

# REPULSIVE LATENT SCORE DISTILLATION FOR SOLVING INVERSE PROBLEMS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Score Distillation Sampling (SDS) has been pivotal for leveraging pre-trained diffusion models in downstream tasks such as inverse problems, but it faces two major challenges: (i) mode collapse and (ii) latent space inversion, which become more pronounced in high-dimensional data. To address mode collapse, we introduce a novel variational framework for posterior sampling. Utilizing the Wasserstein gradient flow interpretation of SDS, we propose a multimodal variational approximation with a *repulsion* mechanism that promotes diversity among particles by penalizing pairwise kernel-based similarity. This repulsion acts as a simple regularizer, encouraging a more diverse set of solutions. To mitigate latent space ambiguity, we extend this framework with an *augmented* variational distribution that disentangles the latent and data. This repulsive augmented formulation balances computational efficiency, quality, and diversity. Extensive experiments on linear and nonlinear inverse tasks with high-resolution images ( $512 \times 512$ ) using pre-trained Stable Diffusion models demonstrate the effectiveness of our approach.

## 1 INTRODUCTION

Diffusion models have recently achieved remarkable success in visual domains. A key application of these models is solving various inverse problems in a *plug-and-play* manner, where diffusion models act as *rich priors* to regularize the search space, ensuring the generation of plausible solutions. Variational samplers (Poole et al., 2022; Mardani et al., 2024) approach sampling as an optimization problem, providing a high degree of control and fidelity in generation. However, they encounter two significant challenges, particularly when dealing with high-dimensional data that requires diverse outputs: (c1) *mode collapse*, and (c2) *inversion of the latent space*, such as that seen in the adversarial autoencoder of Stable Diffusion (Rombach et al., 2021).

There have been a few recent attempts to address these challenges separately in the context of text-to-image/3D generation and inverse problems. To mitigate (c1), for text-to-3D generation, ProlificDreamer (VSD) (Wang et al., 2024) introduces data-driven dispersion with independent particles. Still, the combination of independence and unimodal approximation per particle renders an optimization that collapses to the same local minimum, limiting diversity. Collaborative Score Distillation (CSD) (Kim et al., 2023) seeks to diversify the variational approximation using Stein Variational Gradient Descent (SVGD) (Liu and Wang, 2016), but smoothing particle gradients with SVGD is problematic in high-dimensional spaces (D’Angelo and Fortuin, 2021a; Ba et al., 2021). To address (c2), recent samplers using latent diffusion models (Rout et al., 2024; Chung et al., 2024; Song et al., 2024) remain computationally demanding, similar to earlier pixel-based methods (Chung et al., 2022a; Song et al., 2022), due to multiple correction steps required for deviations from the image manifold, a challenge arising from adversarial training of autoencoders, akin to GAN inversion (Xia et al., 2022; Daras et al., 2021). Thus, no current solution effectively handles both mode collapse and latent space issues in inverse problems.

We hypothesize that the primary issue with (c1) arises from collapse in high-dimensional spaces. To address this, we employ an ensemble of interactive particles with repulsion to prevent collapse. Inspired by kernel-density estimation (D’Angelo and Fortuin, 2021b), we introduce a particle-based multimodal variational approximation that incorporates repulsive forces. These forces are defined through pairwise interactions based on similarity, such as using a radial basis kernel of DINO features (Caron et al., 2021).

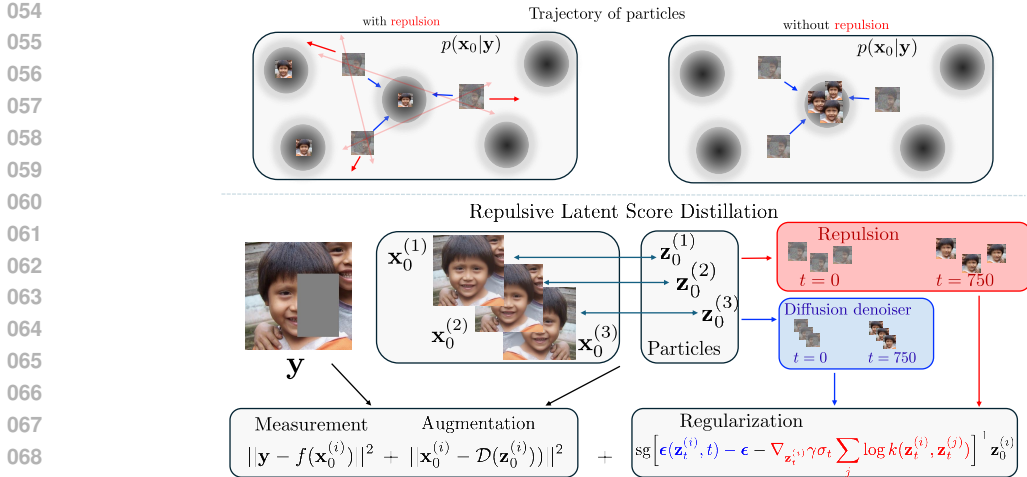


Figure 1: Illustration of Repulsive Latent Score Distillation (RLSD): It propagates a set of particles by adding noise and applying two levels of regularization: (i) **Denoising**, via score-matching regularization, which directs particles toward modes of the distribution  $p(\mathbf{x}_0|\mathbf{y})$  (blue arrows); and (ii) **Repulsion**, which pushes particles apart (red arrows) to explore other regions of the posterior density. During sampling, the repulsion gradient ensures particles remain separated, leading to different modes, as shown in the upper-right box.

To address (c2), we propose a variational augmented distribution that jointly optimizes the latent and data variables, similar to half-quadratic splitting (Geman and Yang, 1995). This method disentangles the latent (prior) from the data (measurement), yielding solutions with sharper details. We show that KL minimization of our augmented interactive particle approximation leads to score-matching regularization with two gradient terms: (a) denoising regularization along the entire diffused trajectory, and (b) repulsion regularization to encourage diversity in the latent diffusion’s trajectory; see Fig. 1. We refer to our method as *Repulsive Latent Score Distillation (RLSD)*.

We validate the advantages of RLSD through extensive experiments on both linear and nonlinear tasks, using Stable Diffusion as the prior. In diversity-critical cases such as inpainting and phase retrieval, our method provides a solid trade-off between diversity and quality. For tasks where diversity is less essential, like deblurring, our augmented formulation offers a fast solver by avoiding score Jacobian computations and performs efficiently on high-resolution ( $512 \times 512$ ) images. Overall, RLSD combines the strengths of variational samplers (memory and compute efficiency) and posterior samplers (diversity), enabling control over speed and diversity by adjusting the scalar weights between denoising and repulsion regularizations. A detailed comparison of RLSD’s properties is summarized in Table 1.

All in all, the main contributions of this paper are summarized as follows:

- We propose **Repulsive Latent Score Distillation (RLSD)**, a variational posterior sampler for general inverse problems with high-resolution images (e.g.,  $512 \times 512$ ), that *trades-off diversity for quality* in a controllable fashion simply via regularization weights.
- We introduce a *repulsion regularization* to boost the diversity via an interactive particle-based variational approximation inspired by Wasserstein gradient flow.
- To handle the latent space inversion, we propose a *distribution augmentation* that decouples the latent and pixel space, rendering a two-step optimization problem.
- We perform extensive experiments for various (non)linear inverse tasks using Stable Diffusion. The results indicate the superior performance of RLSD over existing alternatives such as PLSD (Rout et al., 2024), DPS (Chung et al., 2022a) and RED-Diff (Mardani et al., 2024).

## 2 RELATED WORKS

This paper is primarily related to diffusion models at its core, and two related lines of work: inverse problems and score distillation sampling.



	Posterior samplers		Variational samplers			
	DPS	PSLD	RED-Diff	VSD	CSD	RLSD (Ours)
Diversity (Low-dim)	✓	✓	✗	✓	✓	✓
Diversity (High-dim)	✓	✓	✗	✗	✗	✓
Latent Diffusion	✗	✓	✗	✓	✓	✓
Linear Inv. Problems	✓	✓	✓	-	-	✓
Nonlinear Inv. Problems	✓	-	✓	-	-	✓
No Score Jacobian	✗	✗	✓	✓	✓	✓

Table 1: Comparison of our work with DPS, RED-Diff, PSLD, VSD, and CSD. Our method combines the strengths of variational samplers (no score Jacobian, computational efficiency) and posterior sampling algorithms (diversity at high dimensions). In addition, our formulation enables us to solve nonlinear inverse problems at  $512 \times 512$ .

**Diffusion models for inverse problems.** Several works have used diffusion models as priors to solve inverse problems in various domains (Daras et al., 2024; Kong et al., 2020). A recent approach termed RED-Diff (Mardani et al., 2024) uses variational inference for solving inverse problems with diffusion priors, similar to plug-and-play methods (Venkatakrisnan et al., 2013); see also Zhu et al. (2023); Zhang et al. (2021). This method employs the diffusion model as denoisers at different scales, akin to the RED framework (Romano et al., 2017). Despite successfully balancing quality and runtime, it suffers from mode collapse due to the unimodal approximation. Furthermore, optimizing directly in the pixel domain restricts the ability to leverage latent diffusion models like Stable Diffusion (Rombach et al., 2021) for solving inverse problems at high-resolution. Recent works incorporate latent diffusion models as prior (Rout et al., 2024; Chung et al., 2024; Song et al., 2024). While they partially alleviate the computational demands of pixel-domain solvers, they introduce additional steps to correct the deviations from the image manifold, which arise from the adversarial training of the autoencoder. Thus, it is still an open problem to develop methods that are fast, promote diversity, and optimize in the latent space of the diffusion model.

**Score distillation: diversity and mode collapse.** Recently, SDS enabled the use of pretrained diffusion models for text-to-3D generation (Poole et al., 2022). Although SDS provides an efficient mechanism for the aforementioned task, it often suffers from mode collapse and saturated images; see details about SDS and its formulation in Appendix E. ProlificDreamer (Wang et al., 2024) aims to fix the mode collapse using a data-driven dispersion fine-tuned at each iteration via LoRA (Hu et al., 2021). It is, however, costly, and the independence of particles hinders diversity. Recently, the authors in (Kim et al., 2023) propose to use the well-known Stein variational gradient descent (SVGD) as an update direction, which yields an interactive particle system. However, it is known that SVGD suffers from the curse of dimensionality (D’Angelo and Fortuin, 2021a). Other related works are Armandpour et al. (2023), where the authors leverage the negative prompt to eliminate undesired perspectives, and Wang et al. (2023), where an entropic regularization is proposed.

### 3 BACKGROUND

We review latent diffusion models in Section 3.1, and we briefly discuss how they are incorporated as priors to solve inverse problems in Section 3.2.

#### 3.1 DIFFUSION MODELS IN THE LATENT SPACE

Diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song et al., 2021b) consist of two processes modeled using stochastic differential equations: 1) a forward process that gradually adds noise to a clean image, and 2) a reverse process that learns to generate images by iteratively denoising the diffused data. In latent diffusion models (Vahdat et al., 2021; Rombach et al., 2021), the data  $\mathbf{x}_0$  is encoded into a latent space through an *encoder*  $\mathcal{E}(\mathbf{x}_0) = \mathbf{z}_0$ , and the forward process follows the variance-preserving SDE (Song et al., 2021b) in the latent space:  $d\mathbf{z}_t = -\frac{1}{2}\beta(t)\mathbf{z}_tdt + \sqrt{\beta(t)}d\mathbf{W}_t$ , for  $t \in [0, T]$ . Here,  $\beta(t)$  is a function that defines a step size for each  $t$  from 0 to  $T$ , and is defined as  $\beta(t) := \beta_{\min} + (\beta_{\max} - \beta_{\min})\frac{t}{T}$ , and  $\mathbf{W}_t$  is the standard Brownian motion. The forward

process is designed in such a way that the distribution of  $\mathbf{z}_T$  converges to a standard Gaussian distribution. Given the forward process, the reverse process is defined as  $d\mathbf{z}_t = -\frac{1}{2}\beta(t)\mathbf{z}_t dt - \beta(t)\nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) + \sqrt{\beta(t)}d\mathbf{W}_t$ , where  $\nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t)$  is the *score function*, which is unknown. To map back to the ambient space, we pass the generated sample  $\mathbf{z}_0$  through a *decoder*  $\mathcal{D}(\mathbf{z}_0) = \mathbf{x}_0$ .

Therefore, to solve the reverse process and use it for sampling, the score function ( $\nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t)$ ), encoder ( $\mathcal{E}$ ), and decoder ( $\mathcal{D}$ ) are learned by minimizing the denoising score-matching loss (Vincent, 2011). Diffused samples are generated as  $\mathbf{z}_t = \alpha_t \mathbf{z}_0 + \sigma_t \epsilon$ , where  $\mathbf{z}_0$  encodes  $\mathbf{x}_0 \sim p_{\text{data}}(\mathbf{x})$ , and  $\sigma_t = 1 - e^{-\int_0^t \beta(s) ds}$ , and  $\alpha_t = \sqrt{1 - \sigma_t^2}$ . The score function is approximated by  $\epsilon_{\theta}(\mathbf{z}_t, t) \approx -\sigma_t \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t)$ , and the score-matching loss is minimized. After training, samples are generated using samplers like DDPM (Ho et al., 2020) and DDIM (Song et al., 2020).

### 3.2 INVERSE PROBLEMS WITH DIFFUSION PRIORS

In general, an *inverse problem* aims to find an unknown signal  $\mathbf{x}_0$  given some noisy measurement  $\mathbf{y}$ , related via some forward model  $f(\cdot)$ ,

$$\mathbf{y} = f(\mathbf{x}_0) + \mathbf{v}, \quad \mathbf{v} \sim \mathcal{N}(0, \sigma_v^2 \mathbf{I}), \quad (1)$$

where the forward model is domain-dependent. In a Bayesian framework, the solution boils down to sample from the posterior  $p(\mathbf{x}_0|\mathbf{y}) \propto p(\mathbf{y}|\mathbf{x}_0)p(\mathbf{x}_0)$ , where  $p(\mathbf{y}|\mathbf{x}_0)$  is the measurement model (1) and  $p(\mathbf{x}_0)$  is the prior imposed by the diffusion model.

**Diffusion posterior sampling approaches.** These methods generate a sample from the posterior by running the reverse process (see Section 3.1) using conditional score at  $t$  obtained via Bayes' rule as

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}) = \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t). \quad (2)$$

While the second term uses a pre-trained diffusion model, the first is intractable, as seen from  $p(\mathbf{y}|\mathbf{x}_t) = \int p(\mathbf{y}|\mathbf{x}_0)p(\mathbf{x}_0|\mathbf{x}_t)d\mathbf{x}_0$ . Prior works (Chung et al., 2022a; Song et al., 2022; Kadkhodaie and Simoncelli, 2021; Song et al., 2023) address this with a Gaussian approximation of  $p(\mathbf{x}_0|\mathbf{x}_t)$  using Tweedie's formula  $\mathbb{E}[\mathbf{x}_0|\mathbf{x}_t] = \frac{1}{\alpha_t}(\mathbf{x}_t - \sigma_t \epsilon_{\theta}(\mathbf{x}_t, t))$ . Still, this requires the computation of the *score Jacobian*, which is computationally expensive, especially for pixel-based models at high-resolution. This can be partially addressed by using a suitable latent space (Rombach et al., 2021), alleviating the computational demands of pixel-domain solvers for high-resolution images. However, as discussed in the previous section, it introduces additional steps which arise from the adversarial training of the autoencoder (Rout et al., 2024). We defer more details to Appendix B.

**Variational inference approaches.** Recently, RED-diff was introduced in Mardani et al. (2024), which avoids computing the score Jacobian<sup>1</sup>. RED-diff frames the sampling problem as stochastic optimization by minimizing the KL divergence

$$q(\mathbf{z}_0|\mathbf{y}) = \underset{q(\mathbf{z}_0|\mathbf{y})}{\operatorname{argmin}} \operatorname{KL}(q(\mathbf{z}_0|\mathbf{y})||p(\mathbf{z}_0|\mathbf{y})), \quad (3)$$

where  $\mathbf{x}_0 = \mathcal{D}(\mathbf{z}_0)$ . When  $q(\mathbf{z}_0|\mathbf{y}) \sim \mathcal{N}(\boldsymbol{\mu}_z, \sigma_z^2 \mathbf{I})$ , the KL minimization (3) boils down to a maximum-a-posteriori optimization that leverages the diffusion model's trajectory as a regularizer, resulting in a simple and tractable method. However, this approach shares the same limitations as score distillation regarding diversity and mode collapse. Additionally, applying this formulation directly with latent diffusion models produces blurry results (see Appendix D.8.3).

## 4 REPULSIVE VARIATIONAL DIFFUSION SAMPLING

In Section 4.1, we address (c1) by introducing a repulsion mechanism, promoting diversity through a multimodal variational approximation using interactive particles. Then, in Section 4.2, we tackle (c2) by proposing an augmented variational formulation. Finally, in Section 4.3 we combine both techniques to derive **Repulsive Latent Score Distillation** (RLSD), our proposed solver for inverse problems using latent diffusion models.

<sup>1</sup>RED-Diff was proposed in the pixel-domain. However, for the sake of clarity, here we express it with respect to the latent diffusion models.

#### 216 4.1 TACKLING MODE COLLAPSE: ENHANCING DIVERSITY VIA REPULSION

217  
218 We aim to solve inverse problems characterized by the forward model in (1) by minimizing the reverse  
219 KL divergence (3). However, as explained in Section 3.2, minimizing (3) with a Gaussian variational  
220 distribution leads to a unimodal approximation of a multimodal posterior, which is problematic for  
221 highly ill-posed problems like inpainting. To circumvent this, we propose a *particle approximation*  
222 for defining a *multimodal* variational distribution. More precisely, we incorporate a repulsion force  
223 within the particles to encourage the exploration of multiple modes.

224 **Particle interpretation of SDS.** To facilitate the presentation, throughout this section we consider  
225 the unconditional case of (3), i.e., without measurements:

$$226 q(\mathbf{z}_0) = \operatorname{argmin}_{q(\mathbf{z}_0)} \operatorname{KL}(q(\mathbf{z}_0) || p(\mathbf{z}_0)) \quad (4)$$

228 Following Song et al. (2021a), we rewrite (4) in terms of the diffused trajectory as

$$229 q(\mathbf{z}_0) = \operatorname{argmin}_{q(\mathbf{z}_0)} \mathbb{E}_{t \sim \mathcal{U}[0, T]} [\omega(t) \operatorname{KL}(q(\mathbf{z}_t) || p(\mathbf{z}_t))], \quad (5)$$

231 where  $\omega(t)$  is a weighting function and  $q(\mathbf{z}_t) = \int q(\mathbf{z}_t | \mathbf{z}_0) q(\mathbf{z}_0) d\mathbf{z}_0$  depends on the diffused trajectory  
232  $q(\mathbf{z}_t | \mathbf{z}_0) \sim \mathcal{N}(\alpha_t \mathbf{z}_0, \sigma_t^2 \mathbf{I})$  and the *variational approximation*  $q(\mathbf{z}_0)$ . The optimization in (5)  
233 corresponds to score distillation, which can be formulated as a Wasserstein gradient flow (details of  
234 WGF can be found in Appendix E.1). At the particle level, the WGF is described by the following  
235 ODE

$$236 d\mathbf{z}_{0, \tau}^{(i)} = \mathbb{E}_t \left[ \omega(t) \left( \nabla_{\mathbf{z}_{t, \tau}^{(i)}} \log p \left( \mathbf{z}_{t, \tau}^{(i)} \right) - \nabla_{\mathbf{z}_{t, \tau}^{(i)}} \log q_\tau \left( \mathbf{z}_{t, \tau}^{(i)} \right) \right) \right] d\tau, \quad (6)$$

238 where  $p \left( \mathbf{z}_{t, \tau}^{(i)} \right)$  is the target distribution,  $q_\tau \left( \mathbf{z}_{t, \tau}^{(i)} \right)$  is the marginal distribution of a generic particle  $i$   
239 at time-step  $\tau$ , and  $t$  is the noise level of the diffusion model. In a nutshell, the WGF of  $\mathbf{z}_0$  in (6) is  
240 computed as an expectation over its diffused trajectory (noise levels  $t$ ), involving the gradients of  
241  $\mathbf{z}_t$ . This formulation shields light on how the particles are propagated when optimizing (5). More  
242 precisely, it becomes evident that for an *initial Gaussian variational approximation at  $\tau = 0$ , the*  
243 *marginal for all  $\tau$  is also Gaussian*. The dynamic in (6) yields a deterministic trajectory where *the*  
244 *mode of the Gaussian variational approximation will match one of the modes of  $p(\mathbf{z}_0)$* . Consequently,  
245 assuming the same initial position, *all particles will converge to the same mode*. Although this can be  
246 mitigated ad hoc by changing the particles’ initial positions, we seek a more principled method.

247 In this context, we can enhance diversity by considering 1) a multimodal (but most likely intractable)  
248 variational distribution or 2) interactive particle systems; in this work, we focus on the second one  
249 due to its tractability. When considering an interactive set of particles, the key design factor is the  
250 coupling term, which prevents the ensemble from collapsing to the same mode. In particular, we  
251 propose using a repulsion term.

253 **Repulsive variational distribution.** Inspired by D’Angelo and Fortuin (2021b), we consider an  
254 ensemble of interacting particles coupled via a *repulsive force* that pushes particles away from  
255 collapsing to the same solution. In a nutshell, we introduce a repulsive force that yields the following  
256 modification of the gradient flow in (6)

$$257 d\mathbf{z}_{0, \tau}^{(i)} = \mathbb{E}_t \left[ \omega(t) \left( \nabla_{\mathbf{z}_{t, \tau}^{(i)}} \log p \left( \mathbf{z}_{t, \tau}^{(i)} \right) - \nabla_{\mathbf{z}_{t, \tau}^{(i)}} \log q_\tau \left( \mathbf{z}_{t, \tau}^{(i)} \right) - \nabla_{\mathbf{z}_{t, \tau}^{(i)}} \mathcal{R}(\mathbf{z}_{t, \tau}^{(1)}, \dots, \mathbf{z}_{t, \tau}^{(N)}) \right) \right] d\tau, \quad (7)$$

259 where  $N$  is the number of particles and  $\mathcal{R}(\mathbf{z}_{t, \tau}^{(1)}, \dots, \mathbf{z}_{t, \tau}^{(N)})$  is the coupling between particles such that  
260 its gradient is the repulsive force<sup>2</sup>. Notice that the marginal distribution in (7) at each time-step  $\tau$  is  
261 given by (where  $Z$  is a normalizing constant)

$$262 q_\tau \left( \mathbf{z}_{t, \tau}^{(1)}, \dots, \mathbf{z}_{t, \tau}^{(N)} \right) = \frac{1}{Z} \mathcal{R} \left( \mathbf{z}_{t, \tau}^{(1)}, \dots, \mathbf{z}_{t, \tau}^{(N)} \right) \prod_{i=1}^N q_\tau \left( \mathbf{z}_{t, \tau}^{(i)} \right). \quad (8)$$

266 Throughout this work, we consider a pairwise kernel function  $k$  such that the repulsive force adopts the  
267 form  $\nabla_{\mathbf{z}_{t, \tau}^{(i)}} \mathcal{R} \left( \mathbf{z}_{t, \tau}^{(1)}, \dots, \mathbf{z}_{t, \tau}^{(N)} \right) = \nabla_{\mathbf{z}_{t, \tau}^{(i)}} \sum_{j=1}^N \log \left[ k \left( \mathbf{z}_{t, \tau}^{(i)}, \mathbf{z}_{t, \tau}^{(j)} \right) \right]^\gamma$ ; see the numerical experiments  
268

269 <sup>2</sup>We consider here a repulsive force because we seek diversity. However, an attractive force can be considered within this same framework.

for the particular instances of the kernel  $k$ . The repulsive force allows us to consider simple and flexible variational distributions that can discover multiple modes. For a simple illustration in the Gaussian case, see Appendix D.10.1. Finally, notice that when  $\gamma = 0$ , we recover the i.i.d. (non-repulsive) case.

#### 4.2 TACKLING LATENT INVERSION: AUGMENTATION OF THE VARIATIONAL DISTRIBUTION

As discussed in Section 3.2, solving (3) directly in the latent spaces yields blurry solutions. To tackle this, we propose to solve an augmented version of this problem, allowing us to decouple the data and the latent space of the diffusion model. Formally, we introduce an auxiliary variable  $\mathbf{x}_0$  defined in the data (pixel) space, which yields an augmented variational distribution  $q(\mathbf{z}_0, \mathbf{x}_0|\mathbf{y})$  and an augmented posterior as

$$p(\mathbf{z}_0, \mathbf{x}_0|\mathbf{y}) \propto \exp\left(-\frac{1}{2\sigma_v^2}\|\mathbf{y} - f(\mathbf{x}_0)\|^2 - \lambda g(\mathbf{z}_0) - \frac{1}{2\rho^2}\|\mathbf{x}_0 - \mathcal{D}(\mathbf{z}_0)\|^2\right), \quad (9)$$

where  $\rho$  controls the correlation between the variables  $\mathbf{x}_0$  and  $\mathbf{z}_0$ , and  $\exp(-\lambda g(\mathbf{z}_0))$  represents the prior distribution parameterized by the latent diffusion model. Notice that the definition in (9) implies that  $\mathbf{x}_0|\mathbf{z}_0 \sim \mathcal{N}(\mathcal{D}(\mathbf{z}_0), \rho^2\mathbf{I})$ , i.e., the conditional distribution of the data point ( $\mathbf{x}_0$ ) is centered at the value of the decoder applied to the latent point ( $\mathbf{z}_0$ ). However, this is not a delta but has some variance given by  $\rho^2$ . It can be shown that  $p(\mathbf{x}_0|\mathbf{y}, \lambda, \rho^2)$  converges in total variational to the true posterior  $p(\mathbf{x}_0|\mathbf{y}, \lambda)$  when  $\rho \rightarrow 0$  (Van Dyk and Meng, 2001; Vono et al., 2020) (details can be found in Appendix A.1). We can reformulate the optimization problem in (3) as

$$q(\mathbf{z}_0, \mathbf{x}_0|\mathbf{y}) = \operatorname{argmin}_{q(\mathbf{z}_0, \mathbf{x}_0|\mathbf{y})} \operatorname{KL}(q(\mathbf{z}_0, \mathbf{x}_0|\mathbf{y})\|p(\mathbf{z}_0, \mathbf{x}_0|\mathbf{y})). \quad (10)$$

When considering a diffusion model as data prior, our problem boils down to minimizing the variational lower bound, formalized in Proposition 1.

**Proposition 1** *Assuming we have access to a diffusion model  $\nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t)$  for the prior on  $\mathbf{z}_0$ , then the KL minimization w.r.t  $q$  in (10) is equivalent to minimizing the variational bound, which can be done by solving the following optimization problem*

$$\begin{aligned} \min_{q(\mathbf{x}_0, \mathbf{z}_0|\mathbf{y})} & \mathbb{E}_{q(\mathbf{z}_0|\mathbf{y})}[H(q(\mathbf{x}_0|\mathbf{z}_0, \mathbf{y}))] + \mathbb{E}_{q(\mathbf{x}_0, \mathbf{z}_0|\mathbf{y})} \left[ \frac{1}{2\sigma_v^2} \|\mathbf{y} - f(\mathbf{x}_0)\|^2 \right] \\ & + \mathbb{E}_{q(\mathbf{x}_0, \mathbf{z}_0|\mathbf{y})} \left[ \frac{1}{2\rho^2} \|\mathbf{x}_0 - \mathcal{D}(\mathbf{z}_0)\|^2 \right] + \int_0^T \tilde{\omega}(t) \mathbb{E}_{q(\mathbf{z}_t|\mathbf{y})} \left[ \|\nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t | \mathbf{y}) - \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t)\|_2^2 \right] dt. \end{aligned} \quad (11)$$

The proof is in Appendix A.2. When  $\rho \rightarrow 0$ , then  $\mathbf{x}_0 = \mathcal{D}(\mathbf{z}_0)$  and the augmented KL optimization boils down to the objective (3). To solve the problem in Proposition 1, we need to specify the variational distribution  $q(\mathbf{x}_0, \mathbf{z}_0|\mathbf{y})$ ; we now incorporate our result from Section 4.1.

#### 4.3 REPULSIVE LATENT SCORE DISTILLATION FOR SOLVING INVERSE PROBLEMS

We now derive RLSD, which integrates the techniques from Sections 4.1 and 4.2. Specifically, we apply the repulsive variational distribution introduced in Section 4.1 to instantiate the augmented variational formulation detailed in Proposition 1. By employing the particle approximation defined in (8), we define a multimodal distribution that enables a better exploration of the posterior’s search space, facilitating the discovery of multiple modes and addressing (c1). Notably, this variational approximation yields a tractable gradient, formalized in Proposition 2.

**Proposition 2** *When considering the variational distribution defined in (8), the KL minimization w.r.t  $q(\mathbf{x}_0, \mathbf{z}_0|\mathbf{y})$  defined in Proposition 1 can be approximated with an ensemble of  $N$  particles and admits the following gradient*

$$\frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{u}^{(i)}} \left[ \frac{1}{2\sigma_v^2} \|\mathbf{y} - f(\mathbf{x}_0^{(i)})\|^2 + \frac{1}{2\rho^2} \|\mathbf{x}_0^{(i)} - \mathcal{D}(\mathbf{z}_0^{(i)})\|^2 \right] + \nabla_{\mathbf{z}_0^{(i)}} \operatorname{reg}(\mathbf{z}_{t,\tau}^{(1)}, \dots, \mathbf{z}_{t,\tau}^{(N)}) \quad (12)$$

for  $i = 1, \dots, N$  and where  $\mathbf{u}^{(i)} = [\mathbf{x}_0^{(i)}, \mathbf{z}_0^{(i)}]$ . The regularization term is given by

$$\nabla_{\mathbf{z}_{0,\tau}^{(i)}} \text{reg}(\mathbf{z}_{t,\tau}^{(1)}, \dots, \mathbf{z}_{t,\tau}^{(N)}) = \mathbb{E}_{\epsilon, t} \left[ \lambda_t \left( \epsilon_{\theta}(\mathbf{z}_{t,\tau}^{(i)}, t) - \epsilon - \nabla_{\mathbf{z}_{t,\tau}^{(i)}} \gamma \sigma_t \log \sum_{j=1}^N k(\mathbf{z}_{t,\tau}^{(i)}, \mathbf{z}_{t,\tau}^{(j)}) \right) \right]. \quad (13)$$

where  $\lambda_t := \frac{T\alpha_t}{\sigma_t} \frac{d\omega(t)}{dt}$  and  $\gamma \geq 0$ .

The proof can be found in Appendix A.3. The gradient defined in Proposition 2 comprises three terms: a *measurement matching term*, an *error term* measuring the discrepancy between the variable in the pixel space and the decoded latent, and a *regularization term* that combines a *score-matching regularizer with a diversity-promoting component*. This repulsion term acts as a second regularizer, enhancing diversity during the sampling process. Our approach is not limited to any specific latent diffusion model.

**Practical algorithm for sampling using RLSD.** The underlying optimization of the gradient update in Proposition 2, a particular case of Proposition 1, is highly non-convex (diffusion denoiser and the repulsion term) and challenging to solve. To alleviate this, we adopt a half quadratic splitting technique (Geman and Yang, 1995). The algorithm is shown in Algorithm 1; we define  $\text{sg}[\cdot]$  as stopped-gradient operator to emphasize that the term inside it is not differentiated during the optimization step. We denote  $\tilde{\rho} = \frac{\sigma_t^2}{\rho^2}$ . For the weighting function  $\lambda_t$  ( $\omega(t)$  is embedded) and timesteps, we follow the strategy introduced in Mardani et al. (2024), where  $\lambda_t = \lambda(\sigma_t/\alpha_t)$ , and the timesteps follow a decreasing order (from  $t_{\max}$  to  $t_{\min}$ ); we fix  $t_{\max} = T$  and  $t_{\min} = 0$ . Regarding *computational burden*, our final algorithm only performs one backpropagation through the decoder in the  $\mathbf{z}$ -step and  $N$  backpropagations to compute the repulsive kernel (lines 7 and 9 in Alg. 1 are with respect to all particles). The complexity of the repulsive force depends on the number of particles as well as the domain of the kernel. Notice that the amount of particles is a hyperparameter, allowing us to control the trade-off between diversity and speed. Importantly, in contrast to previous works, we do not backpropagate through the score network.

---

#### Algorithm 1 RLSD for solving inverse problems

---

**Require:**  $\mathbf{y}, f(\cdot), L, \epsilon_{\theta}(\mathbf{z}_t, t), \mathcal{D}(\cdot), \{\lambda, \gamma, \tilde{\rho}, l_{r_x}, l_{r_z}\}$

- 1: Initialize  $\{\mathbf{x}_{i,0}^0\}_{i=1}^N, \{\mathbf{z}_{i,0}^0\}_{i=1}^N$
  - 2: **for**  $\ell = 1$  to  $L$  **do**
  - 3:    $t = T - \frac{\ell}{L}T$  and  $\epsilon \sim \mathcal{N}(0, \mathbf{I})$
  - 4:    $\lambda_t = \lambda(\sigma_t/\alpha_t)$
  - 5:    $\mathbf{z}_{i,t}^{\ell} = \alpha_t \mathbf{z}_{i,0}^{\ell} + \sigma_t \epsilon$
  - 6:    $\mathcal{L}_z = \sum_{i=1}^N \|\mathbf{x}_i^{\ell} - \mathcal{D}(\mathbf{z}_{i,0}^{\ell})\|^2 + \lambda_t \left( \text{sg} \left[ \epsilon_{\theta}(\mathbf{z}_{i,t}^{\ell}, t) - \epsilon - \gamma \nabla_{\mathbf{z}_t^{(i)}} \sigma_t \log \sum_{j=1}^N k(\mathbf{z}_t^{(i)}, \mathbf{z}_t^{(j)}) \right] \right)^{\top} \mathbf{z}_{i,0}^{\ell}$
  - 7:    $\mathbf{z}_0^{\ell} = \text{OptimizerStep}_{\mathbf{z}_0^{\ell}}(\mathcal{L}_z, l_{r_z})$
  - 8:    $\mathcal{L}_x = \sum_{i=1}^N \|\mathbf{y} - f(\mathbf{x}_i^{\ell})\|^2 + \tilde{\rho} \|\mathbf{x}_i^{\ell} - \mathcal{D}(\mathbf{z}_{i,0}^{\ell})\|^2$
  - 9:    $\mathbf{x}_0^{\ell} = \text{OptimizerStep}_{\mathbf{x}_0^{\ell}}(\mathcal{L}_x, l_{r_x})$
  - 10: **end for**
  - 11: **return**  $\{\mathbf{x}_{i,0}^L\}_{i=1}^N$
- 

## 5 EXPERIMENTS

In this section, we compare RLSD against state-of-the-art (SoTA) methods for solving inverse problems using latent diffusion models. We consider 100 samples from the validation set of FFHQ (Karras et al., 2019) used in Chung et al. (2022a). We compute PSNR [dB], LPIPS, and FID as metrics. Throughout the experiments, we seek to show the following:

- Our method generates more diverse solutions, in particular for tasks like inpainting and phase retrieval,

- When diversity is not relevant, our augmented formulation generates high-quality samples and outperforms baseline methods.

**Sampling setup.** Unless we state otherwise, we consider 1000 steps (the full denoising trajectory) for all the cases. We denote by NonAug-RLSD our repulsive method without augmentation (see Alg. 2 in Appendix D.1), and by NonRepuls-RLSD our method with augmentation but without repulsion. We consider Adam (Kingma, 2014) in the optimization steps (lines 7 and 9 in Alg. 1) and set the momentum pair (0.9, 0.99). We randomly initialize variables  $\mathbf{x}$  and  $\mathbf{z}$  and generate a batch of  $N = 4$  particles per noisy measurement. Regarding the pre-trained model, we consider Stable diffusion v2.1, although other latent diffusion models can be used. As diversity metric, we evaluate the pairwise diversity as the  $1 - \text{cosine similarity}$  between the  $N$  particles. Lastly, for the kernel function, we consider a RBF:  $k(\mathbf{z}_i, \mathbf{z}_j) = \exp(-\frac{\|g_{\text{DINO}}(\mathbf{z}_i) - g_{\text{DINO}}(\mathbf{z}_j)\|^2}{h_t})$ , where  $h_t = m_t^2 / \log N$ ,  $m_t$  is the median particle distance (Liu and Wang, 2016) and  $g_{\text{DINO}}$  is a pre-trained neural network (Caron et al., 2021). Details about implementation are in Appendix D.1, and ablation analysis in D.8.

**Baselines methods.** As we focus on methods that leverage large pre-trained models such as Stable diffusion, we compare with the recent PSLD (Rout et al., 2024) and Latent RED-Diff. For completeness, we also include a comparison with SoTA methods in the pixel-domain, namely DPS (Chung et al., 2022a) and RED-diff (Mardani et al., 2024). Details about the implementation of each method are in Appendix D.1. Given that pixel-based diffusion solvers generate images at  $256 \times 256$ , we follow the strategy from Rout et al. (2024) and downsample the results generated by our sampler, which have a  $512 \times 512$  resolution, to do a fair comparison.

## 5.1 INPAINTING

Inpainting, with its inherent ambiguity, provides a suitable benchmark to showcase two key aspects of RLSD: 1) high-quality reconstruction and 2) enhanced diversity achieved through the repulsion term. Additional linear inverse problems such as super resolution and deblurring are detailed in Appendix D. Specifically, we consider a box hiding half of the faces (see Fig. 2 and Appendix D.7). For the hyperparameters, we set  $\lambda = 0.14$ ,  $\tilde{\rho} = 0.075$ ,  $l_{r_x} = 0.4$  and  $l_{r_z} = 0.8$ . Results in Table 2 show that *RLSD outperforms their baselines* in image quality (PSNR and FID); in particular, it outperforms PSLD, the other sampler at a resolution of  $512 \times 512$ . Moreover, it demonstrates that our method can trade-off diversity for quality by modifying the weight  $\gamma$ : while RLSD ( $\gamma = 50$ ) achieves higher diversity than NonRepuls-RLSD, the later achieves better performance.

This highlights RLSD’s ability to combine superior reconstruction with higher diversity, suggesting a *mixed strategy where some particles interact and others do not*. Indeed, this strategy (Hybrid-RLSD), where three particles interact and one particle is propagated independently of the ensemble, combines the best of both worlds, achieving the best balance between quality and diversity.

**Diversity-quality trade off.** Fig. 2 showcases the diversity-quality trade-off. While PSLD generates four diverse samples of lower quality, Non-Repuls RLSD tends to fill the four images with very similar solutions. On the other hand, when considering RLSD ( $\gamma = 50$ ), the results are more diverse while maintaining high quality: in images 1 and 3, the woman has the left eye hidden, while none of the generated images with NonRepuls-RLSD show this. We defer to Appendix D.10 a more exhaustive analysis of diversity-quality trade off.

## 5.2 NON-LINEAR INVERSE PROBLEMS

We consider in this case nonlinear inverse problems. Given that PSLD only works on linear inverse problems, we compare against latent DPS and latent RED-Diff.

**Phase retrieval.** We first consider phase retrieval, which deals with reconstructing the phase from only magnitude observations in the Fourier domain. Phase retrieval is known as a highly ill-posed problem, given that it is invariant to  $180^\circ$  rotation, which yields two equally probable modes. Thus, its posterior has multiple modes, which are discrete and isolated. We follow the strategy from DPS (Chung et al., 2022a), where an oversampling of rate 2 is used. We consider 6 particles for the particle variational approximation, and 6 independent particles for the non-repulsive case. Furthermore, we set  $\gamma = 30$ , and we only consider repulsion between  $t \in [0.4T, T]$ . Results are shown in Table 3. First, we

Table 2: Box inpainting (half face) with  $\sigma_v = 0.001$  - FFHQ 512. For evaluation, and to compare with the other methods, we downsample the estimated images of RLSD and PSLD to 256. The best method for each metric is bolded.

Sampler	PSNR [dB] $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$	Diversity
PSLD	22.72	0.082	57.7	0.03
Latent RED-diff	23.5	0.15	92.59	0.009
RED-Diff (Pixel)	23.1	<b>0.067</b>	<u>29.79</u>	0.004
DPS (Pixel)	23.4	0.14	78.88	<b>0.04</b>
NonAug-RLSD ( $\gamma = 50$ )	23.34	0.164	98.65	<u>0.035</u>
NonRepuls-RLSD	<b>24.98</b>	<u>0.079</u>	<b>29.18</b>	0.004
RLSD ( $\gamma = 50$ )	24.69	0.111	31.41	0.015
Hybrid-RLSD	<u>24.72</u>	0.096	30.48	0.018



Figure 2: Inpainting half face using (from top to bottom): Ground truth and Measurement, PSLD, NonRepuls-RLSD, and RLSD ( $\gamma = 50$ ). We generate four samples for each method, starting from different initializations. For RLSD, the samples interact through the repulsion term. First, both NonRepuls-RLSD and RLSD outperform PSLD across all four images. Second, while images 1, 3 and 2, 4 from RLSD differ noticeably (e.g., images 1 and 3 have the left eye hidden), all samples generated by NonRepuls-RLSD appear quite similar.

observe that for this experiment, the augmented variational approximation (NonRepuls-RLSD and RLSD) entails a more unstable algorithm, and thus, not converging to good modes.

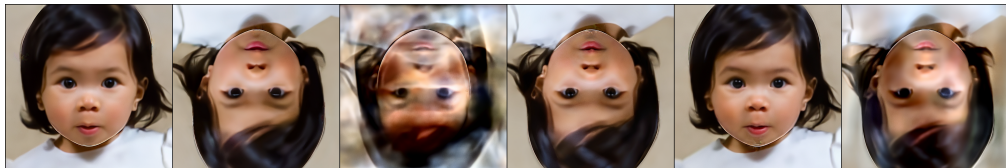
Second, the results show that our proposed method NonAug-RLSD is more stable, and that the particle approximation effectively captures more modes. In particular, in Fig. 3 we show an example where the latent RED-diff generates 5 images that look similar, while NonAug-RLSD generates 6 samples that corresponds to different modes. This showcases that our methods indeed promote diversity, even for a nonlinear inverse problems and at a resolution of  $512 \times 512$ .

**High dynamic range (HDR).** We try HDR, which performs the clipping function  $f(\mathbf{x}) = \text{clip}(2\mathbf{x}, -1, 1)$  on the normalized RGB pixels. HDR is known to be simpler than phase retrieval, and where diversity is not fundamental. Again, we consider Latent DPS as PSLD does not work for nonlinear inverse problems. Results are in Table 4, where NonRepuls-RLSD outperforms all the other baselines. This is aligned with our claim that our method can trade-off quality for diversity: in this case, it is better to focus on quality by  $\gamma = 0$ .



Table 3: Phase-retrieval with  $\sigma_v = 0.001$  on FFHQ 512. The best method for each metric is bolded.

Sampler	PSNR [dB] $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$
Latent-DPS	14.98	0.618	291.68
Latent-RED-diff ( $\gamma = 0$ )	19.65	0.458	173.18
NonAug-RLSD ( $\gamma = 30$ )	<b>24.21</b>	<b>0.359</b>	<b>130.09</b>
NonRepuls-RLSD	18.33	0.495	223.14
RLSD ( $\gamma = 30$ )	20.43	0.449	207

(a) NonAug-RLSD ( $\gamma = 30$ ).

(b) Latent RED-Diff.

Figure 3: Results for Phase Retrieval. Adding repulsion between particles allows to sample from different modes (top row).

Table 4: HDR with  $\sigma_v = 0.001$  on FFHQ 512. The best method for each metric is bolded.

Sampler	PSNR [dB] $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$
Latent-DPS	15.77	0.449	181.19
Latent-RED-diff	25.68	0.200	93.71
NonAug-RLSD ( $\gamma = 30$ )	24.84	0.210	94.15
NonRepuls-RLSD	<b>27.10</b>	<b>0.092</b>	<b>38.57</b>
RLSD ( $\gamma = 30$ )	25.53	0.113	57.89

## 6 CONCLUSIONS AND LIMITATIONS

In this paper, we introduce **Repulsive Latent Score Distillation (RLSD)**, a *plug-and-play* variational sampler that leverages pre-trained latent diffusion models to solve inverse problems, balancing quality, diversity, and computational efficiency. Inspired by the Wasserstein gradient flow of score distillation, **RLSD** mitigates mode collapse and latent diffusion inversion.

To tackle (c1) mode collapse, we introduce a particle-based variational distribution with a *repulsion mechanism* based on kernel similarity. To handle (c2) latent inversion (from adversarial training), we propose *distribution augmentation* to decouple latent and pixel spaces. The algorithm applies two regularizations: **denoising** to enforce the prior and **repulsion** to promote diversity.

Numerical experiments demonstrate that RLSD merges the benefits of variational samplers (memory and compute efficiency) with posterior samplers (diversity), allowing control over speed and diversity through simple regularization weights. The repulsion force significantly boosts diversity in ill-posed problems like inpainting and phase retrieval.

Our method has some limitations. Including the repulsion term increases computational demands, complicating real-time use. The chosen repulsion kernel may not be optimal under high noise levels, suggesting a need for *adaptive kernel learning*. Moreover, the method introduces additional hyperparameters, requiring better coupling for noise levels and deriving repulsion weights based on the forward operator.

## REPRODUCIBILITY STATEMENT

In this work, we have taken several steps to ensure the reproducibility of our results. We provide a comprehensive description of the methodology in Section 4 of the main text and the full algorithm in pseudocode are in Algorithms 1 and 2. In addition to this, in Section 5 as well in Appendix D.1 we give details of all the hyperparameters for each experiment. Additionally, we have included all necessary proofs for theoretical claims in Appendix A. For the experiments, we gave details of all the datasets that we used and how we compute the metrics. We include a link to anonymous downloadable source code in Appendix D.1. In the README.md file in the repository it is explained the steps to run the code. Lastly, in Appendix D.10.3 we included four .gif files, which might require a pdf reader that can reproduce gifs.

## ETHICS STATEMENT

Our method has the potential to cause unintended negative consequences if not handled responsibly. Key ethical and societal risks include the amplification of biases, difficulties in verifying the authenticity of generated content, which could contribute to misinformation, and the economic impact on creative professionals. Additionally, there are concerns over the misuse of this technology for harmful purposes, privacy issues related to the datasets used, cultural insensitivity, and potential intellectual property conflicts surrounding AI-generated creations. Addressing these risks necessitates the development of strong ethical standards, regulatory frameworks, and safeguards to ensure fairness, privacy protection, and respect for cultural and intellectual property rights. Therefore, it is essential that RLSD and other generative models are applied with a clear understanding of their limitations, and that outcomes are validated carefully to reduce these risks.

## REFERENCES

- Mohammadreza Armandpour, Huangjie Zheng, Ali Sadeghian, Amir Sadeghian, and Mingyuan Zhou. Re-imagine the negative prompt algorithm: Transform 2D diffusion into 3D, alleviate janus problem and beyond. *arXiv preprint arXiv:2304.04968*, 2023.
- Jimmy Ba, Murat A Erdogdu, Marzyeh Ghassemi, Shengyang Sun, Taiji Suzuki, Denny Wu, and Tianzong Zhang. Understanding the variance collapse of SVGD in high dimensions. In *Intl. Conf. Learn. Repr. (ICLR)*, 2021.
- David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017.
- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF Int. Conf. Comput. Vis. (ICCV)*, pages 9650–9660, 2021.
- Yifan Chen, Daniel Zhengyu Huang, Jiaoyang Huang, Sebastian Reich, and Andrew M Stuart. Gradient flows for sampling: mean-field models, Gaussian approximations and affine invariance. *arXiv preprint arXiv:2302.11024*, 2023.
- Hyungjin Chung, Jeongsol Kim, Michael Thompson Mccann, Marc Louis Klasky, and Jong Chul Ye. Diffusion posterior sampling for general noisy inverse problems. In *Intl. Conf. Learn. Repr. (ICLR)*, 2022a.
- Hyungjin Chung, Byeongsu Sim, Dohoon Ryu, and Jong Chul Ye. Improving diffusion models for inverse problems using manifold constraints. *Advances in Neural Inf. Process. Syst. (NIPS)*, 35: 25683–25696, 2022b.
- Hyungjin Chung, Byeongsu Sim, and Jong Chul Ye. Come-closer-diffuse-faster: Accelerating conditional diffusion models for inverse problems through stochastic contraction. In *Proceedings of the IEEE/CVF Int. Conf. Comput. Vis. Pattern Recogn. (CVPR)*, pages 12413–12422, 2022c.
- Hyungjin Chung, Jong Chul Ye, Peyman Milanfar, and Mauricio Delbracio. Prompt-tuning latent diffusion models for inverse problems. *Intl. Conf. on Machine Learning (ICML)*, 2024.

- 594 Gabriele Corso, Yilun Xu, Valentin De Bortoli, Regina Barzilay, and Tommi Jaakkola. Particle  
595 guidance: non-IID diverse sampling with diffusion models. *Intl. Conf. Learn. Repr. (ICLR)*, 2024.  
596
- 597 Francesco D’Angelo and Vincent Fortuin. Annealed Stein variational gradient descent. *arXiv preprint*  
598 *arXiv:2101.09815*, 2021a.
- 599 Francesco D’Angelo and Vincent Fortuin. Repulsive deep ensembles are Bayesian. *Advances in*  
600 *Neural Inf. Process. Syst. (NIPS)*, 34:3451–3465, 2021b.  
601
- 602 Giannis Daras, Joseph Dean, Ajil Jalal, and Alex Dimakis. Intermediate layer optimization for  
603 inverse problems using deep generative models. In *Intl. Conf. on Machine Learning (ICML)*, pages  
604 2421–2432. PMLR, 2021.
- 605 Giannis Daras, Hyungjin Chung, Chieh-Hsin Lai, Yuki Mitsufuji, Jong Chul Ye, Peyman Milanfar,  
606 Alexandros G. Dimakis, and Mauricio Delbracio. A survey on diffusion models for inverse  
607 problems. *arXiv preprint arXiv:2410.00083*, 2024.  
608
- 609 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances*  
610 *in Neural Inf. Process. Syst. (NIPS)*, 34:8780–8794, 2021.
- 611 Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real NVP. In *Intl.*  
612 *Conf. Learn. Repr. (ICLR)*, 2017.  
613
- 614 Zehao Dou and Yang Song. Diffusion posterior sampling for linear inverse problem solving: A  
615 filtering perspective. In *Intl. Conf. Learn. Repr. (ICLR)*, 2024.
- 616 Andrew Duncan, Nikolas Nüsken, and Lukasz Szpruch. On the geometry of Stein variational gradient  
617 descent. *J. Mach. Learn. Res.*, 24(56):1–39, 2023.  
618
- 619 Elhadji C Faye, Mame Diarra Fall, and Nicolas Dobigeon. Regularization by denoising: Bayesian  
620 model and Langevin-within-split Gibbs sampling. *arXiv preprint arXiv:2402.12292*, 2024.
- 621 Berthy T Feng, Jamie Smith, Michael Rubinstein, Huiwen Chang, Katherine L Bouman, and  
622 William T Freeman. Score-based diffusion models as principled priors for inverse imaging. In  
623 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10520–10531,  
624 2023.  
625
- 626 Donald Geman and Chengda Yang. Nonlinear image recovery with half-quadratic regularization.  
627 *IEEE transactions on Image Processing*, 4(7):932–946, 1995.
- 628 Yuan-Chen Guo, Ying-Tian Liu, Ruizhi Shao, Christian Laforte, Vikram Voleti, Guan Luo, Chia-  
629 Hao Chen, Zi-Xin Zou, Chen Wang, Yan-Pei Cao, and Song-Hai Zhang. threestudio: A unified  
630 framework for 3d content generation. [https://github.com/threestudio-project/](https://github.com/threestudio-project/threestudio)  
631 [threestudio](https://github.com/threestudio-project/threestudio), 2023.
- 632 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on*  
633 *Deep Generative Models and Downstream Applications*, 2021.  
634
- 635 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*  
636 *Neural Inf. Process. Syst. (NIPS)*, 33:6840–6851, 2020.  
637
- 638 Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J  
639 Fleet. Video diffusion models. *Advances in Neural Inf. Process. Syst. (NIPS)*, 35:8633–8646, 2022.
- 640 Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen,  
641 et al. LoRA: Low-rank adaptation of large language models. In *Intl. Conf. Learn. Repr. (ICLR)*,  
642 2021.
- 643 Ajay Jain, Ben Mildenhall, Jonathan T Barron, Pieter Abbeel, and Ben Poole. Zero-shot text-guided  
644 object generation with dream fields. In *Proceedings of the IEEE/CVF Int. Conf. Comput. Vis.*  
645 *Pattern Recogn. (CVPR)*, pages 867–876, 2022.  
646
- 647 Richard Jordan, David Kinderlehrer, and Felix Otto. The variational formulation of the Fokker–Planck  
equation. *SIAM J. Math. Anal.*, 29(1):1–17, 1998.

- 648 Zahra Kadkhodaie and Eero Simoncelli. Stochastic solutions for linear inverse problems using the  
649 prior implicit in a denoiser. *Advances in Neural Inf. Process. Syst. (NIPS)*, 34:13242–13254, 2021.
- 650
- 651 Ulugbek S Kamilov, Charles A Bouman, Gregory T Buzzard, and Brendt Wohlberg. Plug-and-play  
652 methods for integrating physical and learned models in computational imaging: Theory, algorithms,  
653 and applications. *IEEE Signal Processing Magazine*, 40(1):85–97, 2023.
- 654 Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative  
655 adversarial networks. In *Proceedings of the IEEE/CVF Int. Conf. Comput. Vis. Pattern Recogn.*  
656 *(CVPR)*, pages 4401–4410, 2019.
- 657 Oren Katzir, Or Patashnik, Daniel Cohen-Or, and Dani Lischinski. Noise-free score distillation. In  
658 *Intl. Conf. Learn. Repr. (ICLR)*, 2024.
- 659
- 660 Bahjat Kawar, Gregory Vaksman, and Michael Elad. SNIPS: Solving noisy inverse problems  
661 stochastically. *Advances in Neural Inf. Process. Syst. (NIPS)*, 34:21757–21769, 2021.
- 662 Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising diffusion restoration  
663 models. *Advances in Neural Inf. Process. Syst. (NIPS)*, 35:23593–23606, 2022.
- 664
- 665 Jeongsol Kim, Geon Yeong Park, and Jong Chul Ye. Dream sampler: Unifying diffusion sampling  
666 and score distillation for image manipulation. *arXiv preprint arXiv:2403.11415*, 2024.
- 667 Subin Kim, Kyungmin Lee, June Suk Choi, Jongheon Jeong, Kihyuk Sohn, and Jinwoo Shin.  
668 Collaborative score distillation for consistent visual editing. In *Advances in Neural Inf. Process.*  
669 *Syst. (NIPS)*, volume 36, pages 73232–73257, 2023.
- 670
- 671 Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,  
672 2014.
- 673 Florian Knoll, Kristian Bredies, Thomas Pock, and Rudolf Stollberger. Second order total generalized  
674 variation (tgv) for mri. *Magnetic resonance in medicine*, 65(2):480–491, 2011.
- 675
- 676 Erich Kobler, Teresa Klatzer, Kerstin Hammernik, and Thomas Pock. Variational networks: con-  
677 necting variational methods and deep learning. In *Pattern Recognition: 39th German Conference,*  
678 *GCPR 2017, Basel, Switzerland, September 12–15, 2017, Proceedings 39*, pages 281–293. Springer,  
679 2017.
- 680 Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. DiffWave: A versatile  
681 diffusion model for audio synthesis. In *Intl. Conf. Learn. Repr. (ICLR)*, 2020.
- 682
- 683 Marc Lambert, Sinho Chewi, Francis Bach, Silvère Bonnabel, and Philippe Rigollet. Variational  
684 inference via Wasserstein gradient flows. *Advances in Neural Inf. Process. Syst. (NIPS)*, 35:  
685 14434–14447, 2022.
- 686 Rémi Laumont, Valentin De Bortoli, Andrés Almansa, Julie Delon, Alain Durmus, and Marcelo  
687 Pereyra. Bayesian imaging using plug & play priors: when Langevin meets Tweedie. *SIAM J.*  
688 *Imag. Sciences*, 15(2):701–737, 2022.
- 689 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr  
690 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conf.*  
691 *Comp. Vision (ECCV)*, pages 740–755. Springer, 2014.
- 692
- 693 Qiang Liu. Stein variational gradient descent as gradient flow. *Advances in Neural Inf. Process. Syst.*  
694 *(NIPS)*, 30, 2017.
- 695 Qiang Liu and Dilin Wang. Stein variational gradient descent: A general purpose Bayesian inference  
696 algorithm. *Advances in Neural Inf. Process. Syst. (NIPS)*, 29, 2016.
- 697
- 698 Weijian Luo, Tianyang Hu, Shifeng Zhang, Jiacheng Sun, Zhenguo Li, and Zhihua Zhang. Diff-  
699 instruct: A universal approach for transferring knowledge from pre-trained diffusion models.  
700 *Advances in Neural Inf. Process. Syst. (NIPS)*, 36, 2024.
- 701
- Morteza Mardani, Jiaming Song, Jan Kautz, and Arash Vahdat. A variational perspective on solving  
inverse problems with diffusion models. *Intl. Conf. Learn. Repr. (ICLR)*, 2024.

- 702 Roger B Nelsen. *An introduction to copulas*. Springer, 2006.  
703
- 704 Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3D using 2D  
705 diffusion. In *Intl. Conf. Learn. Repr. (ICLR)*, 2022.  
706
- 707 Javier Portilla, Vasily Strela, Martin J Wainwright, and Eero P Simoncelli. Image denoising using  
708 scale mixtures of gaussians in the wavelet domain. *IEEE Transactions on Image processing*, 12  
709 (11):1338–1351, 2003.
- 710 Yaniv Romano, Michael Elad, and Peyman Milanfar. The little engine that could: Regularization by  
711 denoising (RED). *SIAM J. Imag. Sciences*, 10(4):1804–1844, 2017.  
712
- 713 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
714 resolution image synthesis with latent diffusion models. 2022 ieee. In *Proceedings of the IEEE/CVF*  
715 *Int. Conf. Comput. Vis. Pattern Recogn. (CVPR)*, pages 10674–10685, 2021.
- 716 Litu Rout, Negin Raouf, Giannis Daras, Constantine Caramanis, Alex Dimakis, and Sanjay Shakkottai.  
717 Solving linear inverse problems provably via posterior sampling with latent diffusion models.  
718 *Advances in Neural Inf. Process. Syst. (NIPS)*, 36, 2024.  
719
- 720 Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang,  
721 Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition  
722 challenge. *Int. J. Comp. Vision*, 115:211–252, 2015.
- 723 Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet,  
724 and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022*  
725 *conference proceedings*, pages 1–10, 2022.  
726
- 727 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi  
728 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. LAION-5B: An  
729 open large-scale dataset for training next generation image-text models. *Advances in Neural Inf.*  
730 *Process. Syst. (NIPS)*, 35:25278–25294, 2022.
- 731 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised  
732 learning using nonequilibrium thermodynamics. In *Intl. Conf. on Machine Learning (ICML)*, pages  
733 2256–2265. PMLR, 2015.  
734
- 735 Bowen Song, Soo Min Kwon, Zecheng Zhang, Xinyu Hu, Qing Qu, and Liyue Shen. Solving inverse  
736 problems with latent diffusion models via hard data consistency. *Intl. Conf. Learn. Repr. (ICLR)*,  
737 2024.
- 738 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *Intl. Conf.*  
739 *Learn. Repr. (ICLR)*, 2020.  
740
- 741 Jiaming Song, Arash Vahdat, Morteza Mardani, and Jan Kautz. Pseudoinverse-guided diffusion  
742 models for inverse problems. In *Intl. Conf. Learn. Repr. (ICLR)*, 2022.  
743
- 744 Jiaming Song, Qinsheng Zhang, Hongxu Yin, Morteza Mardani, Ming-Yu Liu, Jan Kautz, Yongxin  
745 Chen, and Arash Vahdat. Loss-guided diffusion models for plug-and-play controllable generation.  
746 In *Intl. Conf. on Machine Learning (ICML)*, pages 32483–32498. PMLR, 2023.
- 747 Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of  
748 score-based diffusion models. *Advances in Neural Inf. Process. Syst. (NIPS)*, 34:1415–1428,  
749 2021a.
- 750 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben  
751 Poole. Score-based generative modeling through stochastic differential equations. In *Intl. Conf.*  
752 *Learn. Repr. (ICLR)*, 2021b.  
753
- 754 Francesco Tonolini, Jack Radford, Alex Turpin, Daniele Faccio, and Roderick Murray-Smith. Vari-  
755 ational inference for computational imaging inverse problems. *Journal of Machine Learning*  
*Research*, 21(179):1–46, 2020.

- 756 Dustin Tran, David Blei, and Edo M Airolidi. Copula variational inference. *Advances in Neural Inf.*  
757 *Process. Syst. (NIPS)*, 28, 2015.
- 758
- 759 Arash Vahdat, Karsten Kreis, and Jan Kautz. Score-based generative modeling in latent space.  
760 *Advances in Neural Inf. Process. Syst. (NIPS)*, 34:11287–11302, 2021.
- 761
- 762 David A Van Dyk and Xiao-Li Meng. The art of data augmentation. *J. Comp. Graph. Stat.*, 10(1):  
763 1–50, 2001.
- 764
- 765 Singanallur V Venkatakrisnan, Charles A Bouman, and Brendt Wohlberg. Plug-and-play priors for  
766 model based reconstruction. In *IEEE Global Conf. Signal and Info. Process. (GlobalSIP)*, pages  
767 945–948. IEEE, 2013.
- 768
- 769 Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computa-*  
770 *tion*, 23(7):1661–1674, 2011.
- 771
- 772 Maxime Vono, Nicolas Dobigeon, and Pierre Chainais. Asymptotically exact data augmentation:  
773 Models, properties, and algorithms. *J. Comp. Graph. Stat.*, 30(2):335–348, 2020.
- 774
- 775 Peihao Wang, Dejia Xu, Zhiwen Fan, Dilin Wang, Sreyas Mohan, Forrest Iandola, Rakesh Ranjan,  
776 Yilei Li, Qiang Liu, Zhangyang Wang, et al. Taming mode collapse in score distillation for  
777 text-to-3D generation. *arXiv preprint arXiv:2401.00909*, 2023.
- 778
- 779 Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Pro-  
780 lificdreamer: High-fidelity and diverse text-to-3D generation with variational score distillation.  
781 *Advances in Neural Inf. Process. Syst. (NIPS)*, 36, 2024.
- 782
- 783 Weihao Xia, Yulun Zhang, Yujiu Yang, Jing-Hao Xue, Bolei Zhou, and Ming-Hsuan Yang. Gan  
784 inversion: A survey. *IEEE Trans. Pattern Anal. Mach. Intell.*, 45(3):3121–3138, 2022.
- 785
- 786 Kai Zhang, Yawei Li, Wangmeng Zuo, Lei Zhang, Luc Van Gool, and Radu Timofte. Plug-and-play  
787 image restoration with deep denoiser prior. *IEEE Trans. Pattern Anal. Mach. Intell.*, 44(10):  
788 6360–6376, 2021.
- 789
- 790 Junzhe Zhu, Peiye Zhuang, and Sanmi Koyejo. HIFA: High-fidelity text-to-3d generation with  
791 advanced diffusion guidance. In *Intl. Conf. Learn. Repr. (ICLR)*, 2024.
- 792
- 793 Yuanzhi Zhu, Kai Zhang, Jingyun Liang, Jiezhong Cao, Bihan Wen, Radu Timofte, and Luc Van Gool.  
794 Denoising diffusion models for plug-and-play image restoration. In *Proceedings of the IEEE/CVF*  
795 *Int. Conf. Comput. Vis. Pattern Recogn. (CVPR)*, pages 1219–1229, 2023.
- 796
- 797 Nicolas Zilberstein, Chris Dick, Rahman Doost-Mohammady, Ashutosh Sabharwal, and Santiago  
798 Segarra. Annealed Langevin dynamics for massive MIMO detection. *IEEE Trans. Wireless*  
799 *Commun.*, 2022.
- 800
- 801 Nicolas Zilberstein, Ashutosh Sabharwal, and Santiago Segarra. Solving linear inverse problems  
802 using higher-order annealed Langevin diffusion. *IEEE Trans. Signal Process.*, 2024.

## 800 A TECHNICAL PROOFS

### 802 A.1 DATA AUGMENTATION

803 In Section 4.2 we consider an augmented variational distribution  $q_\rho(\mathbf{x}_0, \mathbf{z}_0|\mathbf{y})$  such that

$$805 \quad q_\rho(\mathbf{x}_0|\mathbf{y}) = \int q_\rho(\mathbf{x}_0, \mathbf{z}_0|\mathbf{y}) d\mathbf{z}_0. \quad (14)$$

806

807 Therefore, we seek a joint distribution such that Property 1 holds.

808

809 **Property 1** For all  $\mathbf{x}_0 \in \mathbb{R}^N$ , it holds  $\lim_{\rho \rightarrow 0} \pi_\rho(\mathbf{x}_0) = \pi(\mathbf{x}_0)$ .

Notice that our approach using data augmentation resembles some recent methods introduced in the bibliography. In Zhu et al. (2023), the authors propose to decouple data and diffusion model. More recently, a Bayesian version of the RED method was proposed in Faye et al. (2024). This method shares similarities with our method, in the sense that both are augmented versions that resemble RED. However, our method has three main differences: 1) it leverages latent diffusion models, which allows us to solve large-scale inverse problems ( $512 \times 512$  and beyond), 2) it uses the diffused trajectory to regularize the solution, and 3) it promotes diversity via the coupling term.

## A.2 PROOF OF PROPOSITION 1

We expand the KL objective as follows

$$\begin{aligned}
 \text{KL}(q(\mathbf{z}_0, \mathbf{x}_0|\mathbf{y})||p(\mathbf{z}_0, \mathbf{x}_0|\mathbf{y})) &= \int q(\mathbf{z}_0, \mathbf{x}_0|\mathbf{y}) \log \frac{q(\mathbf{z}_0, \mathbf{x}_0|\mathbf{y})}{p(\mathbf{z}_0, \mathbf{x}_0|\mathbf{y})} d\mathbf{z}_0 d\mathbf{x}_0 & (15) \\
 &= \int q(\mathbf{z}_0, \mathbf{x}_0|\mathbf{y}) \log \frac{q(\mathbf{x}_0|\mathbf{z}_0, \mathbf{y})q(\mathbf{z}_0|\mathbf{y})p(\mathbf{y})}{p(\mathbf{y}|\mathbf{x}_0)p(\mathbf{x}_0|\mathbf{z}_0)p(\mathbf{z}_0)} d\mathbf{z}_0 d\mathbf{x}_0 \\
 &= \underbrace{\int q(\mathbf{z}_0, \mathbf{x}_0|\mathbf{y}) \log q(\mathbf{x}_0|\mathbf{z}_0, \mathbf{y}) d\mathbf{z}_0 d\mathbf{x}_0}_{(i)} \\
 &\quad - \underbrace{\int q(\mathbf{z}_0, \mathbf{x}_0|\mathbf{y}) \log p(\mathbf{y}|\mathbf{x}_0) d\mathbf{z}_0 d\mathbf{x}_0}_{(ii)} \\
 &\quad - \underbrace{\int q(\mathbf{z}_0, \mathbf{x}_0|\mathbf{y}) \log p(\mathbf{x}_0|\mathbf{z}_0) d\mathbf{z}_0 d\mathbf{x}_0}_{(iii)} \\
 &\quad + \underbrace{\int q(\mathbf{z}_0, \mathbf{x}_0|\mathbf{y}) \frac{q(\mathbf{z}_0|\mathbf{y})}{p(\mathbf{z}_0)} d\mathbf{z}_0 d\mathbf{x}_0}_{(iv)} + \log p(\mathbf{y}).
 \end{aligned}$$

Based on the augmented posterior, we have for (i) and (iii) that

$$(i) = \int q(\mathbf{z}_0, \mathbf{x}_0|\mathbf{y}) \log p(\mathbf{y}|\mathbf{x}_0) d\mathbf{z}_0 d\mathbf{x}_0 = \mathbb{E}_{q(\mathbf{x}_0, \mathbf{z}_0|\mathbf{y})} \left[ \frac{1}{2\sigma_v^2} \|\mathbf{y} - f(\mathbf{x}_0)\|^2 \right] \quad (16)$$

and

$$(iii) = \int q(\mathbf{z}_0, \mathbf{x}_0|\mathbf{y}) \log p(\mathbf{x}_0|\mathbf{z}_0) d\mathbf{z}_0 d\mathbf{x}_0 = \mathbb{E}_{q(\mathbf{x}_0, \mathbf{z}_0|\mathbf{y})} \left[ \frac{1}{2\rho^2} \|\mathbf{x}_0 - \mathcal{D}(\mathbf{z}_0)\|^2 \right]. \quad (17)$$

Regarding the first term, we can write as

$$(i) = \int q(\mathbf{z}_0|\mathbf{y}) [q(\mathbf{x}_0|\mathbf{z}_0, \mathbf{y}) \log q(\mathbf{x}_0|\mathbf{z}_0, \mathbf{y})] d\mathbf{z}_0 d\mathbf{x}_0 = \mathbb{E}_{q(\mathbf{z}_0|\mathbf{y})} [H(q(\mathbf{x}_0|\mathbf{z}_0, \mathbf{y}))]. \quad (18)$$

Finally, the last term can be obtained by following theorem 2 in Song et al. (2021a), assuming that the score is learned exactly, namely  $\epsilon_\theta(\mathbf{z}_t; t) = -\sigma_t \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t)$ , and under some mild assumptions on the growth of  $\log q(\mathbf{z}_t|\mathbf{y})$  and  $p(\mathbf{z}_t)$  at infinity, we have

$$\text{KL}(q(\mathbf{z}_0|\mathbf{y})||p(\mathbf{z}_0)) = \int_0^T \frac{\beta(t)}{2} \omega(t) \mathbb{E}_{q(\mathbf{z}_t|\mathbf{y})} \left[ \|\nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t|\mathbf{y}) - \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t)\|_2^2 \right] dt \quad (19)$$

over the denoising diffusion trajectory  $\{\mathbf{z}_t\}$  for positive values  $\{\beta(t)\}$ . This essentially implies that a weighted score-matching over the continuous denoising diffusion trajectory is equal to the KL divergence.

## A.3 PROOF OF PROPOSITION 2

The first two terms are straightforward. We focus here on the regularization term. The regularization term in Proposition 1 corresponds to the score matching loss defined in Song et al. (2021a) For general weighting schemes  $\omega(t)$ , we have the following Lemma from Song et al. (2021a)



**Lemma 1** *The time-derivative of the KL divergence at timestep  $t$  obeys*

$$\frac{d \text{KL}(q(\mathbf{z}_t | \mathbf{y}) \| p(\mathbf{z}_t))}{dt} = -\frac{\beta(t)}{2} \mathbb{E}_{q(\mathbf{z}_t | \mathbf{y})} \left[ \|\nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t | \mathbf{y}) - \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t)\|^2 \right].$$

Then, under the condition  $\omega(0) = 0$ , the integral in Proposition 1 can be written as (see Song et al. (2021a))

$$\begin{aligned} & \int_0^T \frac{\beta(t)}{2} \omega(t) \mathbb{E}_{q(\mathbf{z}_t | \mathbf{y})} \left[ \|\nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t | \mathbf{y}) - \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t)\|^2 \right] dt \\ &= - \int_0^T \omega(t) \frac{d \text{KL}(q(\mathbf{z}_t | \mathbf{y}) \| p(\mathbf{z}_t))}{dt} dt \\ &\stackrel{(a)}{=} \underbrace{-\omega(t) \text{KL}(q(\mathbf{z}_t | \mathbf{y}) \| p(\mathbf{z}_t))}_0 \Big|_0^T + \int_0^T \omega'(t) \text{KL}(q(\mathbf{z}_t | \mathbf{y}) \| p(\mathbf{z}_t)) dt \\ &= \int_0^T \omega'(t) \mathbb{E}_{q(\mathbf{z}_t | \mathbf{y})} \left[ \log \frac{q(\mathbf{z}_t | \mathbf{y})}{p(\mathbf{z}_t)} \right] dt, \end{aligned}$$

where  $\omega'(t) := \frac{d\omega(t)}{dt}$ . The equality holds because  $\omega(t) \text{KL}(q(\mathbf{z}_t | \mathbf{y}) \| p(\mathbf{z}_t)) \Big|_0^T$  is zero at  $t = 0$  and  $t = T$ . This is because  $\omega(t) = 0$  by assumption at  $t = 0$ , and  $x_T$  becomes a pure Gaussian noise at the end of the diffusion process which makes  $p(\mathbf{z}_T) = q(\mathbf{z}_T | \mathbf{y})$  and thus  $\text{KL}(q(\mathbf{z}_T | \mathbf{y}) \| p(\mathbf{z}_T)) = 0$ .

Now, we consider our proposed variational distribution defined in (8). with  $N$  particles and the pairwise kernel. For each particle  $i$ , we apply the forward diffusion  $\mathbf{z}_t^{(i)} = \alpha_t \mathbf{z}_0^{(i)} + \sigma_t \epsilon$ , which yields the distribution  $q(\mathbf{z}_t^{(i)} | \mathbf{y}) = \mathcal{N}(\alpha_t \mathbf{z}_0^{(i)}, \sigma_t^2 I)$ , and thus  $\nabla_{\mathbf{z}_t} \log q_t(\mathbf{z}_t | \mathbf{y}) = -(\mathbf{z}_t^{(i)} - \alpha_t \mathbf{z}_0^{(i)}) / \sigma_t^2 = -\frac{\epsilon^{(i)}}{\sigma_t}$ . By applying the re-parameterization trick, we obtain

$$\begin{aligned} & \nabla_{\mathbf{z}_0^{(i)}} \text{reg}(\mathbf{z}_t^{(1)}, \dots, \mathbf{z}_t^{(N)}) = \tag{20} \\ & \int_0^T \omega'(t) \mathbb{E}_{\epsilon \sim \mathcal{N}(0,1)} \left[ \left( -\nabla_{\mathbf{z}_t^{(i)}} \gamma \log \sum_{j=1}^N k(\mathbf{z}_t^{(i)}, \mathbf{z}_t^{(j)}) + \nabla_{\mathbf{z}_t^{(i)}} \log q_t(\mathbf{z}_t^{(i)} | \mathbf{y}) - \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t^{(i)}) \right)^\top \frac{d\mathbf{z}_t^{(i)}}{d\mathbf{z}_0^{(i)}} \right] dt \\ & \int_0^T \omega'(t) \mathbb{E}_{\epsilon \sim \mathcal{N}(0,1)} \left[ \left( -\nabla_{\mathbf{z}_t^{(i)}} \gamma \log \sum_{j=1}^N k(\mathbf{z}_t^{(i)}, \mathbf{z}_t^{(j)}) - \frac{\epsilon}{\sigma_t} + \frac{\epsilon \boldsymbol{\theta}(\mathbf{z}_t^{(i)}; t)}{\sigma_t} \right)^\top \alpha_t \mathbf{I} \right] dt \\ & \int_0^T \omega'(t) \frac{\alpha_t}{\sigma_t} \mathbb{E}_{\epsilon \sim \mathcal{N}(0,1)} \left[ \left( -\nabla_{\mathbf{z}_t^{(i)}} \gamma \sigma_t \log \sum_{j=1}^N k(\mathbf{z}_t^{(i)}, \mathbf{z}_t^{(j)}) - \epsilon + \epsilon \boldsymbol{\theta}(\mathbf{z}_t^{(i)}; t) \right) \right] dt. \end{aligned}$$

We can rearrange terms to arrive at the following compact form

$$\begin{aligned} & \nabla_{\mathbf{z}_0^{(i)}} \text{reg}(\mathbf{z}_t^{(1)}, \dots, \mathbf{z}_t^{(N)}) = \tag{21} \\ & \mathbb{E}_{\lambda_t \sim \mathcal{U}[0, T], \epsilon \sim \mathcal{N}(0,1)} \left[ \lambda_t \left( -\nabla_{\mathbf{z}_t^{(i)}} \gamma \sigma_t \log \sum_{j=1}^N k(\mathbf{z}_t^{(i)}, \mathbf{z}_t^{(j)}) - \epsilon + \epsilon \boldsymbol{\theta}(\mathbf{z}_t^{(i)}; t) \right) \right] \end{aligned}$$

for  $\lambda_t := T\omega'(t)\alpha_t/\sigma_t$ . When considering the measurement matching term and the error between the ambient and the augmented variable, we obtain  $\lambda_t := T\omega'(t)\alpha_t/\sigma_t 4\sigma_v^2 \rho^2$ .

## B DIFFUSION FOR INVERSE PROBLEMS

Diffusion models are powerful generative models. Therefore, they have been used as deep generative priors to solve inverse problems. Given a pre-trained diffusion model, this involves running the

backward process using a guidance (likelihood) term that incorporates the measurement information. Formally, we can sample from the posterior  $p(\mathbf{x}_0|\mathbf{y})$  by running (22)

$$d\mathbf{x}_t = -\frac{1}{2}\beta(t)\mathbf{z}_t dt - \beta(t) [\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t)] + \sqrt{\beta(t)}d\mathbf{W}_t. \quad (22)$$

Early studies used Langevin dynamics for linear problems (Kadkhodaie and Simoncelli, 2021; Kawar et al., 2021; Laumont et al., 2022; Zilberstein et al., 2024; 2022), while others used DDPM (Kawar et al., 2022; Chung et al., 2022c;b; Ho et al., 2022). However, approximating the guidance term remains a challenge. Previous works addressed this with a Gaussian approximation of  $p(\mathbf{x}_0|\mathbf{x}_t)$  around the MMSE estimator via Tweedie’s formula, increasing computational burden (Chung et al., 2022a; Song et al., 2022; Kadkhodaie and Simoncelli, 2021; Song et al., 2023). These methods crudely approximate the posterior score, especially for non-small noise levels. One of the most effective methods is DPS (Chung et al., 2022a), which assumes:

$$p(\mathbf{y}|\mathbf{x}_t) \approx p(\mathbf{y}|\hat{\mathbf{x}}_0 := \mathbb{E}[\mathbf{x}_0|\mathbf{x}_t]) = \mathcal{N}(\mathbf{y}|f(\mathbb{E}[\mathbf{x}_0|\mathbf{x}_t]), \sigma_t^2\mathbf{I}).$$

Essentially, DPS approximates the likelihood with a unimodal Gaussian distribution center around the MMSE estimator  $\mathbb{E}[\mathbf{x}_0|\mathbf{x}_t]$ . Under this approximation, the term  $p(\mathbf{y}|\mathbf{x}_t)$  boils down to the gradient of a multivariate Gaussian. Although it achieves impressive results, the unimodal approximation is far from optimal. Furthermore, its adaptation to use a latent diffusion model is not straightforward, as explained in (Rout et al., 2024). Recent works (Rout et al., 2024; Song et al., 2024; Kim et al., 2024; Chung et al., 2024) extended this by sampling from the latent space of diffusion models but still face limitations due to the intractable model likelihood. In particular, PSLD incorporates an additional term to guide the reconstruction towards a fixed point of the autoencoder process. This yields the following guidance score, where a *gluing* term is added to circumvent the discontinuity issues at the boundary.

$$\begin{aligned} \nabla_{\mathbf{z}_t} \log p(\mathbf{y}|\mathbf{z}_t) = & \nabla_{\mathbf{z}_t} \log p(\mathbf{y}|\hat{\mathbf{x}}_0 = \mathcal{D}(\mathbb{E}[\mathbf{z}_0|\mathbf{z}_t])) + \\ & \gamma_t \nabla_{\mathbf{z}_t} \left\| \mathbb{E}[\mathbf{z}_0|\mathbf{z}_t] - \mathcal{E}(\mathcal{A}^T \mathbf{y} + (\mathbf{I} - \mathcal{A}^T \mathcal{A}) \mathcal{D}(\mathbb{E}[\mathbf{z}_0|\mathbf{z}_t])) \right\|^2, \end{aligned} \quad (23)$$

where

$$p(\mathbf{y}|\mathbf{z}_t) \approx \mathcal{N}(\mathbf{y}|f(\mathcal{D}(\mathbb{E}[\mathbf{z}_0|\mathbf{z}_t])), \sigma_t^2\mathbf{I}).$$

Notice that while the gluing term is effective for linear inverse problems, it cannot handle non-linear cases.

## C DISCUSSION ON VARIATIONAL AUGMENTED DISTRIBUTION

In Section 4.2 we introduced an augmented variational distribution instead of a variational formulation in the latent space directly. This decision stems from the observation that optimizing in the latent space often produces blurry reconstructions. We hypothesize that this is due to the nonlinearity of the decoder  $\mathcal{D}(\cdot)$  and the adversarial training of the encoder-decoder pair. Specifically, the autoencoder tends to *compress* fine details, resulting in reconstructions that capture high-level semantics but fail to reproduce the fine-grained features.

To address this issue, we correct deviations from the image manifold during optimization by using an augmented variational formulation. This introduces a coupling term between  $\mathbf{x}$  and  $\mathbf{z}$  to account for these deviations. First, we define the true augmented posterior distribution. We write  $\mathbf{x}_0 = \mathcal{D}(\mathbf{z}_0) + \sigma_z \epsilon$ , where  $\sigma_z$  is the variance of the posterior of the decoder. Intuitively, when optimizing (10), we are optimizing the following joint target posterior  $p(\mathbf{x}, \mathbf{z}|\mathbf{y})$ , which has the following two conditionals associated

$$p(\mathbf{x}|\mathbf{z}, \mathbf{y}) \sim \mathcal{N}(\mu(\mathbf{z}, \mathbf{y}), \Sigma(\mathbf{z})) \quad (24)$$

$$p(\mathbf{z}|\mathbf{x}) \propto p(\mathbf{x}|\mathbf{z})p(\mathbf{z}) \quad (25)$$

where  $\Sigma(\mathbf{z}) = \frac{1}{\sigma_v^2} \mathbf{A}^\top \mathbf{A} + \frac{1}{\sigma_z^2} \mathcal{D}(\mathbf{z})^\top \mathcal{D}(\mathbf{z})$  and  $\mu(\mathbf{z}, \mathbf{y}) = \Sigma(\mathbf{z})^{-1} \left( \frac{1}{\sigma_v^2} \mathbf{A}^\top \mathbf{y} + \frac{1}{\sigma_z^2} \mathcal{D}(\mathbf{z}) \right)$ , and  $p(\mathbf{z})$  is the diffusion prior. Therefore, the variational inference in the augmented formulation aims to approximate the first Gaussian with another Gaussian  $q(\mathbf{x}|\mathbf{z}, \mathbf{y}) \sim \mathcal{N}(\boldsymbol{\mu}_x, \boldsymbol{\sigma}_x^2 \mathbf{I})$ , with  $\boldsymbol{\sigma}_x \rightarrow 0$ , i.e., a

MAP estimate, and the second one with the particle approximation, promoting diversity throughout the diffused trajectory.

This analysis also helps explain why NonAug-RLSD outperforms its augmented variant in phase retrieval. Phase retrieval is a nonlinear inverse problem, making alters (24) intractable: we cannot express the mean and covariance as explained above. Consequently, the augmentation renders a more difficult problem, which might explain why it is unstable.

## D ADDITIONAL EXPERIMENTS

### D.1 IMPLEMENTATION DETAILS OF BASELINES

To facilitate reproducibility, we share an anonymous link of our source code <https://file.io/iQNq3U5GpsY6>. If the paper is accepted, we will publish in a public repository.

**PSLD/Latent-DPS.** We use the original code from Rout et al. (2024). We use the version with Stable Diffusion, and we select hyperparameters as detailed in the paper.

**RED-Diff.** We use the implementation from the original paper (Mardani et al., 2024). We follow the same weighting scheme, and we use  $\lambda = 0.25$  and  $l_r = 0.1$ . As pretrained models, for FFHQ we use the model from Chung et al. (2022a), while for ImageNet we use the one from Dhariwal and Nichol (2021).

**Latent RED-Diff.** This method is the same as RED-Diff but using a latent diffusion model (as explain in Section 3.2).

**DPS.** We use the implementation from the original paper (Chung et al., 2022a). We follow their configuration of hyperparameters. We use the same pretrained models as RED-Diff.

**FPS-SMC(Dou and Song, 2024).** We use the implementation from the original paper. We follow their configuration of hyperparameters. We use the same pretrained models as RED-Diff.

**PIGDM.** We use the implementation from the original paper (Song et al., 2022). We follow their configuration of hyperparameters. We use the same pretrained models as RED-Diff.

**ReSample** We use the implementation from the original paper (Song et al., 2022). We follow their configuration of hyperparameters. We use LDM-VQ-4 trained on FFHQ. We tried using Stable Diffusion (the original implementation does not support it), but we got worst results.

**NonAug-RLSD.** This method corresponds to our variant using the particle-based variational approximation (8) and without augmentation. For clarity, we show it in Alg. 2.

---

#### Algorithm 2 Non-augmented RLSD for solving inverse problems

---

**Require:**  $\mathbf{y}, f(\cdot), L, \epsilon_{\theta}(\mathbf{z}_t, t), \mathcal{D}(\cdot), \{\lambda, \gamma, \tilde{\rho}, l_{r_z}\}$

**for**  $l = 1$  to  $L$  **do**

Initialize  $\{\mathbf{z}_{i,0}^0\}_{i=1}^N$

$t = T - \frac{\ell}{L}T$  and  $\epsilon \sim \mathcal{N}(0, \mathbf{I})$

$\lambda_t = \lambda(\sigma_t/\alpha_t)$

$\mathbf{z}_{i,t}^{\ell} = \alpha_t \mathbf{z}_{i,0}^{\ell} + \sigma_t \epsilon$

$\mathcal{L}_z = \sum_{i=1}^n \|\mathbf{y} - f(\mathcal{D}(\mathbf{z}_i^{\ell}))\|^2 + \lambda_t \left( \text{sg} \left[ \epsilon_{\theta}(\mathbf{z}_{i,t}^{\ell}, t) - \epsilon - \gamma \nabla_{\mathbf{z}_t^{(i)}} \sigma_t \log \sum_{j=1}^N k(\mathbf{z}_t^{(i)}, \mathbf{z}_t^{(j)}) \right] \right)^{\top} \mathbf{z}_{i,0}^{\ell}$

$\mathbf{z}_0^{\ell} = \text{OptimizerStep}_{\mathbf{z}_0^{\ell}}(\mathcal{L}_z, l_{r_z})$

**end for**

**return**  $\{\mathbf{x}_{i,0}^L = \mathcal{D}(\mathbf{z}_{i,0}^L)\}_{i=1}^n$

---

**RLSD.** For our augmented formulation, we use Stable diffusion trained in the LAION (Schuhmann et al., 2022) dataset as its pre-trained model. For the kernel function, we consider a RBF  $k(\mathbf{z}_i, \mathbf{z}_j) = \exp(-\frac{\|g_{\text{DINO}}(\mathbf{z}_i) - g_{\text{DINO}}(\mathbf{z}_j)\|^2}{h_t})$  where  $h_t = m_t^2 / \log N$ ,  $m_t$  is the median particle distance (Liu and Wang, 2016) and  $g_{\text{DINO}}$  is a pre-trained neural network (Caron et al., 2021). Notice that NonRepuls-RLSD corresponds to  $\gamma = 0$ .

Regarding the similarity metric, we consider the cosine similarity in the range of DINO defined as follow

$$\text{Sim}(\mathbf{x}_1, \dots, \mathbf{x}_N) = \frac{1}{n(n-1)} \sum_{i \neq j} \frac{g_{\text{DINO}}(\mathbf{x}_i)^T g_{\text{DINO}}(\mathbf{x}_j)}{\|g_{\text{DINO}}(\mathbf{x}_i)\|_2 \|g_{\text{DINO}}(\mathbf{x}_j)\|_2}, \quad (26)$$

Base on this metric, we define diversity as

$$\text{Div}(\mathbf{x}_1, \dots, \mathbf{x}_N) = 1 - \text{Sim}(\mathbf{x}_1, \dots, \mathbf{x}_N). \quad (27)$$

Lastly, we consider an decreasing annealing schedule. It has been notice in previous works (Mardani et al., 2024; Zhu et al., 2024) that a decreasing timestep works better than sampling uniformly at random. As a consequence, we consider this scheme.

## D.2 SUPER RESOLUTION

We consider super resolution from  $\times 8$  downsampled images. In this case we use  $\lambda = 0.2$ ,  $\tilde{\rho} = 0.05$ ,  $l_{r_x} = 0.4$  and  $l_{r_z} = 0.6$ . The results are shown in Table 5. For additional comparison, we compare also with solvers that generate samples at  $256 \times 256$ . For this comparison, similar to Section 5.1, we downsample the result of RLSD to 256 and compare at that resolution. The results are in Table 6. This example illustrates the key difference between RLSD and PSLD for solving inverse problems. To be more specific, PSLD generates images of faces that have the typical artifacts when sampling with Stable Diffusion. On the other hand, RLSD leverages Stable Diffusion as multiple denoisers at different scale.

Table 5: SR  $\times 8$  with  $\sigma_v = 0.001$  - FFHQ 512. The best method for each metric and experiment is bolded.

Sampler	PSNR [dB] $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$
PSLD	24.82	0.314	81.31
Latent RED-diff	26.07	0.439	76.07
NonRepuls-RLSD	<b>28.39</b>	<b>0.286</b>	<b>65.42</b>

Table 6: SR  $\times 8$  with  $\sigma_v = 0.001$  - FFHQ 256. The best method for each metric and experiment is bolded.

Sampler	PSNR [dB] $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$
NonRepuls-RLSD	<b>28.4</b>	<b>0.149</b>	<b>65.42</b>
RED-Diff	25.69	0.264	104.59
DPS	23.83	0.175	90.45
IIGDM	24.44	<b>0.128</b>	82.01

For completeness, we also compare with RED-Diff when solving SR  $\times 4$ , as it has the same resolution as input than RLSD with SR  $\times 8$  (RED-Diff handles images of size 256x256). The results is Table 7.

## D.3 MOTION DEBLURRING

We consider Motion Blurring. In this case we use  $\lambda = 0.007$ ,  $\tilde{\rho} = 0.01$ ,  $l_{r_x} = 0.4$  and  $l_{r_z} = 0.3$ , and  $L = 500$ . In particular, we follow Chung et al. (2022a), where we convolve the image with a  $61 \times 61$  motion kernel that is randomly sampled with intensity  $0.3^2$ . The results are shown in Table 8. For additional comparison, we compare also with solvers that generate samples at  $256 \times 256$ . For this

Table 7: SR  $\times$  8 with  $\sigma_v = 0.001$  - FFHQ 256. The best method for each metric and experiment is bolded.

Sampler	PSNR [dB] $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$
NonRepuls-RLSD	28.4	<b>0.149</b>	<b>65.42</b>
RED-Diff (SR $\times$ 4)	<b>28.91</b>	0.157	78.41
DPS (SR $\times$ 4)	27.14	0.128	71.8
FPS-SMC (SR $\times$ 4)	27.36	0.21	120.49

comparison, similar to Section 5.1, we downsample the result of RLSD to 256 and compare at that resolution. The results are in Table 9.

Table 8: Motion Blurring with  $\sigma_v = 0.001$  - FFHQ 512. The best method for each metric and experiment is bolded.

Sampler	PSNR [dB] $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$
PSLD	25.17	0.389	135.22
Latent RED-diff	27.85	0.329	118.09
NonRepuls-RLSD	<b>30.4</b>	<b>0.23</b>	<b>56.79</b>

Table 9: Motion Blurring with  $\sigma_v = 0.001$  - FFHQ 256. The best method for each metric and experiment is bolded.

Sampler	PSNR [dB] $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$
NonRepuls-RLSD	<b>30.47</b>	<b>0.095</b>	<b>56.79</b>
RED-Diff	30.27	0.15	103.17
DPS	24.7	0.22	90.45
ReSample	26.82	0.115	72.74

#### D.4 INPAINTING WITH MASKED BOX (HALF-FACE)

We consider here additional baselines for box inpainting (half-face). See results in Table 10.

Table 10: Box inpainting (half face) with  $\sigma_v = 0.001$  - FFHQ 256. The best method for each metric and experiment is bolded.

Sampler	PSNR [dB] $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$
NonAug-RLSD ( $\gamma = 50$ )	23.34	0.164	98.65
NonRepuls-RLSD	<b>24.98</b>	<b>0.079</b>	<b>29.18</b>
RLSD ( $\gamma = 50$ )	24.69	0.111	31.41
Hybrid-RLSD	<u>24.72</u>	0.096	30.48
IIGDM	23.74	0.077	33.8
FPS-SMC	24.91	0.086	59.59
ReSample	19.44	0.2	146.68

#### D.5 INPAINTING WITH RANDOM MASK

We consider random inpainting where we drop 80% of the pixels. In this case we use  $\lambda = 0.009$ ,  $\tilde{\rho} = 0.08$ ,  $l_{r_x} = 0.4$  and  $l_{r_z} = 0.8$ . The numerical results are shown in Table 12, and visual examples are shown in Fig. 9.

Table 11: Random inpainting with (80%) mask and with ( $\sigma_v = 0.001$ ) - FFHQ 512. In bold is the best method for each metric and experiment.

Sampler	PSNR [dB] $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$
PSLD	28.53	0.212	65.14
NonRepuls-RLSD	<b>30.56</b>	<b>0.145</b>	<b>41.11</b>

Table 12: Random inpainting with (80%) mask and with ( $\sigma_v = 0.001$ ) - FFHQ 256. In bold is the best method for each metric and experiment.

Sampler	PSNR [dB] $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$
NonRepuls-RLSD	<b>30.56</b>	0.073	<b>41.11</b>
RED-Diff	28.55	0.074	62.87
DPS	26.48	0.15	95.44
ReSample	27.49	<b>0.062</b>	54.42

#### D.6 INPAINTING FREE MASK

We use the free masks (10% – 20%) from Saharia et al. (2022). For this experiment, we consider ImageNet (Russakovsky et al., 2015) to demonstrate that our method outperforms its baselines on other datasets. We consider the  $\lambda = 0.15$ ,  $\tilde{\rho} = 0.15$ ,  $l_{r_z} = 0.8$  and  $l_{r_x} = 0.4$ , and we consider 500 steps (instead of the full trajectory of 1000 steps) for both RLSD and PSLD. For PSLD, we use the parameters from their experiments with box inpainting (similar to the case of half face). The quantitative results are shown in Table 13, and qualitative results in Fig. 14.

Table 13: Free mask (Saharia et al., 2022) with ( $\sigma_v = 0.001$ ) - ImageNet 512. In bold is the best method for each metric and experiment.

Sampler	PSNR [dB] $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$
PSLD	22.56	0.154	69.43
NonRepuls-RLSD	<b>26.77</b>	<b>0.075</b>	<b>60.53</b>

**Results when considering the best sample across particles.** If instead of focusing on the average performance, we focus on the performance achieved by the best image among the ensemble of particles, then the conclusion is different in favor of the fill RLSD; see Table 14. This can be explained as follow: while some of the modes obtained with RLSD might not be as good as NonRepuls-RLSD, others might be better.

Table 14: Box inpainting (half face) with  $\sigma_v = 0.001$  - FFHQ 512. We consider the best particle across each batch of them.

Sampler	PSNR [dB] $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$
PSLD (mean)	21.34	0.10	57.7
PSLD (max)	22.72	0.082	57.7
NonRepuls-RLSD (mean)	24.98	0.079	<b>29.18</b>
NonRepuls-RLSD (max)	<u>25.82</u>	<u>0.071</u>	29.18
RLSD ( $\gamma = 50$ ) (mean)	24.69	0.111	31.41
RLSD ( $\gamma = 50$ ) (max)	<b>25.84</b>	<b>0.069</b>	31.41

#### D.7 QUALITATIVE RESULTS

1188  
 1189  
 1190  
 1191  
 1192  
 1193  
 1194  
 1195  
 1196  
 1197  
 1198  
 1199  
 1200  
 1201  
 1202  
 1203  
 1204  
 1205  
 1206  
 1207  
 1208  
 1209  
 1210  
 1211  
 1212  
 1213  
 1214  
 1215  
 1216  
 1217  
 1218  
 1219  
 1220  
 1221  
 1222  
 1223  
 1224  
 1225  
 1226  
 1227  
 1228  
 1229  
 1230  
 1231  
 1232  
 1233  
 1234  
 1235  
 1236  
 1237  
 1238  
 1239  
 1240  
 1241

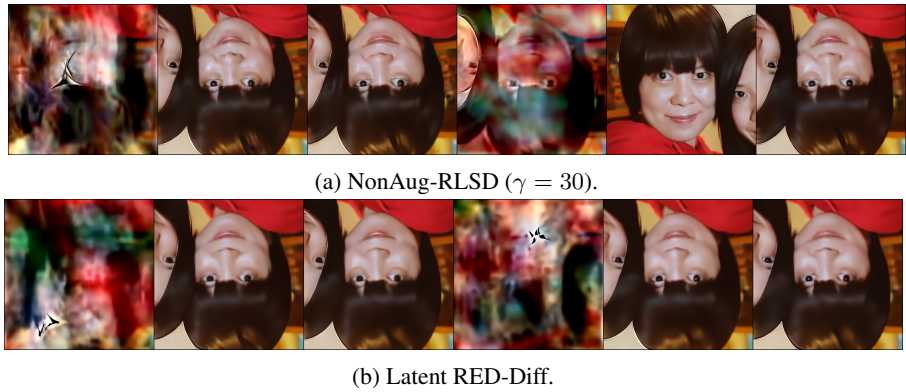


Figure 4: Phase Retrieval. Adding repulsion between particles promotes diversity, and allows to sample from different modes.



Figure 5: Inpainting half face using (from top to bottom): Ground truth and Measurement, PSLD, NonRepuls-RLSD, and RLSD ( $\gamma = 50$ ). We generate four samples for each method from a different initialization; for NonRepuls-RLSD, and RLSD, they interact through the repulsion term. First, NonRepuls-RLSD, and RLSD outperforms PSLD for all four images. Second, while images 1 and 3 from RLSD (last row) look different, the images 1 and 3 of NonRepuls-RLSD are similar; this illustrates that RLSD promotes diversity.



1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264



1265 Figure 6: Inpainting half face using (from top to bottom): Ground truth and Measurement, PSLD, NonRepuls-  
1266 RLSD, and RLSD ( $\gamma = 50$ ). We generate four samples for each method from a different initialization; for  
1267 RLSD, they interact through the repulsion term. First, NonRepuls-RLSD, and RLSD outperforms PSLD for all  
1268 four images. Second, image 4 from RLSD (last row) looks different (has brown eye), while the four samples of  
1269 NonRepuls-RLSD have blue eye; this illustrates that RLSD promotes diversity.

1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283



(a) NonAug-RLSD ( $\gamma = 30$ ).



(b) Latent RED-Diff.

1284 Figure 7: Phase Retrieval. We observe that adding repulsion between particles promotes diversity, and allows  
1285 to sample from different modes.

1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295



(a)

(b)

1293 Figure 8: Two qualitative examples of HDR. Both a) and b) corresponds to NonRepuls-RLSD, and each one has  
1294 Ground-truth; Measurement; Estimation. NonRepuls-RLSD generates an image of high-fidelity in a nonlinear  
1295 problem at  $512 \times 512$ .

1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312



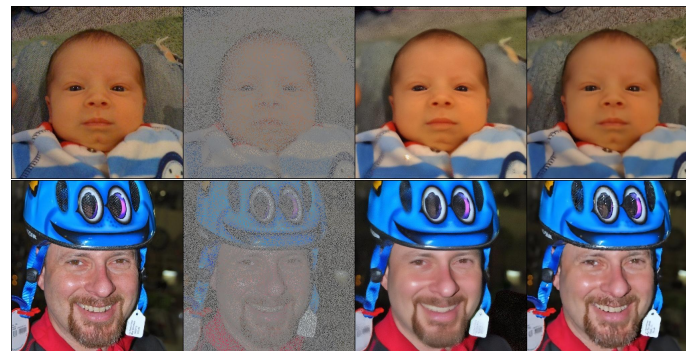
Figure 9: Random inpainting (80%): PSLD (top) and NonRepuls-RLSD (down), and four different samples for each method. Notice that NonRepuls-RLSD generates a better reconstruction of the background.

1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331



Figure 10: Random inpainting (80%): PSLD (top) and NonRepuls-RLSD (down), and four different samples for each method. Notice that NonRepuls-RLSD generates a sharper reconstruction of the face.

1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

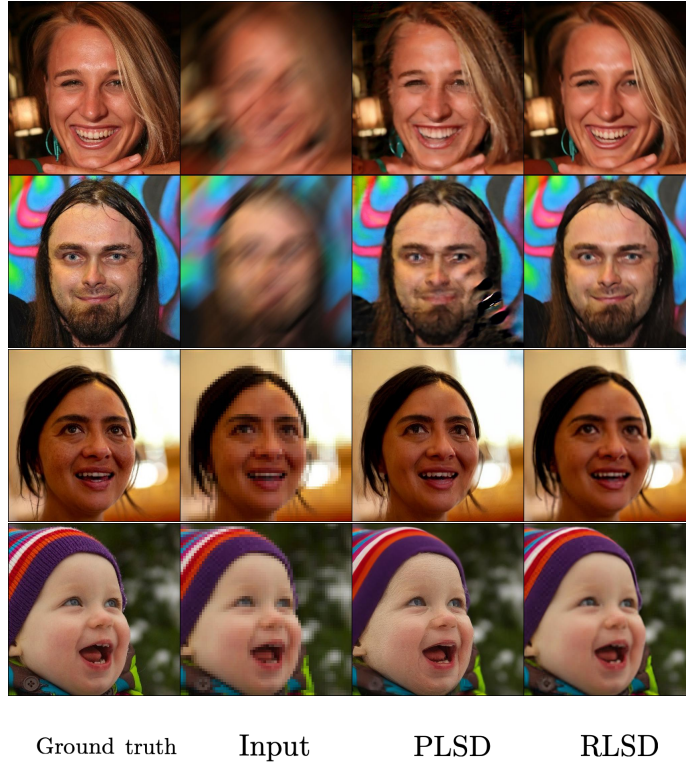


Ground truth      Input      PLSD      RLSD

Figure 11: Additional examples for random inpainting. Notice that RLSD have more details in the background for the first row, while for the second row it can reconstruct the details of the white label.

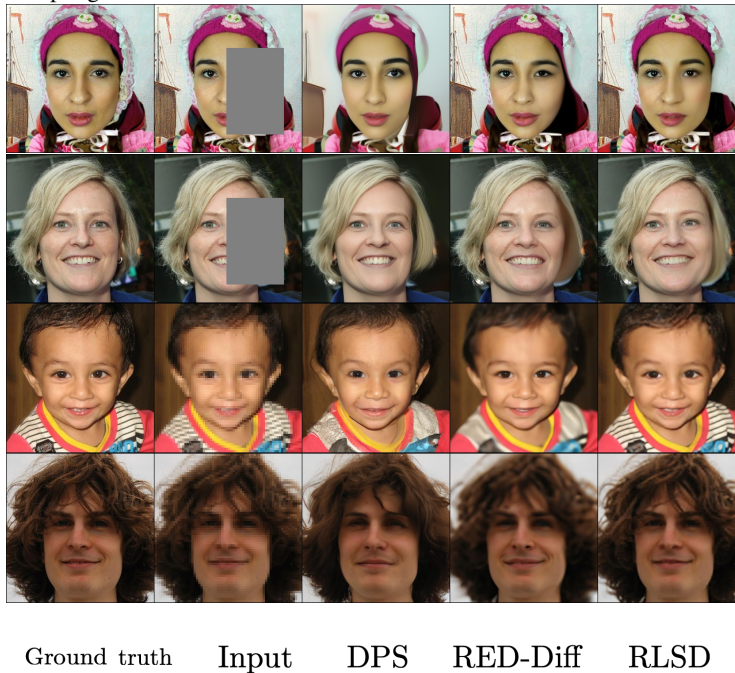


1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373



1374 Figure 12: Qualitative examples of different inverse problems for PSLD (third column) and NonRepuls-RLSD  
1375 (forth column): Rows 1 and 2 are motion blurring, Rows 3 and 4 are  $SR \times 8$  For motion deblurring, RLSL  
1376 clearly outperforms PSLD (PSLD reconstruct a broken image. For  $SR \times 8$ , PSLD introduces artifacts, which  
1377 are typical when sampling faces with Stable Diffusion.

1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398



1399 Figure 13: Qualitative examples of different inverse problems for DPS (third column), RED-Diff (forth column)  
1400 and NonRepuls-RLSD (fifth column): Rows 1 and 2 are half face inpainting, Rows 3 and 4 are  $SR \times 8$ .  
1401 Notice that for  $SR \times 8$ , both DPS and RED-Diff generates images that look different than the original. This is  
1402 expected considering that both operate at a resolution of  $256 \times 256$ . This demonstrates the advantage of using a  
1403 high-resolution model.

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

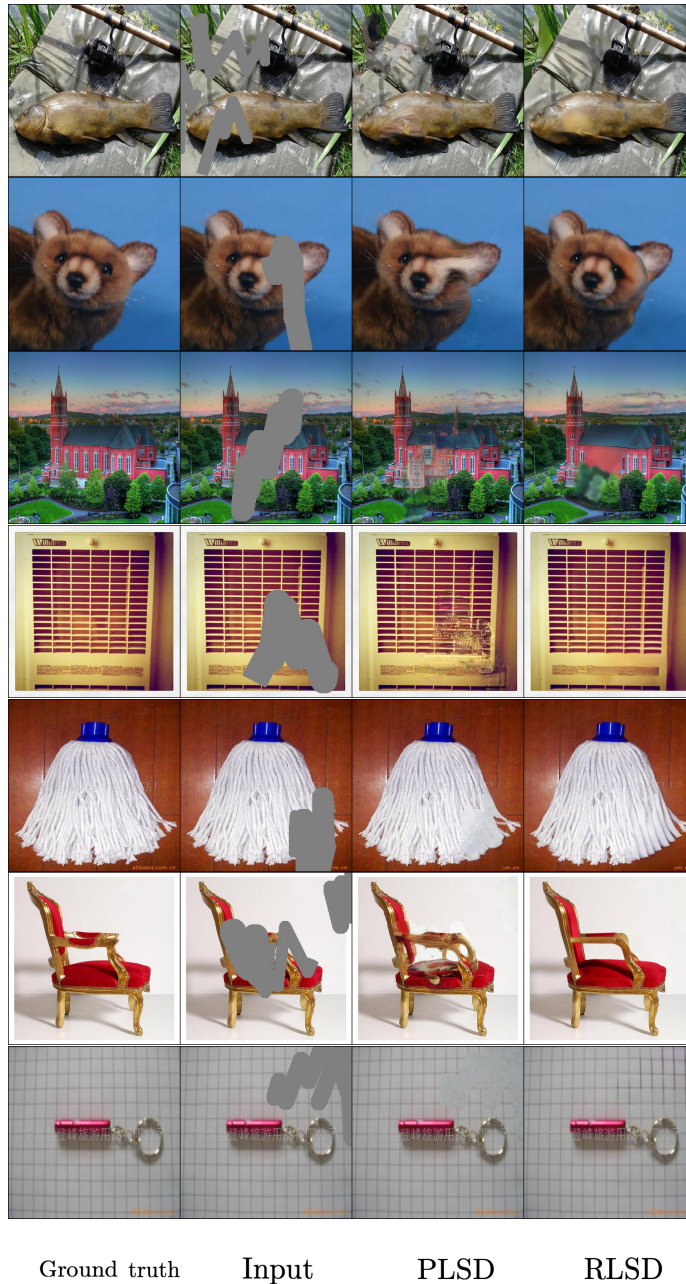


Figure 14: Qualitative examples of Inpainting with free mask on Imagenet. We consider PSLD (third column) and NonRepuls-RLSD (forth column). For rows 2, 4-7, RLSD outperforms PSLD by a large margin. Furthermore, notice that RLSD can reconstruct backgrounds that might be difficult, such as the one row 4. On the other hand, RLSD struggles to reconstruct fine-details, as it is shown in row 3. We expand on this in Appendix D.9.

## 1458 D.8 ABLATION OF RLSD

1459

## 1460 D.8.1 RUNNING TIME/NUMBER OF STEPS

1461

1462 We ablate the running time between RLSD and PSLD when generating two samples for super  
 1463 resolution  $\times 8$ . We ran the experiments on the same NVIDIA A100 GPU with 80GB. When running  
 1464 the full trajectory (1000 steps), the running time of RLSD without particles is 8.8 minutes, while  
 1465 PSLD is 9.2 minutes. When running 200 steps, the running time of RLSD without particles is 1.75  
 1466 minutes, while PSLD is 1.9 minutes. The results for this case are shown in Fig. 15. Notice that  
 1467 although the running time difference between RLSD and PSLD is not large, our formulation leverages  
 1468 the diffused trajectory as a denoiser. Therefore, we require fewer steps to obtain a high-fidelity  
 1469 estimation; for instance, for super resolution we need just 200 steps. This showcases that our method  
 works fine with just a few number of steps; see also Fig. 16.

1470

1471

1472

1473

1474

1475

1476

1477

1478

1479

1480

1481

1482

1483

1484

1485

1486

1487

1488

1489

1490

1491

1492

1493

1494

1495

1496

1497

1498

1499

1500

1501

1502

1503

1504

1505

1506

1507

1508

1509

1510

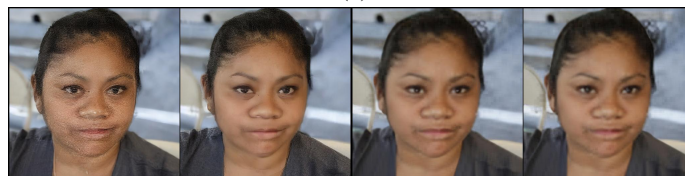
1511



(a)

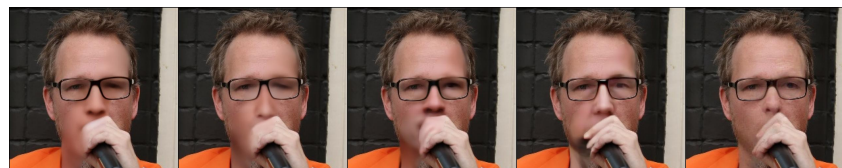


(b)



(c)

Figure 15: Reconstruction for different number of steps. a) Ground-truth (left) and Measurement (right). b) The two on the left are the reconstruction using PSLD with 200 steps (running time = 1.9 min), while the two of the right are with RLSD ( $\gamma = 0$ ) (running time = 1.75 min). c) The two on the left are the reconstruction using PSLD with 1000 steps (running time = 9.2 min), while the two of the right are with RLSD ( $\gamma = 0$ ) (running time = 8.8 min).



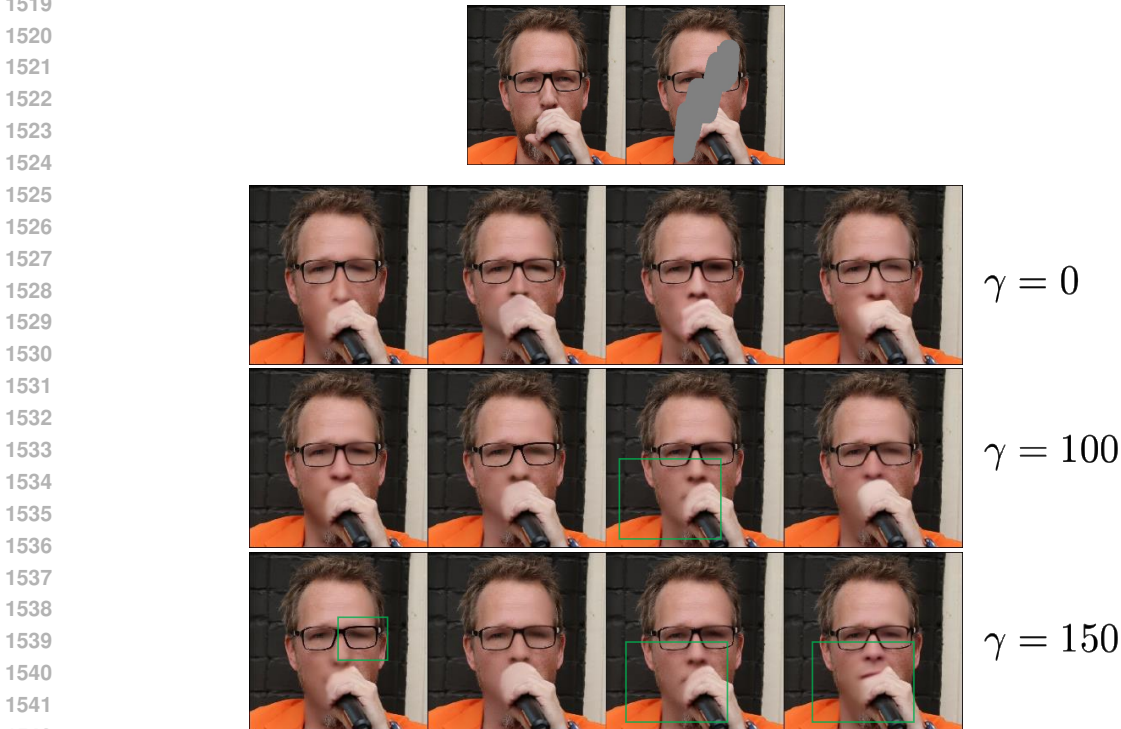
$L = 100$        $L = 200$        $L = 300$        $L = 500$        $L = 999$   
 LPIPS = 0.065    LPIPS = 0.063    LPIPS = 0.062    LPIPS = 0.059    LPIPS = 0.058

Figure 16: Reconstruction as a function of  $L$  (number of steps in the optimization). Increasing  $L$  modifies the number of denoisers that are used (the limits of the interval are the same, and it changes the step-size). Also, we need to decrease  $l_{r_z}$  (from 1 to 0.8) as well the coupling term  $\bar{\rho}$ , from 0.15 to 0.1.



1512 D.8.2 REPULSION TERM

1513  
 1514 **Effect of  $\gamma$ .** We showcase an example of the effect of  $\gamma$ . We consider  $N = 4$  particles and 200  
 1515 steps,  $\tilde{\rho} = 0.1, \lambda_t = 0.15, l_{r_x} = 0.4$  and  $l_{r_z} = 1$ . It is important to remark that when considering  
 1516 fewer steps, we need a higher  $\gamma$  (compared to the example in Section 5.1; the reason why this  
 1517 happens is due to the annealing factor in the repulsion term (13), given by  $\alpha(t)$ . The example for  
 1518  $\gamma = 0, 100, 150, 200$  is shown in Fig. 17.



1543 Figure 17: Reconstruction as a function of  $\gamma$  (number of repulsion), for  $L = 200$  and  $N = 4$ . Increasing  $\gamma$   
 1544 enhances diversity: for  $\gamma = 150$ , 2 out of 4 samples show the mouth, while for  $\gamma = 0$  all samples have the  
 1545 mouth hidden.  
 1546

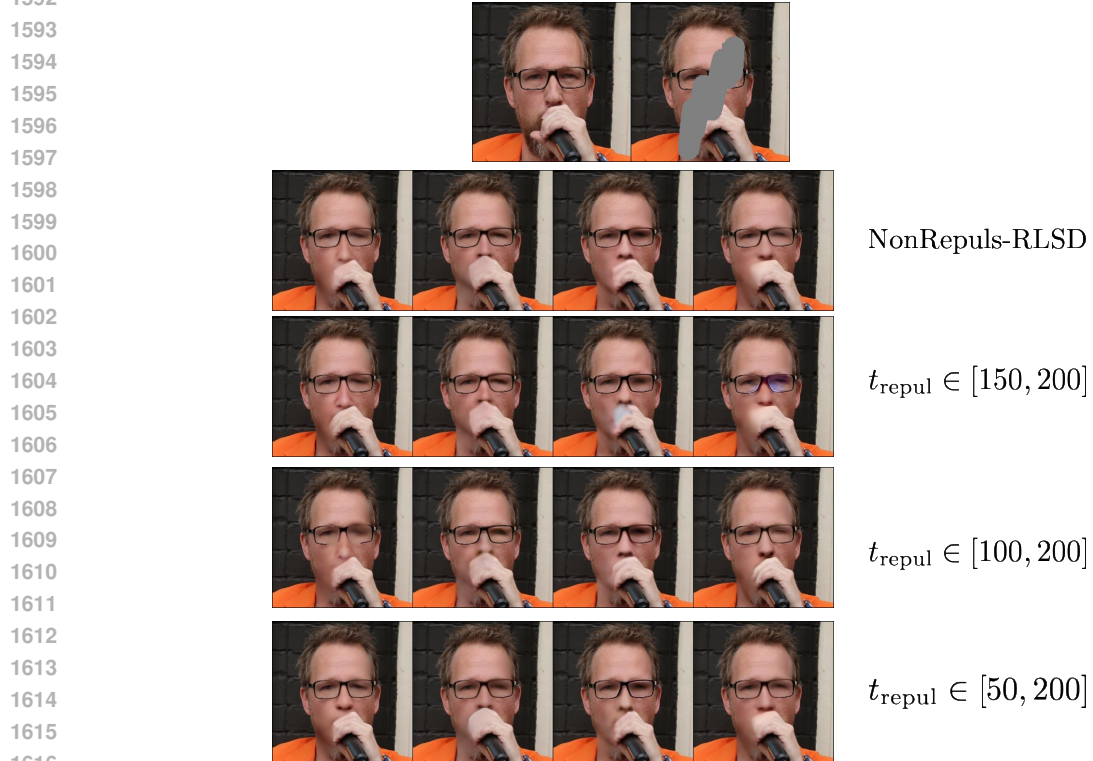
1547 **Effect of repulsion on a fixed interval.** Depending on the downstream task, we can include  
 1548 repulsion only on a fixed interval. For instance, this is the case of Phase Retrieval, where modes are  
 1549 discrete and isolated. Therefore, once each particle get around one of the mode (different), we can  
 1550 turn-off the repulsion, which yields a faster solver. For the case of inpainting, we show in Fig. 18 an  
 1551 ablation changing  $t_{\text{repu}}l$  where  $t \in [0, t_{\text{repu}}l]$  and in Fig. 19 when  $t \in [t_{\text{repu}}l, T]$  (with  $T = 200$ ).

1552 We observe that for diversity, it is more important to have repulsion at the beginnig of the sampling,  
 1553 which corresponds to the higher noise levels. Intuitively, it is more important to impose repulsion in  
 1554 the high-level semantics of the image instead of the fine-details. However, adding some repulsion  
 1555 later in the sampling process might improve some details (see the case  $t_{\text{repu}}l \in [100, 200]$  in Fig. 19).  
 1556

1557 **Effect of number of particles.** Lastly, we do an ablation when considering more particles in the  
 1558 repulsion term. Given four particles, we consider three settings:  $N = 0$  (four independent particles),  
 1559  $N = 2$  (two independent particles and two interacting particles) and  $N = 4$  (four interacting  
 1560 particles). We use again the free mask (Saharia et al., 2022), and  $L = 200$ . In Table 15 we show  
 1561 quantitative results, while in Figs. 20 and 21 we show two qualitative results. In both cases for  $N = 2$ ,  
 1562 the last two columns correspond to the two i.i.d. particles. Consequently, those two images can  
 1563 change when considering  $N = 4$ , while the first two columns can change when moving from  $N = 0$   
 1564 to two. In particular, in Fig. 20 with  $N = 2$  we incerase diversity, while in Fig. 21 we need at least  
 1565  $N = 4$ .



1590 Figure 18: Reconstruction for  $\gamma = 150$ , for  $L = 200$ ,  $N = 4$  and repulsion between 0 and  $t_{\text{repul}}$ : when  
1591  $t_{\text{repul}} = T$ , we have repulsion in all the trajectory, while  $t_{\text{repul}} = 0$  corresponds to NonRepuls-RLSD.  
1592



1617 Figure 19: Reconstruction for  $\gamma = 150$ , for  $L = 200$ ,  $N = 4$  and repulsion for  $t \in [t_{\text{repul}}, 200]$ ; notice that  
1618  $t_{\text{repul}} = T$  corresponds to NonRepuls-RLSD.  
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

Table 15: Ablation when changing the number of particles in the repulsion term with  $\sigma_v = 0.001$  - FFHQ 512.

Sampler	PSNR [dB] $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$	Diversity
NonRepuls-RLSD ( $N = 0$ )	<b>30.03</b>	0.064	<b>65.42</b>	0.002
RLSD ( $N = 2$ )	<u>29.99</u>	0.064	71.8	0.004
RLSD ( $N = 4$ )	29.99	<b>0.063</b>	69.91	0.005

Figure 20: Reconstruction when increasing the number of particles in the repulsion, for  $N = 0, 2$  and  $4$ . Notice that  $N = 0$  corresponds to the NonRepuls-RLSD,  $N = 2$  to Hybrid case, and  $N = 4$  to RLSD.

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

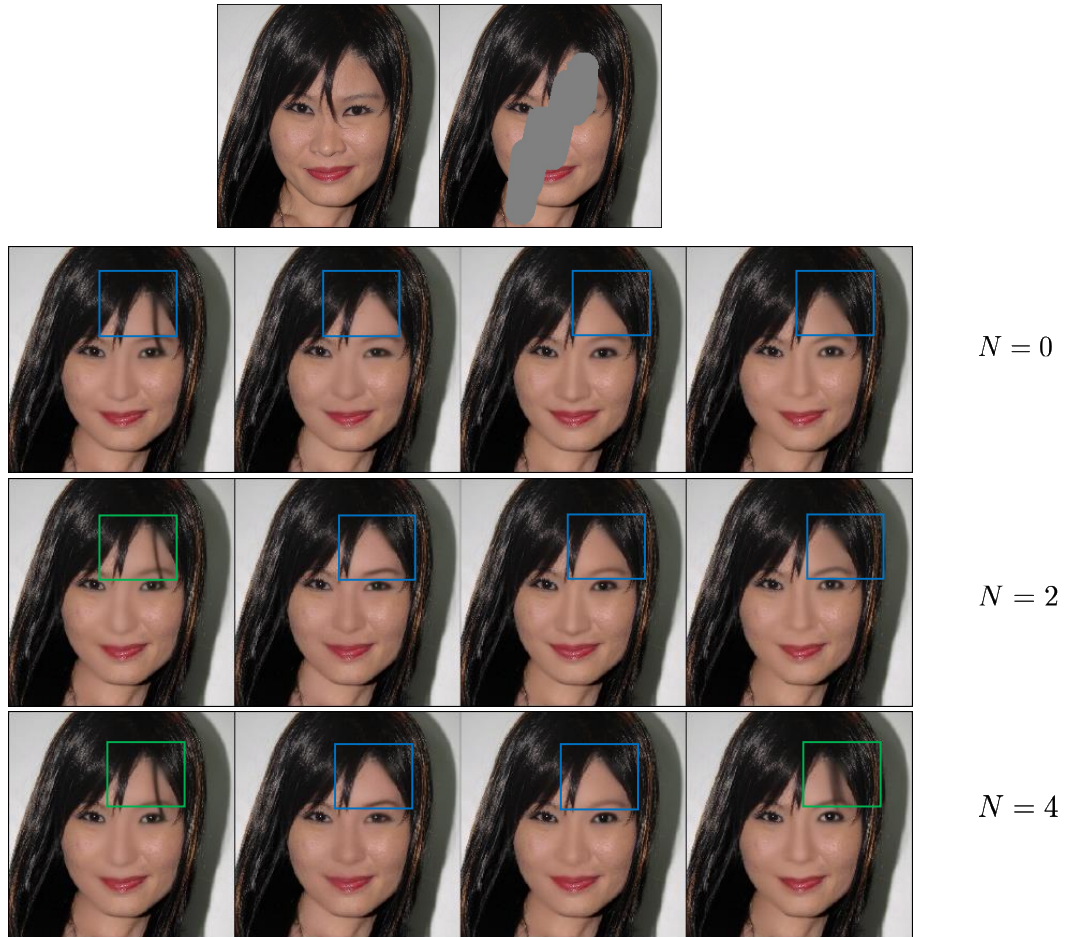
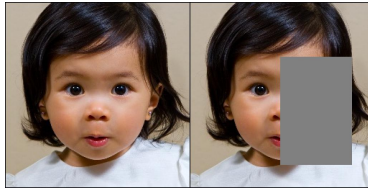


Figure 21: Reconstruction when increasing the number of particles in the repulsion, for  $N = 0, 2$  and  $4$ . Notice that  $N = 0$  corresponds to the NonRepuls-RLSD,  $N = 2$  to Hybrid case, and  $N = 4$  to RLSD. In this case, we need full repulsion to increase diversity w.r.t. the NonRepuls-RLSD case.

### D.8.3 COMPARISON BETWEEN NONAUG-RLSD VS RLSD

We compare in Fig. 22 NonAug-RLSD vs RLSD. Notice that RLSD achieves a sharper reconstructed image; in particular, the shirt for the case of NonAug-RLSD has not fine-detail, while RLSD has.



(a)



(b)

Figure 22: Comparison between RLSD (top row of Fig. b) and NonAug-RLSD (last row of Fig. b). RLSD achieves a sharper reconstructed image; in particular, the shirt for the case of NonAug-RLSD has not fine-detail, while RLSD has.

### D.8.4 COMPARISON WITH IMAGES FROM IMAGENET BETWEEN RLSD AND RED-DIFF WITH DIFFUSION PRIOR TRAINED ON FFHQ

Here we show some examples of reconstruction using out-of-distribution images; we use samples from ImageNet (Russakovsky et al., 2015). In Fig. 23 we show an example from ImageNet, where we compare RLSD with RED-diff using FFHQ. Clearly, the performance of RLSD is better than RED-diff. This is expected given that the diffusion model of RED-diff is with FFHQ. However, this demonstrates that using more powerful diffusion model as prior enables to deploy our model for different type of images.

### D.8.5 COUPLING PARAMETER

We ablate also the effect of the coupling parameter. We consider three values:  $\rho = \{0.03, 0.07, 0.12\}$ . The results for this case are shown in Fig. 24. As expected, increasing the coupling parameter yields a more diverse set of images. The reason why this happens is because the prior (through  $\mathbf{z}$ ) has more weight when estimating the image  $\mathbf{x}$ . However, this increase in diversity is penalized by a lower performance. On the other hand, a lower  $\rho$  generates four images that look very similar. It is important to note that our repulsion term is a better trade-off in terms of diversity/quality.

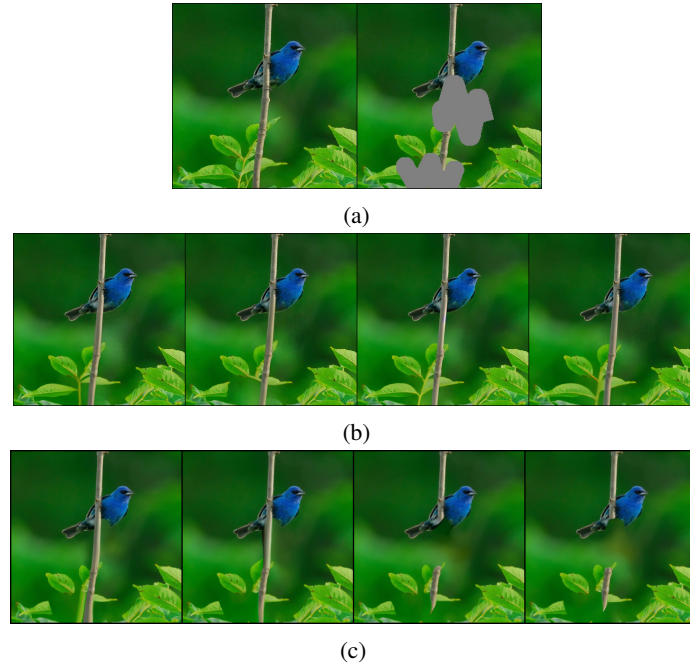
## D.9 LIMITATIONS AND FAILURE SETTINGS

Throughout the experiments, we demonstrated that RLSD outperforms other baselines, particularly PSLD. Additionally, we show that NonAug-RLSD is capable of solving nonlinear inverse problems at high resolutions. However, RLSD still faces challenges in certain settings and inverse problems. Below, we highlight some of these cases, which we aim to address in future work.

**Lack of fine details in ImageNet.** When performing inpainting on ImageNet, we observe that our method struggles to reconstruct fine details.

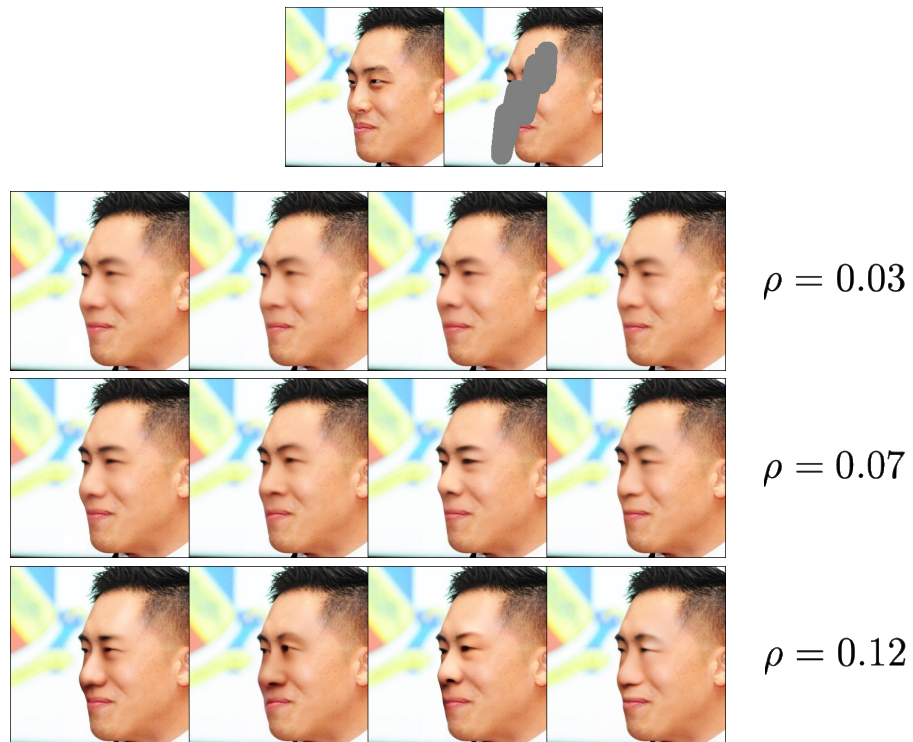


1782  
 1783  
 1784  
 1785  
 1786  
 1787  
 1788  
 1789  
 1790  
 1791  
 1792  
 1793  
 1794  
 1795  
 1796  
 1797  
 1798  
 1799  
 1800  
 1801  
 1802



1803 Figure 23: Results of the reconstruction of an image from the validation set of ImageNet using RLSD (b) and  
 1804 RED-diff with diffusion prior trained on FFHQ (c). This example showchases that using a large-pre trained  
 1805 model such as Stable Diffusion enables to use our method with images from very different classes (Imagenet  
 1806 and FFHQ).

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 Figure 24: Reconstruction as a function of  $\rho$  (coupling parameter). Increasing  $\rho$  enhances diversity but reduces  
 1832 quality: for  $\rho = 0.12$ , the 4 samples looks different but with a low quality, while with a lower  $\rho$ , the diversity is  
 1833 lower the performance is better.

1834  
 1835

1836  
1837  
1838  
1839  
1840  
1841  
1842  
1843  
1844  
1845



1846  
1847  
1848  
1849  
1850

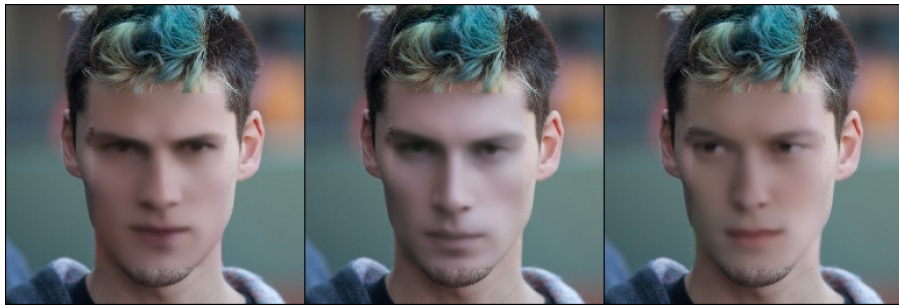
Figure 25: Example of a setting where our method struggles to have a high-quality reconstruction. In particular, RLSD generate images without small details: for instance, the building does not have windows in the walls.

1851  
1852  
1853  
1854  
1855

While this might be circumvent by increasing the weight of the prior, our simulations did not work.

**Full-face on FFHQ.** When we seek to reconstruct the face of a person, we observe that our method might struggle based on the background. For instance, the samples from Fig. 26 look blurry and unnatural. On the other hand, in Fig. 27 we observe a case that has a more natural reconstruction. Still, it has a lack of details in some parts of the face (eyebrow for example).

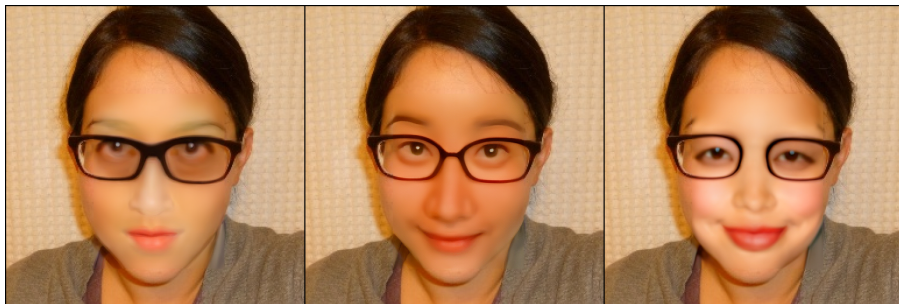
1856  
1857  
1858  
1859  
1860  
1861  
1862  
1863  
1864  
1865



1866  
1867  
1868

Figure 26: Example of a setting where our method struggles to have a high-quality reconstruction. In particular, RLSD generate a face that looks unnatural.

1869  
1870  
1871  
1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879



1880  
1881  
1882

Figure 27: Example of a more natural reconstruction for full-face example. In particular, RLSD generate a face that looks more natural than Fig. 26. Still, the eyebrow looks blurry.

1883

When applying the same box to ImageNet, we observe the same as it is shown in Fig. 28.

1884  
1885  
1886  
1887

**Repulsion too high.** Lastly, an important setting where RLSD fails is when  $\gamma$  is too high. As an example, we show in Fig. 29 a case where  $\gamma = 200$ , and the same setting as in Appendix D.8. While in Appendix D.8.2 we illustrated that increasing  $\gamma$  enhances diversity, if  $\gamma$  is too high, then the reconstruction fails.

1888  
1889

**Fixed resolution given by the diffusion prior.** The dimension of the images generated with the diffusion prior is constrained by the existing pretrained diffusion priors (currently, Stable Diffusion

1890  
1891  
1892  
1893  
1894  
1895  
1896  
1897  
1898  
1899  
1900

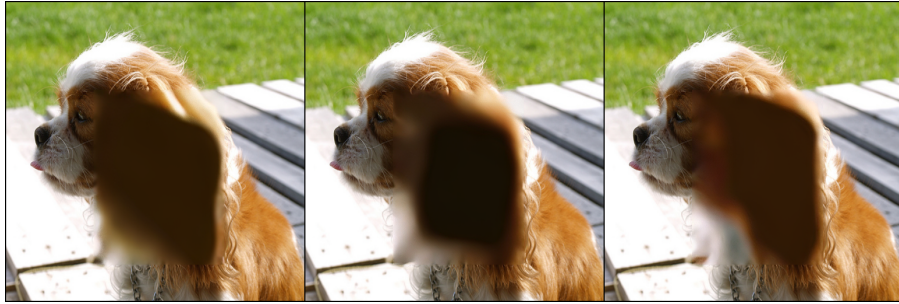


Figure 28: Example from ImageNet, where we see that there is a lack of details.

1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909  
1910  
1911  
1912  
1913  
1914



Figure 29: Example of failure when considering  $\gamma$  too high. We consider  $N = 4$  (the second row),  $L = 200$ , and the same setting as in Appendix D.8.

1915  
1916  
1917  
1918  
1919  
1920  
1921  
1922

provides resolutions up to  $512 \times 512$ ). This is still considered a high dimensional data, compared with existing methods that work with resolutions up to  $256 \times 256$ . Extending our framework to handle higher dimensions than the one of the diffusion prior is an interesting research direction, where ideas from compositional generation might be helpful. We leave this as future work.

1923  
1924  
1925  
1926  
1927  
1928  
1929  
1930

## D.10 TRADE-OFF BETWEEN DIVERSITY AND QUALITY IN UNCONSTRAINED GENERATION

To demonstrate the generality of RLSD and our particle-based variational approximation, we also include additional experiments when considering unconstrained sampling: we include results for the unconstrained case, i.e., text-to-image and text-to-3D. Details about Score distillation and the methods used in this section can be found in Appendix E.

1931  
1932  
1933  
1934  
1935

### D.10.1 TOY EXAMPLE WITH A BIMODAL GAUSSIAN DISTRIBUTION

We consider here a toy example to showcase how the  $\gamma$  parameter (the amount of repulsion) dictates the trade-off between diversity and quality. We consider a mixture of two Gaussians of parameters  $\mathcal{N}_1([1, 0]^\top, 0.005\mathbf{I})$  and  $\mathcal{N}_2([-1, 0]^\top, 0.005\mathbf{I})$ , and we consider two settings:

1936  
1937  
1938  
1939  
1940

1. Two independent Gaussians with  $\sigma \rightarrow 0$  where we fit the mean parameters
2. Two dependent Gaussians with  $\sigma \rightarrow 0$  where we fit the mean parameters and coupled via an Euclidean RBF kernel.

1941  
1942  
1943

For each setting, we compute 200 realizations with the same seed. In Fig. 30a we show how many realizations suffers from mode collapse ( $\gamma = 0$  corresponds to the first setting), while in Figs. 30b and 30c we show a realization for  $\gamma = 1$  and  $\gamma = 2000$ , where we observe that the "quality" of the estimation for this high value is poorer compared to  $\gamma = 1$ .



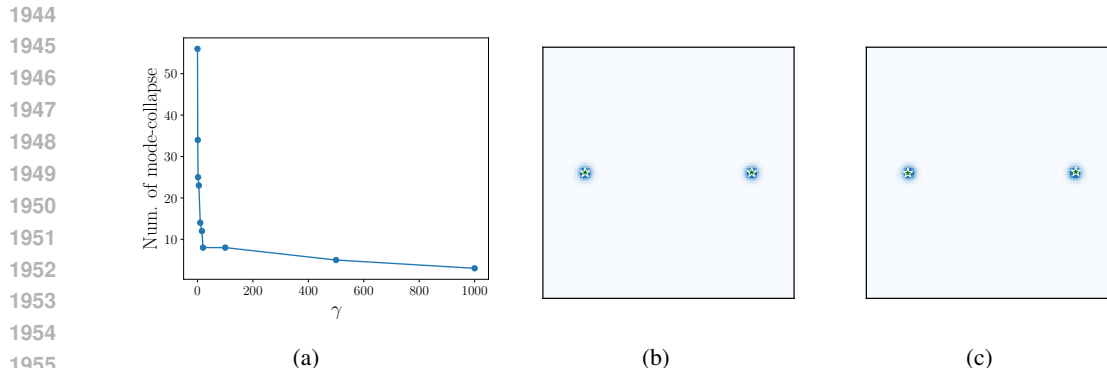


Figure 30: a) Number of realizations that has a particle collapse to the same mode as a function of the amount of repulsion that we consider. Estimation when considering b)  $\gamma = 1$ , c)  $\gamma = 2000$ .

#### D.10.2 TEXT-TO-IMAGE GENERATION: TRADE-OFF PLOT AND ABLATION OF KERNEL FUNCTION

**Implementation details.** Here we describe the details of the experiments on 2D images with ProlificDreamer (VSD) (Wang et al., 2024), which is described in Appendix E. VSD is a particular case of score distillation for unconstrained sampling. We optimize 500 optimization steps, follow their setup, and we train a U-Net from scratch to estimate the variational score; we used the code from the public repository [https://github.com/yuanzhi-zhu/prolific\\_dreamer2d](https://github.com/yuanzhi-zhu/prolific_dreamer2d). We consider ADAM optimizer and set the learning rate of particle images is 0.03 and the learning rate of U-Net is 0.0001. We consider only 4 particles as it is enough to promote diversity (in particular when considering DINO as feature extractor). The images and parameters are initialized at random. We run all the experiments in a single NVIDIA A100 GPU of 80GB.

**Trade-off plot.** We study the trade-off between diversity and quality when increasing the amount of repulsion. We consider 75 images from the COCO dataset and compute average qualities and diversities across all images. We also include a comparison with stochastic sampling using the Euler discretization of the backward process, with 30 steps (denoted by ‘Ancestral’). This serves as an upper bound on the performance of distillation techniques. We consider 4 particles for both the distillation optimization and the sampling via Euler. The results are shown in Fig. 31, where we compute FID (lower is better), Aesthetic (higher is better) and CLIP (higher is better) scores as a function of a diversity score (27); a higher diversity score corresponds to more diverse image generation. For the plot, we fix all the hyperparameters, and we sweep  $\gamma$  between 0 (ProlificDreamer) and 40, and we consider Ancestral sampling as an upper bound in terms of trade-off. Clearly, adding repulsion increases the diversity, while slightly decreasing the quality metrics for low values of  $\gamma$ . However, for these lower values of  $\gamma$ , the diversity is still far from the stochastic sampling. When increasing the values of  $\gamma$  above 30, we close the gap in terms of diversity at the cost of decreasing the quality (FID and Aesthetic) as well as the text-alignment (CLIP).

**Using different domains/distances for the repulsion force** We compare the RBF when considering different domains, namely Euclidean, DINO and LPIPS.

- RBF in the Euclidean domain.** We compare between DINO and RBF for a repulsion with  $\gamma = 10$  with two qualitative examples in Figs. 32 from COCO validation set (Lin et al., 2014) constructed in (Jain et al., 2022).
- RBF using LPIPS.** We also tried LPIPS as metric in RBF. However, we observe that when combined with ProlificDreamer, the generated images have some artifacts or are not well aligned with the text prompt; two examples are shown in Figs. 33.

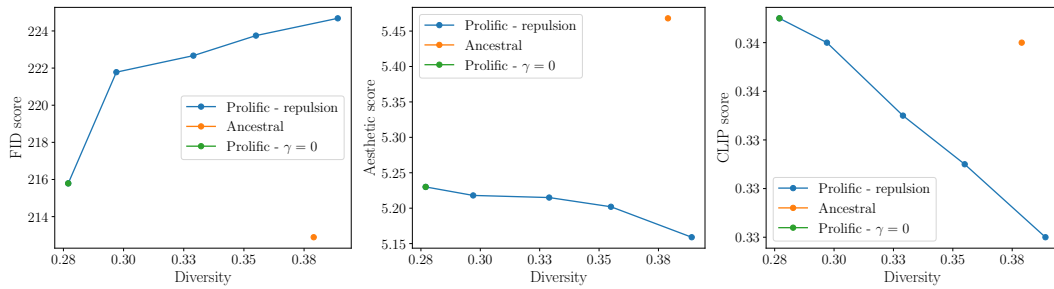
1998  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007

Figure 31: We consider the trade-off between diversity and quality for text-to-2D image generation using score distillation. We consider  $\gamma = \{0, 10, 20, 30, 40\}$ . From left to right: Diversity (27) vs left) FID, middle) Aesthetic score and right) CLIP score. We consistently observe that adding repulsion increases diversity, while slightly decreasing the quality of the images. Notice that we need a repulsion with  $\gamma > 30$  to reach the diversity level of stochastic sampling using Euler.

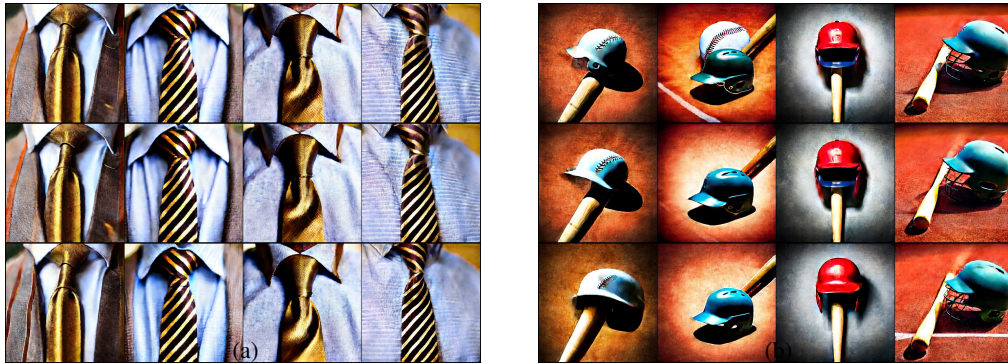
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025

Figure 32: Comparison between RBF with Euclidean distance and DINO. We generate samples using the prompt a) "a gold tie is tied under a brown dress shirt with stripes..", b) "a baseball bat with a batting helmet upsidedown.". From the top to bottom:  $\gamma = \{0, 10(\text{Euclidean}), 10(\text{DINO})\}$  with CFG = 7.5 and 500 steps.

2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037

Figure 33: Comparison between RBF with LPIPS (top row) and DINO (bottom row). We generate samples using the prompt a) "a gold tie is tied under a brown dress shirt with stripes..", b) "a baseball bat with a batting helmet upsidedown.". From the top to bottom:  $\gamma = \{10(\text{LPIPS}), 10(\text{DINO})\}$  with CFG = 7.5 and 500 steps.

2042  
2043

### D.10.3 TEXT-TO-3D GENERATION

2044  
2045  
2046

Lastly, we demonstrate how our method improves the mode collapse phenomena in text-to-3D generation. Our main motivation is to show that including the repulsion entails a more diverse set of scenes when changing the seed.

2047  
2048  
2049  
2050  
2051

**Implementation details.** We consider DreamFusion (Poole et al., 2022) as base method and the implementation from the Threestudio framework (Guo et al., 2023). All 3D models are optimized for 10000 iterations using Adam optimizer with a learning rate of 0.01, and we use the same configuration as DreamFusion for the rendering. For the NeRF architecture we use an MLP, and we use the same configuration from the setting in Threestudio. For the diffusion model, we use DeepFloyd-IF-XL-v1.0



with a guidance scale of 20. For the repulsion, we use a repulsion weight of  $\gamma = 200\sigma_t$ . We use one Nvidia A100 of 80GB, and we consider a batch of 20 particles; this presents a limitation for using more computationally-demand methods such as ProlificDreamer. For the RBF, we use DINO.

**Results.** In Figs. 34 and 39 we show two qualitative examples. For each case, we show two different views of the generated object. Furthermore, we include the corresponding .gif files for the prompt "An ice cream sundae" in Figs. 31 to 34. Notice that with repulsion we can generate two different colors of glass for the ice cream (dark and white), something that we could not achieve without repulsion. Furthermore, while the case without repulsion generates two scenes that have a similar perspective, our method generates two scenes that show the ice cream at different distances. Lastly, notice that the quality of the generated scene is bounded by the performance of the base method (in this case, DreamFusion). Therefore, the results look saturated, similar to what happen with SDS.

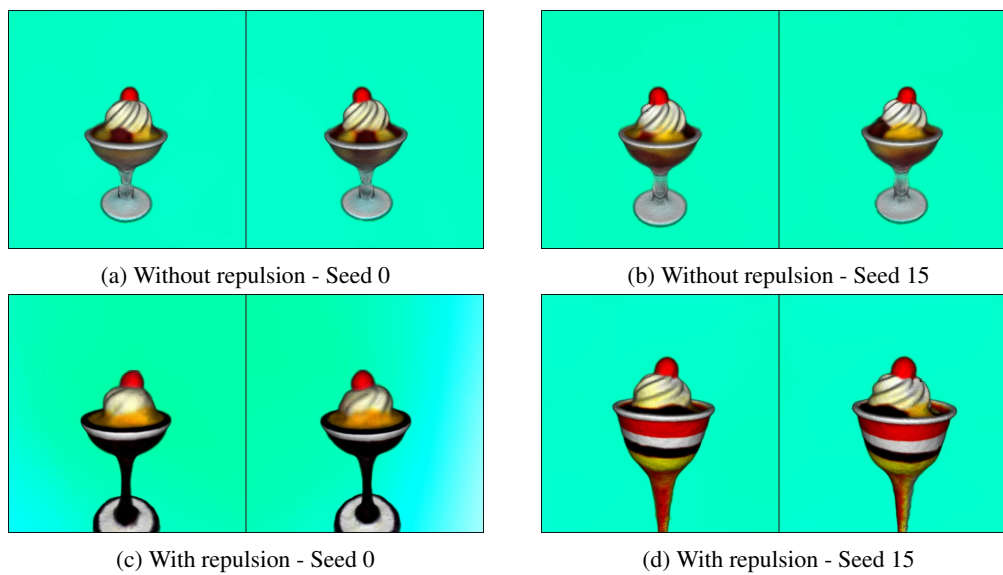


Figure 34: Text-to-3D generation using DreamFusion with the prompt "An ice cream sundae", and considering a batch of 20 samples (particles) a, b) without repulsion with seed 0 and seed 15, c, d) with repulsion with seed 0 and 15 respectively. For each case, we show two different views of the object. Clearly, the two cases without repulsion look very similar, generating the same type of ice cream sundae. On the other hand, adding repulsion increase diversity of the scene (a different glass, and color). However, this comes at the cost of less details in the ice cream.

Figure 35: Without repulsion - Seed 0

Figure 36: With repulsion - seed 0

Figure 37: Without repulsion - seed 15

Figure 38: With repulsion - seed 15

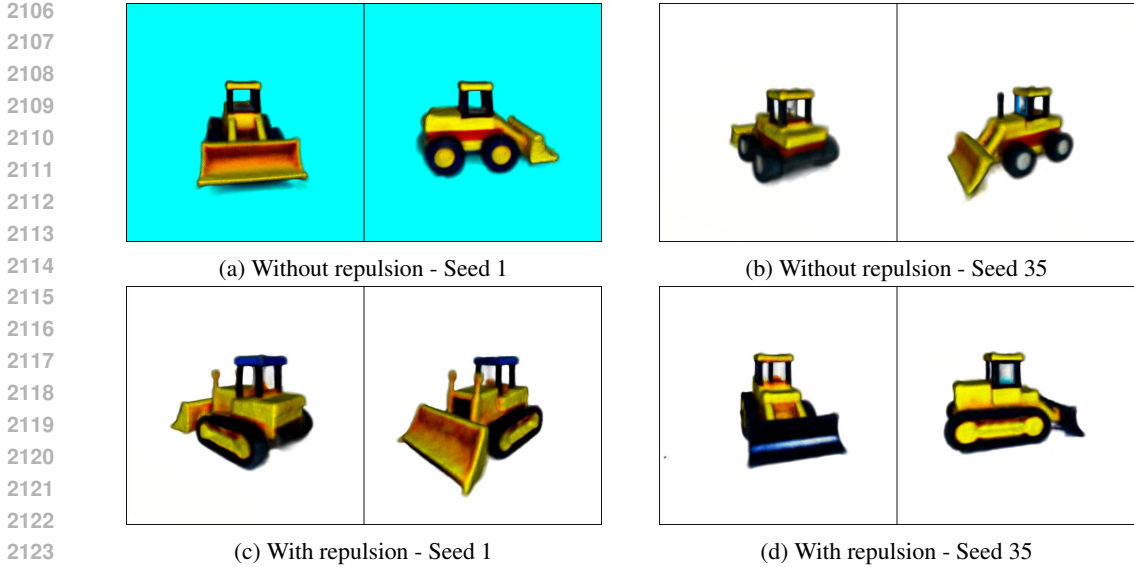


Figure 39: Text-to-3D generation using DreamFusion with the prompt "a bulldozer made out of toy bricks", and considering a batch of 20 samples (particles) a, b) without repulsion with seed 1 and seed 35, c, d) with repulsion with seed 1 and 34 respectively. For each case, we show two different views.

## E SCORE DISTILLATION

Score distillation sampling (Poole et al., 2022) was the first proposed method to optimize a generator by using distillation. Consider a pretrained score function  $\epsilon_{\theta}(\mathbf{x}_t, t) \approx -\sigma_t \nabla_{\mathbf{x}_t} \log p_{\theta}(\mathbf{x}_t)$  representing a distribution of interest  $p_{\theta}$ . The idea of *score distillation sampling* (SDS) is to train a generator  $g_{\phi}$ , such that the output of the generator given an input  $\mathbf{m}$  is a sample  $\mathbf{x}_0 = g_{\phi}(\mathbf{m}) \sim p_{\theta}$ . This is achieved by optimizing the distillation loss (Poole et al., 2022)

$$\mathcal{L}_{\text{SDS}}(\mathbf{x}_0 = g_{\phi}(\cdot)) = \mathbb{E}_{t \sim \mathcal{U}[0, T], \epsilon \sim \mathcal{N}(0, \mathbf{I})} \left[ \omega(t) \frac{\sigma_t}{\alpha_t} \text{KL} (q(\mathbf{x}_t | \mathbf{x}_0 = g_{\phi}(\cdot)) || p_{\theta}(\mathbf{x}_t)) \right], \quad (28)$$

where  $\omega(t)$  is a weighting function and  $q(\mathbf{x}_t | \mathbf{x}_0 = g_{\phi}(\cdot))$  is the variational distribution; when  $g_{\phi}(\cdot) = \phi$ , i.e., identity mapping, then we are in the case our proposed method in Section 4.

Although its remarkable success in generating 3D scenes, these methods suffer from a mode collapse problem. This is mainly driven by the choice of the (reverse) KL divergence as the loss in (28) and the fact that a *unimodal variational distribution*  $q$  is used. In particular, they need to consider a high CFG in the classifier-free guidance score (Ho and Salimans, 2021)

$$\epsilon_{\theta}^w(\mathbf{z}_t; c, t) = \epsilon_{\theta}(\mathbf{z}_t; c = \emptyset, t) + w(\epsilon_{\theta}(\mathbf{z}_t; c, t) - \epsilon_{\theta}(\mathbf{z}_t; c = \emptyset, t)) \quad (29)$$

where  $\emptyset$  indicates a null condition and the  $w$  is the CFG weight; they consider  $w = 100$ .

The mode collapse phenomenon is clear when observing SDS through the lens of gradient flows 4.1, which entails a mode-seeking optimization. Therefore, if we want to avoid mode collapse, we need to modify either the loss function (the divergence) or the variational approximation. While the former renders intractable formulations, the latter can be modified easily. Therefore, previous works focused on modifying the variational distribution. See, for instance, ProlificDreamer (Wang et al., 2024), and more recently, in Luo et al. (2024); also other works (Katzir et al., 2024; Wang et al., 2023) have studied this problem, and proposed alternative approaches to circumvent this problem. In ProlificDreamer, the authors consider a particle approximation of the variational distribution by randomizing the parameters  $\theta$ , which yields the following update rule

$$\nabla_{\theta} \mathcal{L}_{\text{VSD}} = \mathbb{E}_{t, \epsilon, c} \left[ \omega(t) (\epsilon_{\phi}(\mathbf{x}_t; c, t) - \sigma_t \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t | c)) \frac{\partial \mathbf{x}_0}{\partial \theta} \right]. \quad (30)$$

Notice that  $\sigma_t \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t | c)$  is unknown, so they fine-tune a pre-trained diffusion model for each particle (which represents each rendered image). Therefore, they incorporate a second optimization

2160 problem that minimizes the DSM loss; they parametrize the score network using LoRA (Hu et al.,  
2161 2021). However, adding an auxiliary neural network does not necessarily promote diversity.

## 2163 E.1 GRADIENT FLOW PERSPECTIVE OF SCORE DISTILLATION

2164 We can interpret the variational diffusion sampling optimization procedure as a Wasserstein gradient  
2165 flow when we constraint to Bures–Wasserstein manifold.

2166 **Wasserstein gradient flow.** WGF describes how probability distributions change over time by  
2167 minimizing a functional on  $\mathcal{P}(\mathbb{R}^n)$ , representing the space of probability distributions over  $\mathbb{R}^n$  with  
2168 finite second moments. Denoted as  $(\mathcal{P}(\mathbb{R}^n), W_2)$ , this space employs the Wasserstein-2 distance as  
2169 its metric, termed the Wasserstein space. Before delving into the WGF, we define the Wasserstein  
2170 gradient of a functional  $\mathcal{F}(q)$  as

$$2171 \nabla_{W_2} \mathcal{F}(q) = \nabla_{\mathbf{x}} \frac{\delta \mathcal{F}(q)}{\delta q}. \quad (31)$$

2172 where  $\frac{\delta \mathcal{F}(q)}{\delta q} = \lim_{\epsilon \rightarrow 0} \frac{\mathcal{F}(q + \epsilon \sigma) - \mathcal{F}(q)}{\epsilon}$  is the first variation defined for any direction in the tangent  
2173 space of  $\mathcal{P}$ . Given this definition and two boundary conditions  $\rho_0 = q_0(\mathbf{x})$  and  $\rho_\infty = p(\mathbf{x})$ , we can  
2174 define a path of densities  $q_t$  where its evolution is described by the Liouville equation (also known as  
2175 continuity equation)

$$2176 \frac{\partial q_\tau}{\partial \tau} = \operatorname{div}(q_\tau \nabla_{W_2} \mathcal{F}(q_\tau)) \quad (32)$$

2177 At the particle level, for a given particle  $\mathbf{x}_\tau \sim q_\tau$  in  $\mathbb{R}^n$ , the gradient flows defines a dynamical system  
2178 drive a vector a field  $\{v_\tau\}_{\tau \geq 0}$  in the Euclidean space  $\mathbb{R}^n$  given by

$$2179 d\mathbf{x}_\tau = v_\tau(\mathbf{x}_\tau) d\tau = -\nabla_{W_2} \mathcal{F}(q_t)(\mathbf{x}_\tau) d\tau.$$

2180 Therefore, this ODE describes the evolution of the particle  $\mathbf{x}_\tau$  where the associated marginal  $q_\tau$   
2181 evolves to decrease  $\mathcal{F}(q_\tau)$  along the direction of steepest descent according to the continuity equation  
2182 in (32).

2183 Notice that the WGF is defined in a continuous domain for  $\tau$ . We can discretized via the following  
2184 movement minimization scheme with step size  $h$ , also known as the Jordan-Kinderlehrer-Otto (JKO)  
2185 scheme (Jordan et al., 1998),

$$2186 q_{k+h} = \operatorname{argmin}_{q \in \mathcal{P}(\mathbb{R}^n)} \left\{ \mathcal{F}(q) + \frac{1}{2h} W_2^2(q, q_k) \right\}, \quad (33)$$

2187 Notice that the JKO scheme has two terms: the first seeks to minimize the functional  $\mathcal{F}(q)$  and the  
2188 second is a regularization term, that penalize to stay close to  $q_k$  in Wasserstein-2 distance as much  
2189 as possible. It can be shown that as  $h \rightarrow 0$ , the limiting solution of (33) coincides with the path  
2190  $\{q_\tau\}_{\tau \geq 0}$  defined by the continuity equation (32).

2191 Throughout this work, we consider the KL divergence as the function, i.e.,  $\mathcal{F}_{\text{kl}}(q) = \text{KL}(q||p)$ . Then,  
2192 the Liouville equation boils down to the Fokker-Planck equation

$$2193 \frac{\partial q_\tau}{\partial \tau} = \operatorname{div}(q_\tau (\nabla_{\mathbf{x}} \log q_\tau - \nabla_{\mathbf{x}} \log p)) \quad (34)$$

2194 where the Wasserstein gradient is  $\nabla_{W_2} \mathcal{F}_{\text{kl}}(q_\tau) = \nabla_{\mathbf{x}} \log(q_\tau/p)$  and the probability flow ODE  
2195 follows (6). Lastly, although here we focus on the Wasserstein metric, we can consider other metrics  
2196 that yield different flows (Chen et al., 2023). This is a future avenue to explore.

2197 **Bures-Wasserstein gradient flow.** We now show that RED-diff (Mardani et al., 2024) can be de-  
2198 rived from the gradient flow perspective when considering the flow in in the Bures-Wasserstein space  
2199  $(\mathcal{BW}(\mathbb{R}^n), W_2)$ , i.e., the subspace of the Wasserstein space consisting of Gaussian distributions.

2200 The Wasserstein-2 distance between two Gaussian distributions  $q = \mathcal{N}(\mu_q, \Sigma_q)$  and  $p = \mathcal{N}(\mu_p, \Sigma_p)$   
2201 has a closed form,

$$2202 W_2^2(q, p) = \|\mu_q - \mu_p\|_2^2 + \mathcal{B}^2(\Sigma_q, \Sigma_p),$$

where  $\mathcal{B}^2(\Sigma, \Sigma_p) = \text{tr}\left(\Sigma + \Sigma_p - 2\left(\Sigma^{\frac{1}{2}}\Sigma_p\Sigma^{\frac{1}{2}}\right)^{\frac{1}{2}}\right)$  is the squared Bures distance. By restricting the JKO scheme in (33) to the Bures-Wasserstein space, the authors in Lambert et al. (2022) showed that the discretization entails a solution given by the limiting curve  $\{q_t : \mathcal{N}(\mu_t, \Sigma_t)\}_{t \geq 0}$ . Therefore, the gradient flow of the KL divergence in the Bures-Wasserstein space boils down to the evolution of the means and covariance matrices of Gaussians, described by the following ODEs:

$$\frac{d\mu_\tau}{d\tau} = \mathbb{E}_{\mathbf{x} \sim q_\tau} \left[ \nabla_{\mathbf{x}} \log \frac{p(\mathbf{x})}{q_\tau(\mathbf{x})} \right], \quad (35)$$

$$\frac{d\Sigma_\tau}{d\tau} = \mathbb{E}_{\mathbf{x} \sim q_\tau} \left[ \left( \nabla_{\mathbf{x}} \log \frac{p(\mathbf{x})}{q_\tau(\mathbf{x})} \right)^T (\mathbf{x} - \mu_\tau) \right] + \mathbb{E}_{\mathbf{x} \sim q_\tau} \left[ (\mathbf{x} - \mu_\tau)^T \nabla_{\mathbf{x}} \log \frac{p(\mathbf{x})}{q_\tau(\mathbf{x})} \right]. \quad (36)$$

In essence, RED-diff (and DreamFusion as well) follows the ODE corresponding to the mean with the diffusion as regularizer and the likelihood as measurement term. The key of the proposed methods is that the variational distribution uses the diffused trajectory (and therefore, the denoiser) as regularizer.

## F COMPARISON WITH OTHER PARTICLE METHODS USING DIFFUSION MODELS

### F.1 STEIN VARIATIONAL GRADIENT DESCENT (SVGD)

SVGD is a deterministic particle-based variational inference method (Liu, 2017), and is the core method of (Kim et al., 2023). Through the lens of gradient flow, it optimizes the same functional (KL divergence) but considers a different metric induced by the Stein operator (Duncan et al., 2023). In particular, given a pairwise repulsion as  $\mathcal{R}(\mathbf{z}_{t,\tau}^{(1)}, \dots, \mathbf{z}_{t,\tau}^{(N)}) = k(\mathbf{z}_{t,\tau}^{(i)}, \mathbf{z}_{t,\tau}^{(j)})$ , the update direction in SVGD for the particle  $\mathbf{z}_{t,\tau}^{(i)}$  is given by

$$\sum_{j=1}^N k(\mathbf{z}_{t,\tau}^{(i)}, \mathbf{z}_{t,\tau}^{(j)}) \nabla_{\mathbf{z}_{t,\tau}^{(j)}} \log p(\mathbf{z}_{t,\tau}^{(j)}) - \nabla_{\mathbf{z}_{t,\tau}^{(j)}} k(\mathbf{z}_{t,\tau}^{(j)}, \mathbf{z}_{t,\tau}^{(j)}). \quad (37)$$

When comparing the gradient update (37) with (7), the main difference resides in how is computed the gradient of the log target: while in SVGD, all the ensembles follow the same averaged gradient direction weighted by the similarity kernel function, in our method each particle uses its score direction. This difference in the score has important implications: SVGD suffers from the curse of dimensionality and poor exploration of the space in high-dimensional multimodal distributions, suffering the mode collapse phenomenon. Some alternatives have been proposed, like adding an annealing schedules (D’Angelo and Fortuin, 2021a; Ba et al., 2021), but the method still discovers only a few modes in multimodal distributions. Conversely, the repulsive method achieves a better exploration, and therefore, it finds more modes.

### F.2 PARTICLE GUIDANCE

In the context of sampling, a recent work termed Particle guidance (Corso et al., 2024) proposed a method to incorporate an interactive particle sampler based on diffusion. They proposed a guidance term that couples a set of particles to guide the backwards diffusion process towards different modes of the target distribution. Although they leverage a similar idea, the scope of our work is different. Here we focus on diffusion model in the context of distillation, and therefore, as a regularizer, which allows to deploy in constrained problems like inverse problems, or unconstrained cases like text-to-3D.

### F.3 COPULA MODELS

Copula models Nelsen (2006); Tran et al. (2015) are methods for defining variational distributions that incorporate dependencies between variables. While typically this is done for modeling the relationship between each dimension, in this case, we use the copula model to define dependencies between a batch of particles and promote diversity via repulsion.

## G VARIATIONAL SAMPLERS FOR INVERSE PROBLEMS

Variational inference methods Blei et al. (2017) are one of the most important techniques for sampling from intractable distributions. In a nutshell, VI defines a parametric distribution where the parameters are learned via stochastic gradient methods. A key difference between variational samplers and others, such as those based on the Monte Carlo Markov Chain, is that VI defines a parametric distribution. Thus, once the parameters are learned, we can easily generate samples.

In inverse problems, different variational samplers have been proposed in Knoll et al. (2011); Portilla et al. (2003); Kobler et al. (2017). The difference between these approaches lies in the definition of the variational distribution. For instance, Kobler et al. (2017) leverages the algorithm-unrolling framework, while Portilla et al. (2003) uses a mixture of Gaussians in the wavelet domain. More recently, the authors in Tonolini et al. (2020) proposed a two-stage semi-supervised method for learning variational samplers: a first stage that learns the forward model and a second stage that trains a conditional variational auto-encoder that learns to solve the inverse problem.

Related to our work, both Feng et al. (2023) and Mardani et al. (2024) combine variational inference with diffusion priors. Our method has three main differences with Mardani et al. (2024): 1) we define our variational inference problem in a latent space, allowing us to exploit the latent diffusion model, 2) we consider a multi-modal variational distribution with a repulsion term to couple the particles, enabling a better exploration, and 3) we decouple the data and prior term to handle the challenges of latent inversion. On the other hand, Feng et al. (2023) also considers a variational perspective using diffusion priors. First, they consider pixel-based diffusion models, which simplifies the formulation. Second, they incorporate the diffusion prior by computing the log probabilities: this requires solving the underlying ODE and estimating a divergence, which is computationally expensive. Lastly, they consider a normalizing flow as a variational distribution, i.e., they define a distribution  $q_\phi(\mathbf{x}|\mathbf{y})$  parameterized by  $\phi$ ; in particular, they consider RealNVP (Dinh et al., 2017). While this differs from our definition in (8), our framework allows incorporating a normalizing flow as a variational distribution. Exploring this type of parametric distribution opens an interesting research direction, where we define an amortized distribution that can generate solutions in only a few steps; we leave this as future work. A pseudo-code version of this method is shown in Alg 3.

---

**Algorithm 3** Score-Based Diffusion Models as Principled Priors for Inverse Imaging Feng et al. (2023)

---

```

Require:  $\mathbf{y}, f(\cdot), s_\theta(\mathbf{x}_t, t)$ 
for  $n = 1$  to  $N_{iter}$  do
    Compute  $\log p_\theta(\mathbf{x}_0) = \log p_T(\mathbf{x}_T) + \int_0^T \nabla \cdot (-\frac{1}{2}\beta(t)\mathbf{z}_t dt - \beta(t)s_\theta(\mathbf{x}_t, t)) dt$  via ODE solver
    and divergence estimation
    Minimize  $D_{KL}(q_\phi \| p_\theta(\cdot | \mathbf{y}))$ 
end for
return  $q_\phi(\mathbf{x}|\mathbf{y})$ 

```

---

Lastly, it is important to note the connection of RLSD with plug-and-play methods Kamilov et al. (2023); Zhang et al. (2021). However, none of these works incorporate latent diffusion model as denoisers.