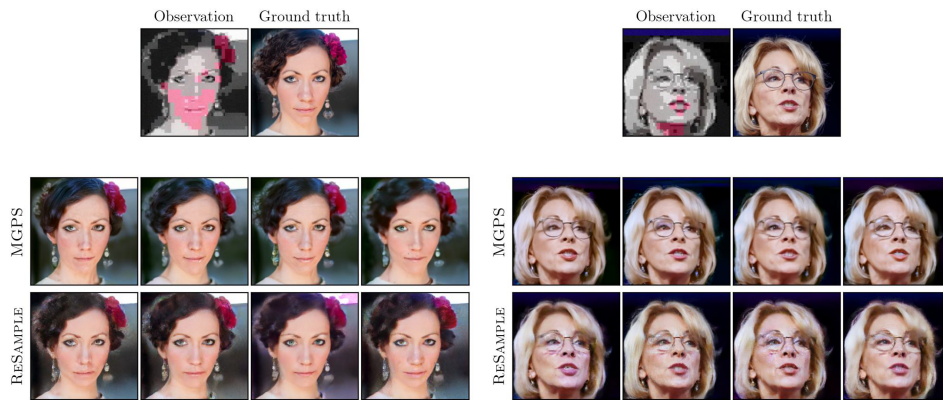
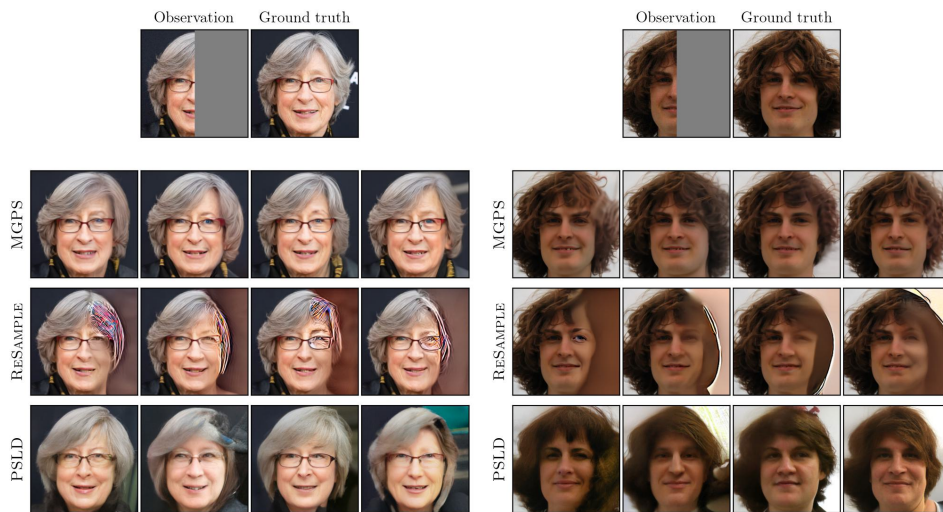


Figure 23: More sample reconstructions on SR (16 \times) task.

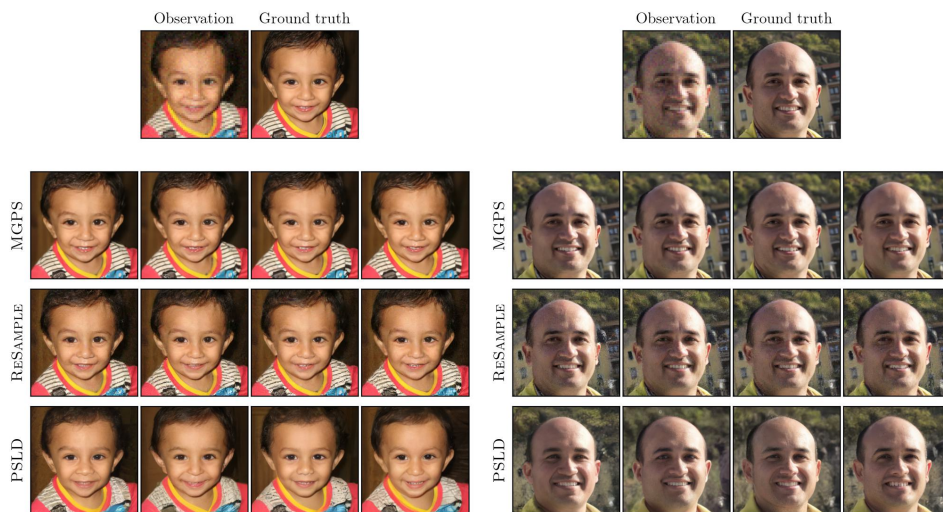
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213



(a) JPEG 2



(b) Outpainting with half mask



(c) SR 4x

Figure 24: Sample reconstructions with latent diffusion models on FFHQ dataset.

2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267

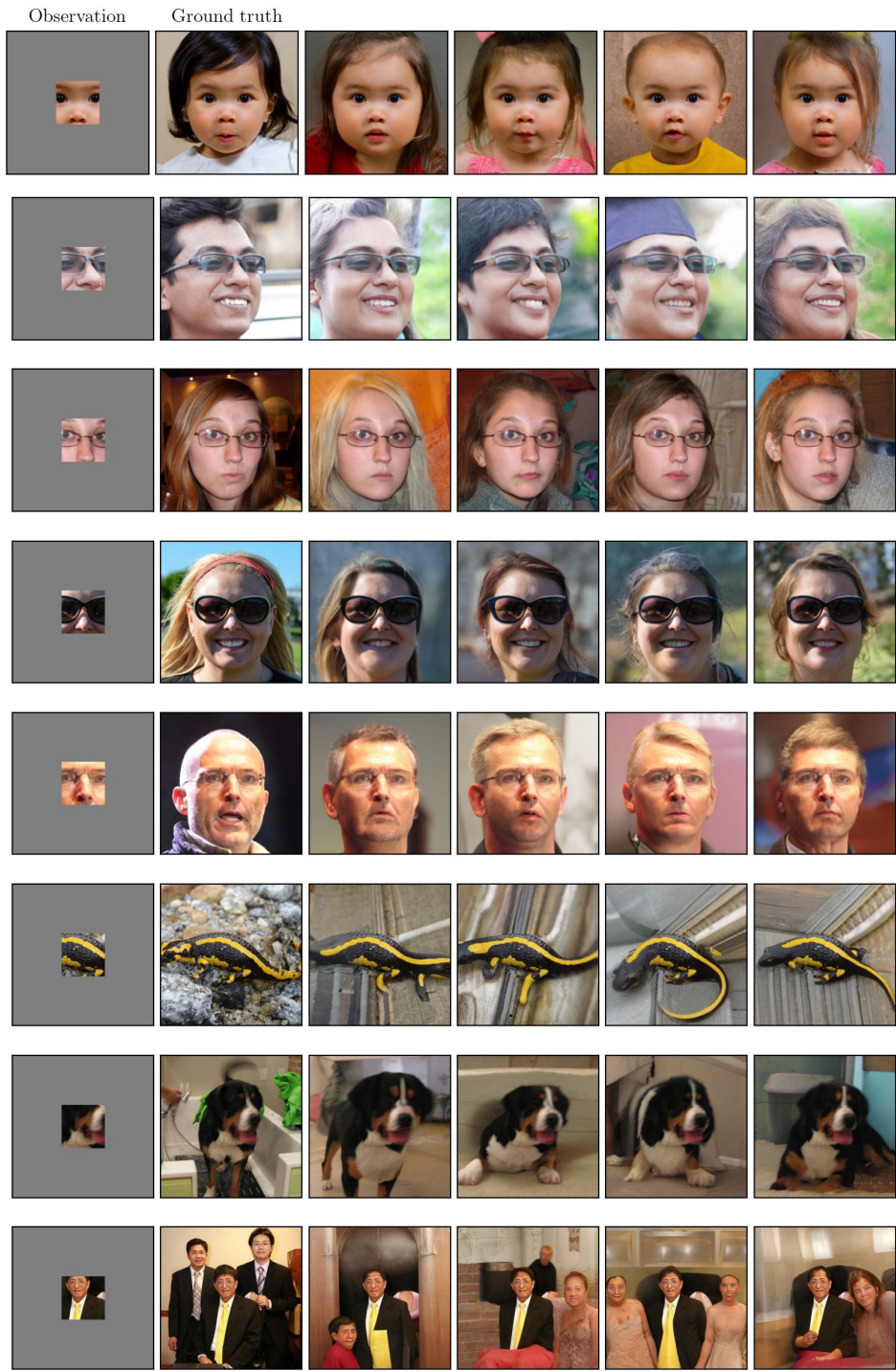


Figure 25: More sample reconstructions with MGPS on the expand task.

2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321

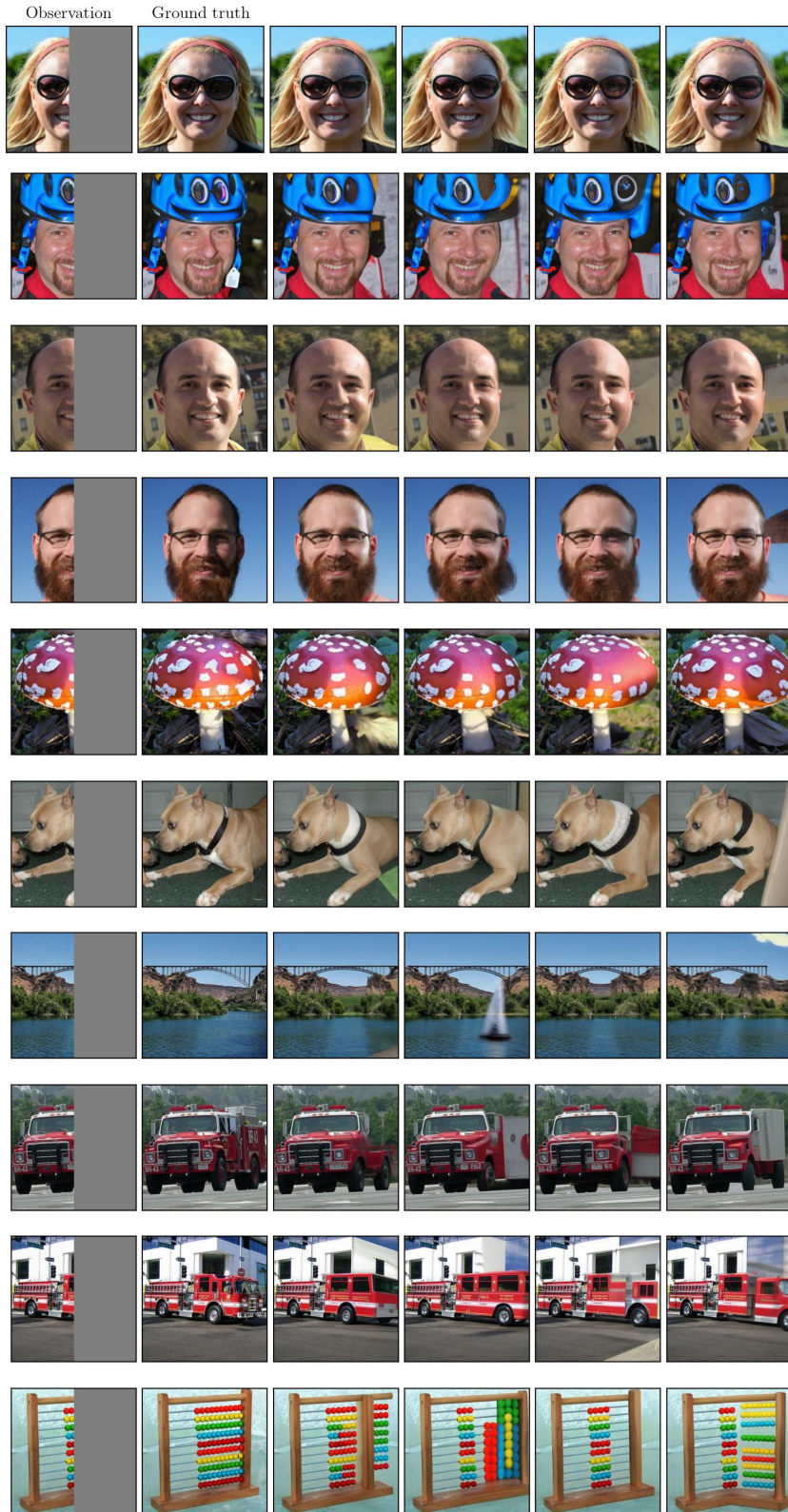


Figure 26: More sample reconstructions with MGPS and half mask.

2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375

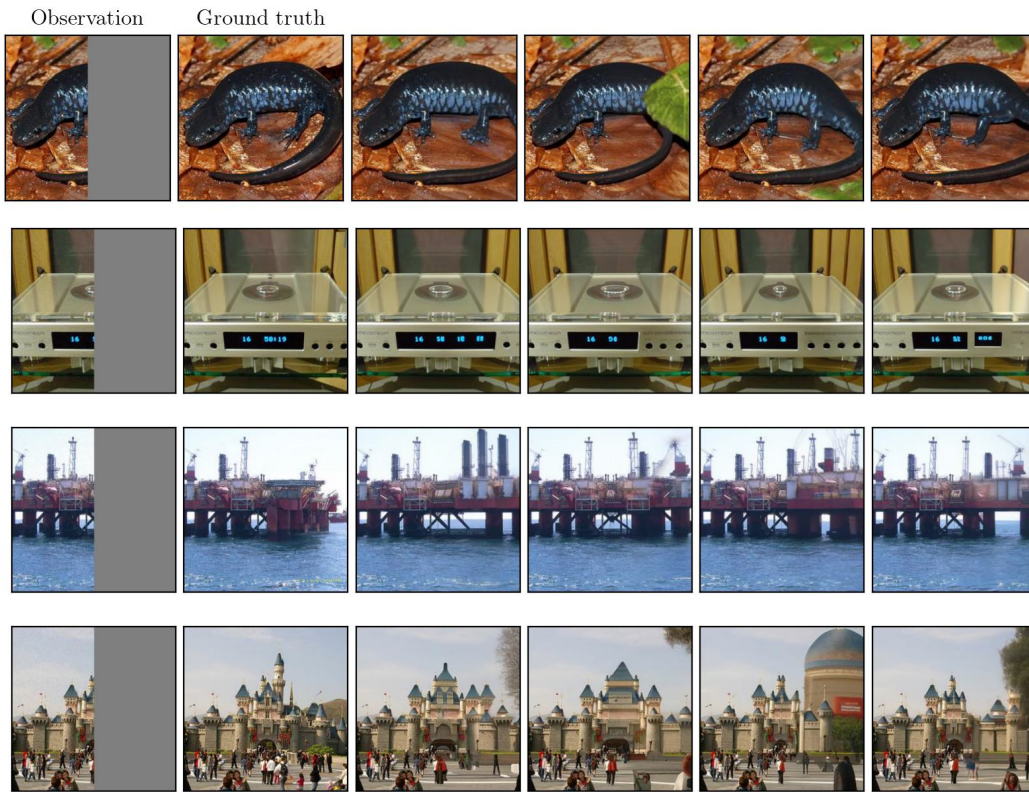


Figure 27: More sample reconstructions with MGPS and half mask.