

FREE HUNCH: DENOISER COVARIANCE ESTIMATION FOR DIFFUSION MODELS WITHOUT EXTRA COSTS

Severi Rissanen
Aalto University
severi.rissanen@aalto.fi

Markus Heinonen
Aalto University

Arno Solin
Aalto University

ABSTRACT

The covariance for clean data given a noisy observation is an important quantity in many training-free guided generation methods for diffusion models. Current methods require heavy test-time computation, altering the standard diffusion training process or denoiser architecture, or making heavy approximations. We propose a new framework that sidesteps these issues by using covariance information that is *available for free* from training data and the curvature of the generative trajectory, which is linked to the covariance through the second-order Tweedie’s formula. We integrate these sources of information using (i) a novel method to transfer covariance estimates across noise levels and (ii) low-rank updates in a given noise level. We validate the method on linear inverse problems, where it outperforms recent baselines, especially with fewer diffusion steps.

1 INTRODUCTION

Diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song et al., 2021) have emerged as a robust class of generative models in machine learning, adept of producing high-quality samples across diverse domains. These models function by progressively denoising data through an iterative process, learning to reverse a predefined forward diffusion process that systematically adds noise. Conditional generation extends the capabilities of diffusion models by allowing them to generate samples based on specific input conditions or attributes. This conditioning enables more controlled and targeted generation, making diffusion models applicable to a wide range of tasks such as text-to-image synthesis or linear inverse problems such as deblurring, inpainting, or super-resolution.

A strand of recent research has concentrated on applying pretrained diffusion models to accommodate user-defined conditions, enhancing the flexibility and control of a single model to an arbitrary number of tasks. These methods guide the sampler towards regions whose denoisings $p(x_0 | x_t)$ are compatible with the condition or constraint, which requires efficient denoising mean $\mathbb{E}[x_0 | x_t]$ and covariance $\text{Cov}[x_0 | x_t]$ estimates (Ho et al., 2022; Song et al., 2023a;b; Boys et al., 2023; Peng et al., 2024). While estimating the mean is straightforward through the denoiser, accurately determining the covariance has proven more challenging. Consequently, efficient approaches have been proposed with heavy approximations (Chung et al., 2023; Song et al., 2023a).

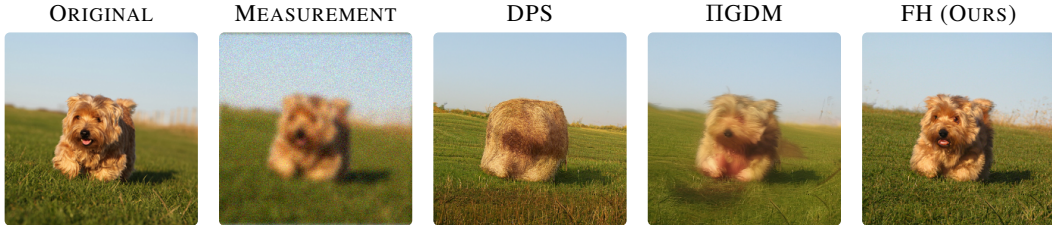


Figure 1: Comparison of different conditional diffusion methods for deblurring, with a low number of solver steps (15 Heun iterations). DPS (Chung et al., 2023) and IIGDM (Song et al., 2023a) work well with many steps, but accurate covariance estimates matter more for small step counts.

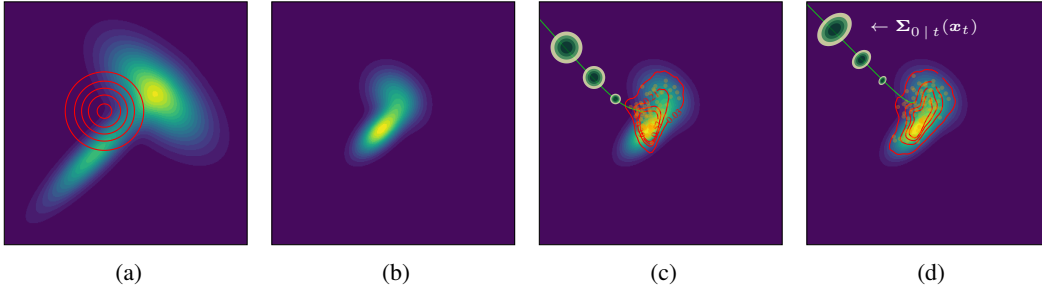


Figure 2: (a) A distribution $p(\mathbf{x}_0)$ represented by a pretrained diffusion model, and a Gaussian likelihood $p(\mathbf{y} | \mathbf{x}_0)$. (b) The (exact) posterior $p(\mathbf{x}_0 | \mathbf{y}) \sim p(\mathbf{x}_0)p(\mathbf{y} | \mathbf{x}_0)$. (c) Generated samples from a model with a heuristic diagonal denoiser covariance $\Sigma_{0 | t}(\mathbf{x}_t)$, and a generative ODE trajectory with approximated $p(\mathbf{x}_0 | \mathbf{x}_t)$ shapes represented as ellipses along the trajectory. (d) Generated samples with our denoiser covariance.

In this paper, we propose a new method for denoiser covariance estimation, which we refer to as *Free Hunch* (FH). The name stems from the core insight that much of the required guiding covariance information is, in fact, freely available from the training data and the generative process itself. FH significantly improves accuracy over baselines, it is directly applicable to all standard diffusion models and does not require significant additional compute. This is achieved by integrating two sources of information into a unified framework: (i) the covariance of the data distribution and (ii) the implicit covariance information available in the denoiser evaluations along the generative trajectory itself. We apply the method to linear inverse problems, where we show mathematically that accurate covariance estimates are crucial for unbiased conditional generation, and achieve significant improvements over recent methods (see Fig. 1). In summary, our contributions are:

- **Methodological:** We propose a novel, efficient method for estimating denoiser covariances in diffusion models. It (i) does not require additional training, (ii) avoids the need for expensive score Jacobian computations, (iii) adapts to the specific input and noise level, and (iv) is applicable to all standard diffusion models.
- **Analytical:** We give a theoretical analysis of why accurate covariance estimation is crucial for reconstruction guidance in linear inverse problems.
- **Practical:** Our improved covariance estimates result in significant improvements over baselines in linear inverse problems, especially with small diffusion step counts.

2 BACKGROUND

Diffusion models are a powerful framework for generative modelling. Given a data distribution $p(\mathbf{x}_0)$, we consider the following sequence of marginal distributions:

$$p(\mathbf{x}_t) = \int \mathcal{N}(\mathbf{x}_t | \mathbf{x}_0, \sigma(t)^2 \mathbf{I}) p(\mathbf{x}_0) d\mathbf{x}_0, \quad (1)$$

and corresponding reverse processes (Song et al., 2021; Karras et al., 2022)

$$\text{Reverse SDE:} \quad d\mathbf{x}_t = -2\dot{\sigma}(t)\sigma(t)\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) dt + \sqrt{2\dot{\sigma}(t)\sigma(t)} d\omega_t, \quad (2)$$

$$\text{PF-ODE:} \quad d\mathbf{x}_t = -\dot{\sigma}(t)\sigma(t)\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) dt. \quad (3)$$

Here, the $\dot{\sigma}(t) = \frac{d}{dt}\sigma(t)$ and ω_t is a Brownian motion. The *score* $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$ can be learned through score matching methods (Hyvärinen & Dayan, 2005; Vincent, 2011; Song et al., 2021). Starting at a sample $\mathbf{x}_t \sim \mathcal{N}(\mathbf{x}_t | \mathbf{x}_0, \sigma_{\max}^2 \mathbf{I})$ at a sufficiently high σ_{\max} and integrating either differential equation backwards in time, we recover the data distribution $p(\mathbf{x}_0)$ if the score is accurate.

In conditional generation, we need to define the conditional score

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}) = \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t), \quad (4)$$

which decomposes into an unconditional score and the conditional adjustment through Bayes’ rule. If we train a classifier to estimate the condition \mathbf{y} given the noisy images \mathbf{x}_t , we get *classifier guidance* (Song et al., 2021; Dhariwal & Nichol, 2021). Using additional training compute for each conditioning task, however, may be prohibitive in some applications. A more modular way to do conditional generation is to define a constraint $p(\mathbf{y} | \mathbf{x}_0)$ only on the clean data points \mathbf{x}_0 , and estimate

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t) = \nabla_{\mathbf{x}_t} \log \int \underbrace{p(\mathbf{y} | \mathbf{x}_0)}_{\text{constraint}} \underbrace{p(\mathbf{x}_0 | \mathbf{x}_t)}_{\text{denoise}} d\mathbf{x}_0 = \nabla_{\mathbf{x}_t} \log \mathbb{E}_{p(\mathbf{x}_0 | \mathbf{x}_t)} [p(\mathbf{y} | \mathbf{x}_0)]. \quad (5)$$

That is, we need to integrate over all possible denoisings of \mathbf{x}_t and their constraints. We want to avoid costly repeated sampling $\mathbf{x}_0 | \mathbf{x}_t$ from the reverse and seek practical approximations to $p(\mathbf{x}_0 | \mathbf{x}_t)$. A common approach is a Gaussian approximation

$$p(\mathbf{x}_0 | \mathbf{x}_t) \approx \mathcal{N}(\mathbf{x}_0 | \boldsymbol{\mu}, \boldsymbol{\Sigma}), \quad (6)$$

which is appealing since the so-called *Tweedie’s formula* (Efron, 2011) links the score function $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$ to exact moments of the posterior $p(\mathbf{x}_0 | \mathbf{x}_t) = \frac{p(\mathbf{x}_0)p(\mathbf{x}_t | \mathbf{x}_0)}{p(\mathbf{x}_t)}$,

$$\mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t] = \mathbf{x}_t + \sigma_t^2 \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t), \quad (7)$$

$$\text{Cov}[\mathbf{x}_0 | \mathbf{x}_t] = \sigma_t^2 \left(\sigma_t^2 \underbrace{\nabla_{\mathbf{x}_t}^2 \log p(\mathbf{x}_t)}_{\text{Hessian}} + \mathbf{I} \right). \quad (8)$$

The correct mean (Eq. (7)) of the denoiser is directly implied by the score function, as long as our estimate of the score is accurate. Estimating the Tweedie covariance through the full Hessian in Eq. (8) is very expensive for high-dimensional data, however, and multiple methods have been proposed.

2.1 RELATED WORK

Denoiser covariance estimation in diffusion models Previous attempts to improve denoiser covariance estimation in diffusion models can be broadly categorized into four categories:

1. **Heuristic methods:** Many methods (Ho et al., 2022; Song et al., 2023a;b) use a heuristic scaled identity covariance or can be seen as a special case where the covariance is zero (Chung et al., 2023). The methods are simple to implement but may result in biases in the conditional distributions.
2. **Training-based methods:** These involve training neural networks to directly output covariance estimates (Nichol & Dhariwal, 2021; Meng et al., 2021; Bao et al., 2022a; Peng et al., 2024). While potentially powerful, these approaches are not directly applicable to many existing diffusion models.
3. **Gradient-based methods:** These techniques estimate covariances by computing gradients of the denoiser (Finzi et al., 2023; Boys et al., 2023; Rozet et al., 2024; Manor & Michaeli, 2024b). However, getting the full covariance is computationally expensive and memory-intensive, making it challenging to apply to high-dimensional data without additional approximations. Manor & Michaeli (2024b) propose a method to compute the covariance along the top principal components efficiently.
4. **Post-hoc constant variance methods:** These approaches optimize constant variances for each time step based on pre-trained diffusion model scores (Bao et al., 2022b; Peng et al., 2024). While they do not require training or significant extra compute, they are limited in their ability to adapt to different inputs.

Diffusion models for inverse problems and training-free conditional generation Recent reviews can be found in Daras et al. (2024); Luo et al. (2024). Many works explicitly train conditional diffusion models for different tasks (Li et al., 2022; Saharia et al., 2022; Whang et al., 2022).

Many other methods adapt pre-trained diffusion models for inverse problems at inference time. DPS (Chung et al., 2022), IIGDM (Song et al., 2023a), TMPD (Boys et al., 2023) and Peng et al. (2024); Song et al. (2023b); Ho et al. (2022) use backpropagation to explicitly approximate Eq. (5). We focus on these methods because they share a common framework with varying covariance approximations, enabling straightforward comparisons. Other methods modify the generative process to reduce the

residual $\mathbf{y} - \mathbf{A}\mathbf{x}_t$ in linear inverse problems (Song et al., 2021; Jalal et al., 2021; Choi et al., 2021). DDS (Chung et al., 2024) and DiffPIR (Zhu et al., 2023) determine guidance direction by optimizing for an \mathbf{x}_0 that balances proximity to both the measurement and denoiser output. DDNM (Wang et al., 2023) projects the denoised \mathbf{x}_0 to the null-space of the measurement operator during sampling. Peng et al. (2024) show that DDNM and DiffPIR can be framed in a similar framework. Rout et al. (2024) use second-order corrections to reconstruction guidance to mitigate biases in first-order Tweedie. Kawar et al. (2021; 2022) decompose the linear measurement operator with SVD to create specialized conditional samplers. Methods based on variational inference optimize for \mathbf{x}_0 that match with the observations while having high diffusion model likelihood (Mardani et al., 2024; Feng et al., 2023). Ben-Hamu et al. (2024); Wang et al. (2024b) optimize the noise latent \mathbf{x}_T such that it matches with the observation. The methods by Wu et al. (2024); Dou & Song (2024); Trippe et al. (2023) frame conditional generation and inverse problems with a Bayesian filtering perspective, giving asymptotic guarantees with increasing compute.

Other applications of Hessians and denoiser covariances Linhart et al. (2024) demonstrate that Gaussian approximation of $p(\mathbf{x}_0 | \mathbf{x}_t)$ enables compositional generation between diffusion models. Higher-order ODE solvers leverage the Hessian $\nabla_{\mathbf{x}_t}^2 \log p(\mathbf{x}_t)$ for efficient sampling (Dockhorn et al., 2022), while Sanchez et al. (2022) use it for causal discovery. Lu et al. (2022) improved model likelihoods by explicitly matching higher-order score gradients. Song & Lai (2024) identified the Hessian as equivalent to Fisher information for measuring step informativeness in conditional generation. Recent advances include efficient Hessian computation using training data (Wang et al., 2024a) and unsupervised audio editing via perturbing generation along denoiser covariance principal components (Manor & Michaeli, 2024a).

3 METHODS

We present our framework for incorporating prior data covariance information with curvature information observed during sampling. We define $\boldsymbol{\mu}_{0|t}(\mathbf{x}_t)$ and $\boldsymbol{\Sigma}_{0|t}(\mathbf{x}_t)$ as our approximations of $\mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t]$ and $\text{Cov}[\mathbf{x}_0 | \mathbf{x}_t]$ at time t and location \mathbf{x} . As we move from point (\mathbf{x}, t) to $(\mathbf{x} + \Delta\mathbf{x}, t + \Delta t)$ in the diffusion process, the denoiser covariance changes but remains similar for small steps. We develop methods to transfer this information across time steps (Sec. 3.1), incorporate additional curvature information (Sec. 3.2), and combine these updates (Sec. 3.3). For high-dimensional data, we propose an efficient algorithm using diagonal and low-rank structures (Sec. 3.4). We discuss covariance initialization (Sec. 3.5) and introduce reconstruction guidance with a linear-Gaussian observation model (Sec. 3.6). Finally, we analyze why diagonal denoiser covariance overestimates guidance for correlated data at large diffusion times and demonstrate this issue with image data, showing that the problem is resolved with correct covariance estimation (Sec. 3.7).

Notation In the following, we interchangeably use $p(\mathbf{x}, t)$ in place of $p(\mathbf{x}_t)$ where we want to emphasise the possibility to change either \mathbf{x} or t , but not the other. However, in contexts where we talk about the posterior, we use $p(\mathbf{x}_t)$ and $p(\mathbf{x}_0 | \mathbf{x}_t)$ to emphasise the difference between the two random variables \mathbf{x}_0 and \mathbf{x}_t .

3.1 TIME UPDATE

Our goal is to obtain the evolution of the denoiser moments $\boldsymbol{\mu}_{0|t+\Delta t}(\mathbf{x}_t)$ and $\boldsymbol{\Sigma}_{0|t+\Delta t}(\mathbf{x}_t)$ (Eqs. (7) and (8)). The evolution of the moments under the diffusion process is characterised by the Fokker-Planck equation, and in practise intractable. We approximate the evolution with a second-order Taylor expansion of $\log p(\mathbf{x}_t)$ around point \mathbf{x}_t , which leads to a Gaussian form for $p(\mathbf{x}_t)$:

$$p(\mathbf{x}'_t) \approx \mathcal{N}(\mathbf{x}'_t | \mathbf{m}(\mathbf{x}_t, t), \mathbf{C}(\mathbf{x}_t, t)), \quad (9)$$

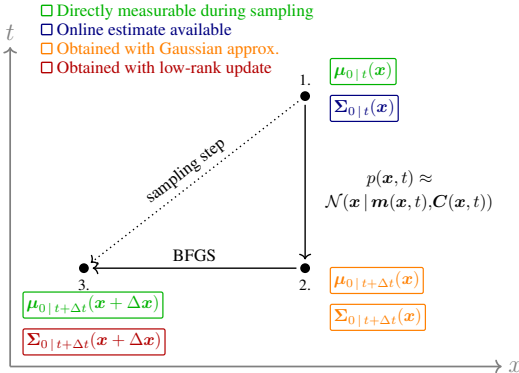


Figure 3: Sketch of our method during sampling.

where (temporarily dropping out the subscript from \mathbf{x}_t for clarity)

$$\mathbf{m}(\mathbf{x}, t) = \mathbf{x} - \nabla_{\mathbf{x}}^2 \log p(\mathbf{x}, t)^{-1} \nabla_{\mathbf{x}} \log p(\mathbf{x}, t), \quad (10)$$

$$\mathbf{C}(\mathbf{x}, t) = -[\nabla_{\mathbf{x}}^2 \log p(\mathbf{x}, t)]^{-1}. \quad (11)$$

The evolution of the Gaussian (Eq. (9)) under the linear forward SDE (Eq. (1)) has a closed form (Särkkä & Solin, 2019). With the forward process induced by Eq. (1), this results in

$$\mathbf{m}(\mathbf{x}, t + \Delta t) = \mathbf{m}(\mathbf{x}, t), \quad (12)$$

$$\mathbf{C}(\mathbf{x}, t + \Delta t) = \mathbf{C}(\mathbf{x}, t) + \Delta \sigma^2 \mathbf{I}, \quad (13)$$

where $\Delta \sigma^2 = \sigma^2(t + \Delta t) - \sigma^2(t)$. That is, the forward keeps the mean intact while increasing the covariance. Using the above equations we can derive Tweedie moment updates:

$$\boldsymbol{\mu}_{0|t+\Delta t}(\mathbf{x}_t) = \mathbf{x}_t + \sigma(t + \Delta t)^2 \left(\sigma(t + \Delta t)^2 \mathbf{I} - \frac{\Delta \sigma^2}{\sigma(t)^2} \boldsymbol{\Sigma}_{0|t}(\mathbf{x}_t) \right)^{-1} (\boldsymbol{\mu}_{0|t}(\mathbf{x}_t) - \mathbf{x}_t), \quad (14)$$

$$\boldsymbol{\Sigma}_{0|t+\Delta t}(\mathbf{x}_t) = \left(\boldsymbol{\Sigma}_{0|t}(\mathbf{x}_t)^{-1} + \Delta \sigma^{-2} \mathbf{I} \right)^{-1}. \quad (15)$$

where $\Delta \sigma^{-2} = \sigma(t + \Delta t)^{-2} - \sigma(t)^{-2}$. The complete derivations can be found in App. A. As Δt approaches zero, the Gaussian approximation becomes increasingly accurate. This is because the solution to the Fokker–Planck equation (which simplifies to the heat equation for variance-exploding diffusion) is a convolution with a small Gaussian $\mathcal{N}(\mathbf{x} | \mathbf{0}, \sigma(t)^2 \mathbf{I})$. The integral $\int p(\mathbf{x}_t) \mathcal{N}(\mathbf{x}_t - \mathbf{x}_s | \mathbf{0}, \sigma(t)^2 \mathbf{I}) d\mathbf{x}_s$ is dominated by values near \mathbf{x}_t , as the Gaussian rapidly diminishes further away.

3.2 SPACE UPDATE FOR ADDING NEW LOW-RANK INFORMATION DURING SAMPLING

We take inspiration from quasi-Newton methods (*e.g.*, BFGS, see Luenberger et al., 1984) in optimization, where repeated gradient evaluations at different points are used for low-rank updates to the Hessian of the function to optimize. The diffusion sampling process is similar: we gather gradient evaluations $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$ at different locations, and could use them to update the Hessian $\nabla_{\mathbf{x}_t}^2 \log p(\mathbf{x}_t)$. The Hessian is then connected to the denoiser covariance via Eq. (8). Here, we derive an even more convenient method to update $\boldsymbol{\Sigma}_{0|t}(\mathbf{x})$ directly.

To use update rules like BFGS, $\boldsymbol{\Sigma}_{0|t}(\mathbf{x})$ should be the Jacobian of some function. Thankfully, we notice that $\text{Cov}[\mathbf{x}_0 | \mathbf{x}_t]$ is proportional to the Jacobian of expectation $\mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t]$:

$$\mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t] \sigma(t)^2 = (\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) \sigma(t)^2 + \mathbf{x}) \sigma(t)^2, \quad (\text{Eq. (7)}) \quad (16)$$

$$\nabla_{\mathbf{x}_t} (\mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t] \sigma(t)^2) = (\nabla_{\mathbf{x}_t}^2 \log p(\mathbf{x}_t) \sigma(t)^2 + \mathbf{I}) \sigma(t)^2 = \text{Cov}[\mathbf{x}_0 | \mathbf{x}_t]. \quad (\text{Eq. (8)}) \quad (17)$$

We can then directly formulate the finite difference update condition equation for our estimate $\boldsymbol{\Sigma}_{0|t}(\mathbf{x})$:

$$\sigma(t)^2 (\boldsymbol{\mu}_{0|t}(\mathbf{x} + \Delta \mathbf{x}) - \boldsymbol{\mu}_{0|t}(\mathbf{x})) \approx [\boldsymbol{\Sigma}_{0|t}(\mathbf{x} + \Delta \mathbf{x})] \Delta \mathbf{x}. \quad (18)$$

This allows us to use a BFGS-like update procedure for the covariance and inverse covariance

$$\boldsymbol{\Sigma}_{0|t}(\mathbf{x} + \Delta \mathbf{x}) = \boldsymbol{\Sigma}_{0|t}(\mathbf{x}) - \frac{\boldsymbol{\Sigma}_{0|t}(\mathbf{x}) \Delta \mathbf{x} \Delta \mathbf{x}^\top \boldsymbol{\Sigma}_{0|t}(\mathbf{x})}{\Delta \mathbf{x}^\top \boldsymbol{\Sigma}_{0|t}(\mathbf{x}) \Delta \mathbf{x}} + \frac{\Delta \mathbf{e} \Delta \mathbf{e}^\top}{\Delta \mathbf{e}^\top \Delta \mathbf{x}}, \quad (19)$$

$$\boldsymbol{\Sigma}_{0|t}(\mathbf{x} + \Delta \mathbf{x})^{-1} = (\mathbf{I} - \gamma \Delta \mathbf{x} \Delta \mathbf{e}^\top) \boldsymbol{\Sigma}_{0|t}(\mathbf{x})^{-1} (\mathbf{I} - \gamma \Delta \mathbf{e} \Delta \mathbf{x}^\top) + \gamma \Delta \mathbf{x} \Delta \mathbf{x}^\top, \quad (20)$$

where

$$\Delta \mathbf{e} = \sigma(t)^2 (\boldsymbol{\mu}_{0|t}(\mathbf{x} + \Delta \mathbf{x}) - \boldsymbol{\mu}_{0|t}(\mathbf{x})) \quad \text{and} \quad \gamma = \frac{1}{\Delta \mathbf{e}^\top \Delta \mathbf{x}}. \quad (21)$$

While other update rules exist, BFGS has the advantage that it preserves the positive-definiteness of the covariance matrix. We provide further discussion in App. I.

Algorithm 1: Time update**Input:** $\Sigma_{0|t}(\mathbf{x}), \sigma(t + \Delta t), \sigma(t), \boldsymbol{\mu}_{0|t}(\mathbf{x})$

- 1 $\Delta(\sigma^{-2}) = \sigma(t + \Delta t)^{-2} - \sigma(t)^{-2}$
 $\Sigma_{0|t+\Delta t}(\mathbf{x})^{-1} = \Sigma_{0|t}(\mathbf{x})^{-1} + \Delta(\sigma^{-2})\mathbf{I}$
- 2 $\Sigma_{0|t+\Delta t}(\mathbf{x}) = (\Sigma_{0|t+\Delta t}(\mathbf{x})^{-1})^{-1}$
- 3 $\boldsymbol{\mu}_{0|t+\Delta t}(\mathbf{x}) = \text{Eq. (14)}$
- 4 **return** $\Sigma_{0|t+\Delta t}(\mathbf{x}), \boldsymbol{\mu}_{0|t+\Delta t}(\mathbf{x})$

Algorithm 2: Space update**Input:** $\Sigma_{0|t}(\mathbf{x}_t), \boldsymbol{\mu}_{0|t}(\mathbf{x} + \Delta\mathbf{x}), \boldsymbol{\mu}_{0|t}(\mathbf{x}), \sigma(t), \Delta\mathbf{x}$

- 1 $\Delta\mathbf{e} = \text{Eq. (21)}$
- 2 $\gamma = \text{Eq. (21)}$
- 3 $\Sigma_{0|t}(\mathbf{x} + \Delta\mathbf{x}) = \text{Eq. (19)}$
- 4 **return** $\Sigma_{0|t}(\mathbf{x} + \Delta\mathbf{x})$

3.3 COMBINING THE UPDATES FOR PRACTICAL SAMPLERS

Note that the update step requires $\boldsymbol{\mu}_{0|t}(\mathbf{x})$ evaluated at two \mathbf{x} locations but with the same t . Usually, we only have $\boldsymbol{\mu}_{0|t}(\mathbf{x})$ at *different* t during sampling, however. The solution is that we can combine the time updates with the space updates in any diffusion model sampler as follows: Let’s say we have two consecutive score evaluations $\nabla_{\mathbf{x}} \log p(\mathbf{x}, t)$ and $\nabla_{\mathbf{x}} \log p(\mathbf{x} + \Delta\mathbf{x}, t + \Delta t)$. We first update the denoiser mean and covariance with the time update to get estimates of $\Sigma_{0|t+\Delta t}(\mathbf{x})$ and $\boldsymbol{\mu}_{0|t+\Delta t}(\mathbf{x})$. Then we can update $\Sigma_{0|t+\Delta t}(\mathbf{x})$ with the BFGS update since we have $\boldsymbol{\mu}_{0|t+\Delta t}(\mathbf{x})$ from the time update and $\boldsymbol{\mu}_{0|t+\Delta t}(\mathbf{x} + \Delta\mathbf{x})$ from the second score function evaluation and Eq. (7). This is visualized in Fig. 3, and the algorithms for updating the covariance are given in Alg. 1 and Alg. 2.

3.4 PRACTICAL IMPLEMENTATION FOR HIGH-DIMENSIONAL DATA

While the method described so far works well for low-dimensional data, storing entire covariance matrices in memory is difficult for high-dimensional data. Luckily, this is not necessary since we only perform low-rank updates to the covariance matrix. In practice, we keep track of the following representation of the denoiser covariance:

$$\Sigma_{0|t}(\mathbf{x}) = \mathbf{D} + \mathbf{U}\mathbf{U}^\top - \mathbf{V}\mathbf{V}^\top, \quad (22)$$

where \mathbf{D} is diagonal and \mathbf{U}, \mathbf{V} are low-rank $N \times k$ matrices. This structure comes from the two outer products in the the BFGS update (positive and negative). The vectors $\frac{\Delta\mathbf{e}}{\sqrt{\Delta\mathbf{e}^\top \Delta\mathbf{x}}}$ and $\frac{\Sigma_{0|t}(\mathbf{x})\Delta\mathbf{x}}{\sqrt{\Delta\mathbf{x}^\top \Sigma_{0|t}(\mathbf{x})\Delta\mathbf{x}}}$ become new columns in \mathbf{U} and \mathbf{V} respectively. In App. B, we show that inverting this matrix structure yields another matrix of the same form: $\Sigma_{0|t}(\mathbf{x})^{-1} = \mathbf{D}' + \mathbf{U}'\mathbf{U}'^\top - \mathbf{V}'\mathbf{V}'^\top$. Using two applications of the Woodbury identity, this computation only requires inverting $k \times k$ matrices rather than $N \times N$ ones, enabling efficient calculation of both $\Sigma_{0|t}(\mathbf{x} + \Delta\mathbf{x})^{-1}$ and the time update inverse.

3.5 INITIALISATION OF THE COVARIANCE

Having established methods for representing and updating denoiser covariances, we address initialization. While one might consider the limit $t \rightarrow \infty$ where $p(\mathbf{x}_t) \rightarrow \mathcal{N}(\mathbf{x}_t | \mathbf{0}, \sigma(t)^2\mathbf{I})$ and $\nabla_{\mathbf{x}_t}^2 \log p(\mathbf{x}_t) \rightarrow -\frac{\mathbf{I}}{\sigma(t)^2}$, this is suboptimal: although the Hessian approaches identity at high t , the denoiser covariance approaches the data covariance. We estimate this from the data and initialise the covariance to it. For high-dimensional data, we approximate this covariance as diagonal in the DCT basis: $\Sigma_t(\mathbf{x}_t) = \Gamma_{\text{DCT}}\mathbf{D}\Gamma_{\text{DCT}}^\top$. This is justified by natural images being approximately diagonal in frequency bases (Hyvärinen et al., 2009). While alternatives like PCA could be used, we found the DCT-based method sufficient. We provide additional discussion on the DCT basis in App. I.

3.6 GUIDANCE WITH A LINEAR-GAUSSIAN OBSERVATION MODEL

If the observation model $p(\mathbf{y} | \mathbf{x}_0)$ is linear-Gaussian, the reconstruction guidance becomes

$$\begin{aligned} \nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t) &\approx \nabla_{\mathbf{x}_t} \log \int \mathcal{N}(\mathbf{y} | \mathbf{A}\mathbf{x}_0, \sigma_y^2\mathbf{I})\mathcal{N}(\mathbf{x}_0 | \boldsymbol{\mu}_{0|t}(\mathbf{x}_t), \Sigma_{0|t}(\mathbf{x}_t)) d\mathbf{x}_0 \\ &= (\mathbf{y} - \mathbf{A}\boldsymbol{\mu}_{0|t}(\mathbf{x}_t))^\top (\mathbf{A}\Sigma_{0|t}(\mathbf{x}_t)\mathbf{A}^\top + \sigma_y^2\mathbf{I})^{-1} \mathbf{A}\nabla_{\mathbf{x}_t}\boldsymbol{\mu}_{0|t}(\mathbf{x}_t), \end{aligned} \quad (23)$$

where \mathbf{A} is the linear measurement operator (e.g., blurring). $\boldsymbol{\mu}_{0|t}(\mathbf{x}_t)$ is obtained using Tweedie’s formula, and $\Sigma_{0|t}(\mathbf{x}_t) = \Sigma_{0|t}$ is assumed constant with respect to \mathbf{x}_t when taking the derivative.

The linear-Gaussian setting is valuable as it both represents many real-world problems (deblurring, inpainting) and provides analytic insights into $\Sigma_{0|t}$ choices. We will show that simplistic denoiser covariance approximations lead to severe overestimation of the guidance scale in Eq. (23).

3.7 ISSUES WITH DIAGONAL DENOISER COVARIANCE

In this section, we will focus on DPS (Chung et al., 2023) and IIGDM (Song et al., 2023a) for the linear inverse problem case. Both can be cast as using the same formula with $\Sigma_{0|t} = r_t^2 \mathbf{I}$ and different post-processing steps on the resulting $\nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t)$ approximation:

1. In DPS, $r_t^2 = 0$. The resulting $\nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t)$ is scaled with $\frac{\xi \sigma_y^2}{\|\mathbf{y} - \mathbf{A}\mathbf{x}_0\|}$ (ξ is a hyperparameter).
2. In IIGDM, $r_t^2 = \frac{\sigma(t)^2}{1 + \sigma(t)^2}$. The resulting $\nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t)$ is further scaled with r_t^2 .

Importantly, in the case where the resulting $\mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t, \mathbf{y}] = \mathbf{x}_t + \sigma(t)^2 \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y})$ approximation is outside the data range $[-1, 1]$, the modified score in both DPS and IIGDM is clipped to keep the denoiser mean within this range. Other choices include $r_t^2 = \sigma(t)^2$ (Ho et al., 2022).

The mismatch between simplistic covariance and the denoiser Jacobian Notice that according to Eq. (8), $\nabla_{\mathbf{x}_t} \mu_{0|t}(\mathbf{x}_t) \approx \frac{\text{Cov}[\mathbf{x}_0 | \mathbf{x}_t]}{\sigma(t)^2}$, where $\text{Cov}[\mathbf{x}_0 | \mathbf{x}_t]$ is the real denoiser covariance. For real data like images, this denoiser covariance is highly non-diagonal due to pixel correlations. This creates tension with the inverse $(\mathbf{A}\Sigma_{0|t}\mathbf{A}^\top + \sigma_y^2 \mathbf{I})^{-1}$ in Eq. (23), which assumes diagonal $\Sigma_{0|t}$.

A toy model For the denoising task $\mathbf{A} = \mathbf{I}$, consider images with perfectly correlated pixels (same color), giving $\text{Cov}[\mathbf{x}_0] = \mathbf{J}$ where \mathbf{J} is all ones. As $t \rightarrow \infty$, $\text{Cov}[\mathbf{x}_0 | \mathbf{x}_t] \rightarrow \text{Cov}[\mathbf{x}_0]$. Assume that the observation \mathbf{y} and the denoiser mean $\mu_{0|t}(\mathbf{x}_t)$ are similarly vectors of ones $\vec{\mathbf{1}}$ scaled by a constant, and thus $\mathbf{y} - \mu_{0|t}(\mathbf{x}_t) = a\vec{\mathbf{1}}$. Note that $\nabla_{\mathbf{x}_t} \mu_{0|t}(\mathbf{x}_t) = \frac{\mathbf{J}}{\sigma(t)^2}$. Now, the guidance terms with $r_t^2 = \frac{\sigma(t)^2}{1 + \sigma(t)^2}$ read:

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t) = a\vec{\mathbf{1}}^\top \frac{1}{1 + \sigma_y^2} \frac{\mathbf{J}}{\sigma(t)^2} = \frac{aN}{(1 + \sigma_y^2)\sigma(t)^2} \vec{\mathbf{1}}^\top, \quad (24)$$

Here N is the data dimensionality. Two key issues emerge: (1) the per-pixel guidance term scales with the total pixel count, and (2) for typical values ($a \approx 1$, $\sigma_y^2 \ll 1$), the guidance becomes implausibly large. For a 1000×1000 image, $\sigma(t)^2 \nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t) \approx N\vec{\mathbf{1}}$, yielding values around 10^6 per pixel—far beyond the $[-1, 1]$ data range. This issue is even worse in DPS where $r_t^2 = 0$. For IIGDM, clipping the denoiser mean to $[-1, 1]$ prevents trajectory blow-up but loses information and introduces biases. The scaling factors in DPS also reduce but do not eliminate the problem.

Solution with the correct covariance In contrast, the same issue does not occur if we use the correct denoiser covariance in the formula:

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t) \approx (\mathbf{y} - \mu_{0|t}(\mathbf{x}_t))^\top (\text{Cov}[\mathbf{x}_0 | \mathbf{x}_t] + \sigma_y^2 \mathbf{I})^{-1} \frac{\text{Cov}[\mathbf{x}_0 | \mathbf{x}_t]}{\sigma_t^2}. \quad (25)$$

Clearly, if $\sigma_y \rightarrow 0$, the covariances cancel out. Thus, the scale of the calculated guidance does not cause issues. In App. C, we repeat this analysis without assuming $\sigma_y = 0$.

Fig. 4 showcases the issue in practice for a Gaussian blur operator \mathbf{A} in Imagenet 256×256 and a denoiser from Dhariwal & Nichol (2021). In comparison, the problem is less severe for a more sophisticated DCT-diagonal covariance approximation. However, even the DCT-diagonal method does cause the adjusted denoiser mean to diverge at high noise levels.

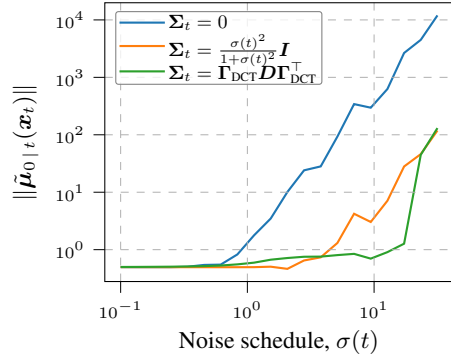


Figure 4: Norm of $\mu_{0|t}(\mathbf{x}) + \sigma(t)^2 \nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t)$ for different covariance estimation methods on ImageNet 256×256 . Values > 1 indicate overestimation since the data is normalized to $[-1, 1]$.

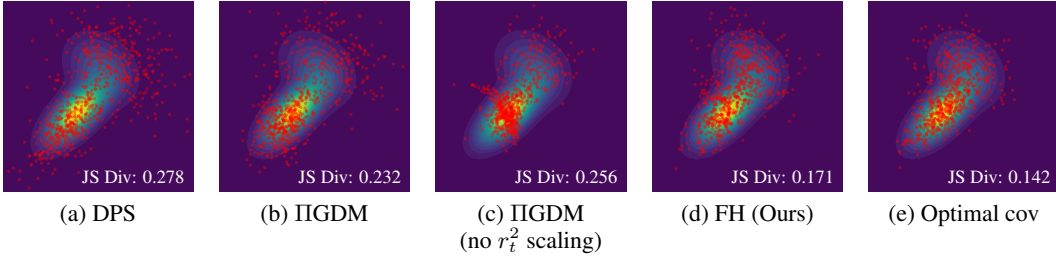


Figure 5: Different methods for posterior inference in the example in Fig. 2 and Jensen–Shannon divergences to the true posterior.

Approximating $\nabla_{\mathbf{x}_t} \boldsymbol{\mu}_{0|t}(\mathbf{x}_t)$ Given that the problem stems from the mismatch between $\nabla_{\mathbf{x}_t} \boldsymbol{\mu}_{0|t}(\mathbf{x}_t)$ and our covariance approximation, a solution needs to harmonize these two terms. Thus, we propose to further approximate $\nabla_{\mathbf{x}_t} \boldsymbol{\mu}_{0|t}(\mathbf{x}_t)$ by our denoiser covariance estimate when the scale of the guidance by calculating the full Jacobian is too large. In practice, we first calculate the adjusted denoiser covariance estimate $\boldsymbol{\mu}_{0|t}(\mathbf{x}_t)$, and fall back to approximating $\nabla_{\mathbf{x}_t} \boldsymbol{\mu}_{0|t}(\mathbf{x}_t) \approx \frac{\boldsymbol{\Sigma}_{0|t}(\mathbf{x}_t)}{\sigma(t)^2}$ in case our initial approximation $\|\sigma(t)^2 \nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t)\| > 1$ (which would push the trajectory in a direction that is outside the data range $[-1, 1]$).

Full algorithm. In App. D, we show full algorithms for linear inverse problems with our covariance estimation method, one with the Euler ODE solver and another that works with any solver.

4 EXPERIMENTS

We validate our method using synthetic Gaussian mixture model data and compare it against baselines on linear imaging inverse problems. Our experiments demonstrate that our more sophisticated covariance approximations reduce bias and improve results, particularly at lower diffusion step counts. We use a linear schedule $\sigma(t) = t$, as advocated by Karras et al. (2022), and follow their settings for our image diffusion models otherwise as well. For the image experiments, we used $\sigma_{\max} = 80$ and $\sigma_{\max} = 20$ for the synthetic data. We use a simple Euler sampler for the synthetic data experiments and a 2nd order Heun method (Karras et al., 2022) for the image experiments.

4.1 SYNTHETIC DATA EXPERIMENTS

Toy data We first showcase the performance of different methods on a toy problem using a mixture of Gaussians distribution, which admits a closed-form formula for the score (see App. M). The results in Fig. 5 show that our method clearly outperforms DPS and IIGDM, approaching the method using optimal covariance obtained by backpropagation and Eq. (8). Note that the example favours DPS, since we tuned the guidance hyperparameter for this particular task.

The effect of dimensionality and correlation In Sec. 3.7, we noticed that the guidance scale is overestimated the larger the dimensionality is. A practical consequence is that the variance of the generative distribution can be underestimated. In App. E, we directly showcase this with synthetic data and show that it does not happen with our method.

Approximation error in the covariance In App. G, we analyse the error in the covariance approximation for a low-dimensional example, and empirically show that the error approaches zero with a large amount of steps and a stochastic sampler.

4.2 IMAGE DATA AND LINEAR INVERSE PROBLEMS

We experiment on ImageNet 256×256 (Deng et al., 2009) with an unconditional denoiser from Dhariwal & Nichol (2021). We evaluate the models on four linear inverse problems: Gaussian deblurring, motion deblurring, random inpainting, and super-resolution. We evaluate our models with peak signal-to-noise ratio (PSNR), structural similarity index measure (SSIM, Wang et al., 2004) and learned perceptual image patch similarity (LPIPS, Zhang et al., 2018) on the ImageNet test set. We use the same set of 1000 randomly selected images for all models.

We solve the inverse in Eq. (23) using conjugate gradient, following Peng et al. (2024). Our custom PyTorch implementation uses GPU acceleration and adjusts solver tolerance based on noise levels. This optimization removes the solver bottleneck without noticeable performance loss. Details are in App. J.

Proposed models We introduce four new methods. The first, ‘Identity’, is initialized with identity covariance. ‘Identity+Online’ also uses the space updates. ‘FH’ is initialized with data covariance projected to a DCT-diagonal basis. Finally, ‘FH+Online’ enhances ‘FH’ with online updates.

Baselines We compare against several methods using Eq. (23) for linear imaging inverse problems: DPS (Chung et al., 2023), IIGDM (Song et al., 2023a), TMPD (Boys et al., 2023), and two methods from Peng et al. (2024) - Peng (Convert) and Peng (Analytic). TMPD uses vector-Jacobian product $\tilde{\mathbf{I}}^\top \nabla_{\mathbf{x}_t} \mathbf{x}_0(s_t) \sigma(t)^2$ for denoiser covariance. Convert employs neural network-output pixel-space diagonal covariance, while Analytic determines optimal constant pixel-diagonal covariances per timestep through moment matching. These methods were selected as they represent reconstruction guidance with different covariances, enabling analysis of our covariance approximation approach. For DPS, we optimized guidance scale via ImageNet validation set sweeps. Non-identity covariance models used SciPy’s conjugate gradient method for solving Eq. (23). For TMPD, we adjusted tolerance at higher noise levels to reduce generation time (see App. J).

Scaling the guidance term We investigated how covariance approximation affects the need for post-hoc changes to the estimated gradient $\nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t)$, as shown in Fig. 6. For deblurring, the cruder identity initialization required scaling slightly below 1, indicating an initial overestimation of the guidance scale. The more sophisticated DCT-diagonal covariance (FH) showed no systematic over- or underestimation, with optimal scaling at 1. We determined optimal guidance strength for identity covariance through a small sweep of 100 ImageNet validation samples at different solver step counts. No scaling was applied for DCT-diagonal covariance. Additional analysis with PSNR and SSIM is provided in App. K.

Baseline comparisons Our experiments focus on the low ODE sampling step regime to ensure practical applicability. Results in Table 1 show that adding online updates during sampling improves performance, with even greater gains when using DCT-diagonal covariance instead of the identity base covariance. On low step counts, our FH models consistently outperform baselines across all metrics, particularly on LPIPS scores. Visual comparisons in Fig. 7 and Fig. 9 confirm the effective fine detail preservation of FH at 15 and 30 steps. Extended results with 50 and 100 steps and the Euler solver in App. H, including comparisons to non-reconstruction guidance methods DDNM+ (Wang et al., 2023) and DiffPIR (Zhu et al., 2023), show FH and FH+Online almost always outperforming others at low step counts and typically achieving the best LPIPS scores even at higher step counts.

5 CONCLUSIONS

We introduced Free Hunch (FH), a framework for denoiser covariance estimation in diffusion models that leverages training data and trajectory curvature. FH provides accurate covariance estimates without additional training, architectural changes or ODE/SDE solver modifications. Our theoretical analysis showed that incorrect denoiser covariances significantly bias linear inverse problem solutions. Experiments on ImageNet demonstrated strong performance in linear inverse problems, especially at low step counts, with excellent LPIPS scores and fine detail preservation.

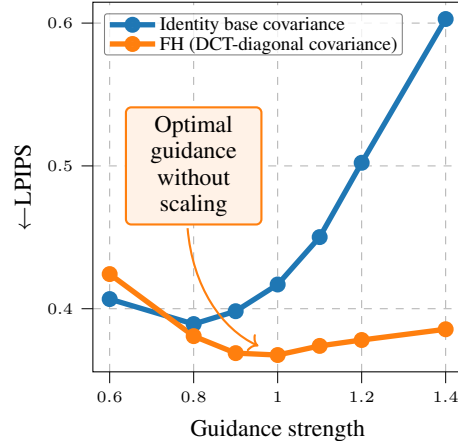


Figure 6: LPIPS w.r.t. guidance strength for the ImageNet validation set and the Gaussian blurring task. With a better covariance approximation, the usefulness of adjusting the approximated guidance $\nabla_{\mathbf{x}_t}$ with post-hoc tricks becomes smaller.

While FH introduces some additional complexity compared to simpler approaches, its efficiency and adaptability make it promising for various conditional generation tasks. Limitations of our work include the focus of the theoretical and experimental analysis on linear inverse problems, and future work could investigate nonlinear inverse problems and other types of conditional generation. Another open question is whether we can derive error bounds on the accuracy of the estimated covariance matrix in the entire process, or within individual ‘time updates’ or ‘space updates’. While our DCT-diagonal base covariance works well for image data, the application to other data domains is another open question. A low-rank estimate of the covariance matrix with a PCA decomposition seems like a generally applicable approach, but this remains to be validated in practice.

Code for our approach is available at <https://github.com/AaltoML/free-hunch>.

Table 1: Comparison of image restoration methods for 15- and 30-step Heun iterations for deblurring (Gaussian), inpainting (random), deblurring (motion), and super-resolution (4x) tasks. Our method (FH) excels overall, especially in the descriptive LPIPS metric. The top scores per category are **bolded**, with runners-up underlined. Close scores share rankings.

Method	Deblur (Gaussian)			Inpainting (Random)			Deblur (Motion)			Super res. (4x)			
	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	
15 steps	DPS	19.94	0.444	0.572	20.68	0.494	0.574	17.02	0.354	0.646	19.85	0.460	0.590
	IIGDM	20.30	0.475	0.574	19.87	0.468	0.598	19.21	0.429	0.602	20.17	0.474	0.582
	TMPD	23.08	0.597	0.420	18.99	0.481	0.539	20.80	0.491	0.514	21.88	0.545	0.476
	Peng (Convert)	22.53	0.563	0.490	22.23	0.579	0.489	20.46	0.475	0.556	21.92	0.541	0.517
	Peng (Analytic)	22.53	0.563	0.490	22.14	0.574	0.494	20.46	0.475	0.556	21.92	0.541	0.517
	Identity	22.91	0.594	0.384	18.83	0.397	0.590	20.06	0.393	0.506	22.65	0.589	0.412
	Identity+online	23.08	0.606	0.385	18.86	0.397	0.590	20.31	0.418	0.492	22.76	0.597	0.414
	FH	<u>23.39</u>	<u>0.624</u>	0.372	<u>24.73</u>	<u>0.701</u>	<u>0.327</u>	<u>21.69</u>	<u>0.534</u>	<u>0.446</u>	<u>23.30</u>	0.624	0.390
	FH+online	23.54	0.634	<u>0.378</u>	25.25	0.728	0.317	21.84	0.549	0.441	23.39	<u>0.632</u>	<u>0.394</u>
	30 steps	DPS	21.76	0.527	0.463	24.84	0.678	0.387	18.22	0.389	0.582	23.00	0.593
IIGDM		22.27	0.559	0.468	21.24	0.518	0.517	21.16	0.508	0.503	22.11	0.553	0.478
TMPD		23.16	0.602	0.415	18.85	0.481	0.537	20.91	0.500	0.507	21.94	0.549	0.472
Peng (Convert)		<u>23.61</u>	0.627	0.405	23.74	0.648	0.403	21.99	0.553	<u>0.463</u>	23.22	0.608	0.430
Peng (Analytic)		23.61	0.626	0.405	23.59	0.640	0.411	<u>21.99</u>	0.552	<u>0.463</u>	23.21	0.608	0.430
Identity		23.15	0.602	0.374	18.75	0.402	0.578	20.14	0.406	0.494	22.82	0.588	0.405
Identity+online		23.38	0.621	<u>0.359</u>	20.07	0.443	0.529	20.47	0.420	0.467	<u>23.38</u>	0.622	0.383
FH		23.55	<u>0.630</u>	0.353	<u>26.00</u>	<u>0.757</u>	0.256	21.80	0.538	0.411	<u>23.38</u>	<u>0.623</u>	0.372
FH+online		23.62	0.635	<u>0.358</u>	26.18	0.767	0.268	<u>21.88</u>	0.547	0.410	23.44	0.628	<u>0.375</u>

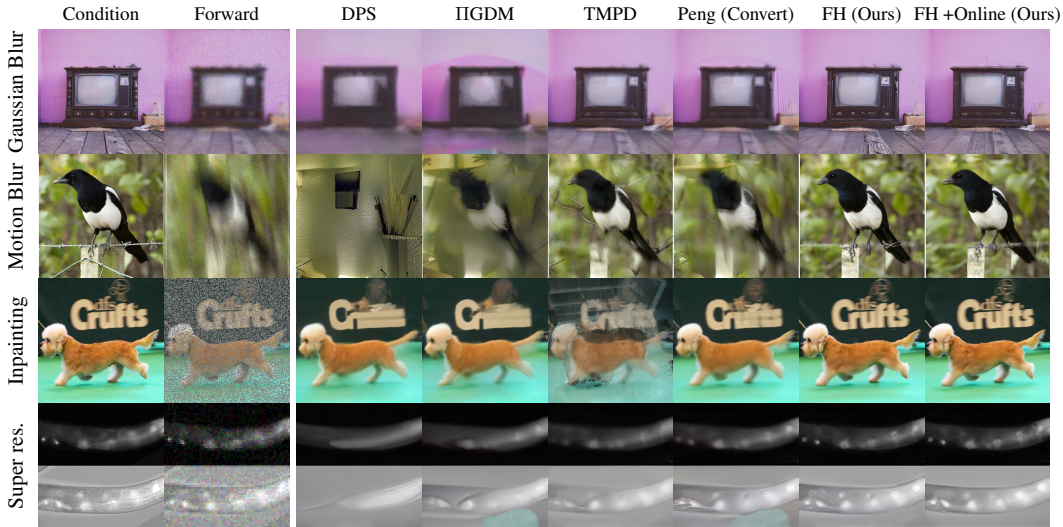


Figure 7: Qualitative examples using the 15-step Heun sampler for image restoration methods for deblurring (Gaussian), inpainting (random), deblurring (motion), and super-resolution (4x) tasks. Quantitative metrics in Table 1. Our method manages to restore the corrupted (‘Forward’) to match well with the original (‘Condition’).

ACKNOWLEDGMENTS

We acknowledge funding from the Research Council of Finland (grants 339730, 362408, 334600). We acknowledge CSC – IT Center for Science, Finland, for awarding this project access to the LUMI supercomputer, owned by the EuroHPC Joint Undertaking, hosted by CSC (Finland) and the LUMI consortium through CSC. We acknowledge the computational resources provided by the Aalto Science-IT project.

REFERENCES

- Fan Bao, Chongxuan Li, Jiacheng Sun, Jun Zhu, and Bo Zhang. Estimating the optimal covariance with imperfect mean in diffusion probabilistic models. In *International Conference on Machine Learning (ICML)*, pp. 1555–1584. PMLR, 2022a.
- Fan Bao, Chongxuan Li, Jun Zhu, and Bo Zhang. Analytic-DPM: An analytic estimate of the optimal reverse variance in diffusion probabilistic models. In *International Conference on Learning Representations (ICLR)*, 2022b.
- Heli Ben-Hamu, Omri Puny, Itai Gat, Brian Karrer, Uriel Singer, and Yaron Lipman. D-Flow: Differentiating through flows for controlled generation. In *International Conference on Machine Learning (ICML)*, pp. 3462–3483. PMLR, 2024.
- Benjamin Boys, Mark Girolami, Jakiw Pidstrigach, Sebastian Reich, Alan Mosca, and O Deniz Akyildiz. Tweedie moment projected diffusions for inverse problems. *arXiv preprint arXiv:2310.06721*, 2023.
- Jooyoung Choi, Sungwon Kim, Yonghyun Jeong, Youngjune Gwon, and Sungroh Yoon. Ilvr: Conditioning method for denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 14367–14376, 2021.
- Hyungjin Chung, Byeongsu Sim, Dohoon Ryu, and Jong Chul Ye. Improving diffusion models for inverse problems using manifold constraints. In *Advances in Neural Information Processing Systems 35 (NeurIPS)*, pp. 25683–25696. Curran Associates, Inc., 2022.
- Hyungjin Chung, Jeongsol Kim, Michael Thompson Mccann, Marc Louis Klasky, and Jong Chul Ye. Diffusion posterior sampling for general noisy inverse problems. In *International Conference on Learning Representations (ICLR)*, 2023.
- Hyungjin Chung, Suhyeon Lee, and Jong Chul Ye. Decomposed diffusion sampler for accelerating large-scale inverse problems. In *International Conference on Learning Representations (ICLR)*, 2024.
- Giannis Daras, Hyungjin Chung, Chieh-Hsin Lai, Yuki Mitsufuji, Peyman Milanfar, Alexandros G. Dimakis, Chul Ye, and Mauricio Delbracio. A survey on diffusion models for inverse problems. 2024. URL https://giannisdaras.github.io/publications/diffusion_survey.pdf.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 248–255. IEEE, 2009.
- John E Dennis, Jr and Jorge J Moré. Quasi-newton methods, motivation and theory. *SIAM Review*, 19(1):46–89, 1977.
- Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. In *Advances in Neural Information Processing Systems 34 (NeurIPS)*, pp. 8780–8794. Curran Associates, Inc., 2021.
- Tim Dockhorn, Arash Vahdat, and Karsten Kreis. Genie: Higher-order denoising diffusion solvers. In *Advances in Neural Information Processing Systems 35 (NeurIPS)*, pp. 30150–30166. Curran Associates, Inc., 2022.

- Zehao Dou and Yang Song. Diffusion posterior sampling for linear inverse problem solving: A filtering perspective. In *International Conference on Learning Representations (ICLR)*, 2024.
- Bradley Efron. Tweedie’s formula and selection bias. *Journal of the American Statistical Association*, 106(496):1602–1614, 2011.
- Berthy T Feng, Jamie Smith, Michael Rubinstein, Huiwen Chang, Katherine L Bouman, and William T Freeman. Score-based diffusion models as principled priors for inverse imaging. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 10520–10531, 2023.
- Marc Anton Finzi, Anudhyan Boral, Andrew Gordon Wilson, Fei Sha, and Leonardo Zepeda-Núñez. User-defined event sampling and uncertainty quantification in diffusion models for physical dynamical systems. In *International Conference on Machine Learning (ICML)*, pp. 10136–10152. PMLR, 2023.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems 33 (NeurIPS)*, pp. 6840–6851. Curran Associates, Inc., 2020.
- Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. In *Advances in Neural Information Processing Systems 35 (NeurIPS)*, pp. 8633–8646. Curran Associates, Inc., 2022.
- Aapo Hyvärinen and Peter Dayan. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(4), 2005.
- Aapo Hyvärinen, Jarmo Hurri, and Patrick O Hoyer. *Natural Image Statistics: A Probabilistic Approach to Early Computational Vision*, volume 39. Springer Science & Business Media, 2009.
- Ajil Jalal, Marius Arvinte, Giannis Daras, Eric Price, Alexandros G Dimakis, and Jon Tamir. Robust compressed sensing mri with deep generative priors. In *Advances in Neural Information Processing Systems 34 (NeurIPS)*, pp. 14938–14954. Curran Associates, Inc., 2021.
- Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. In *Advances in Neural Information Processing Systems 35 (NeurIPS)*, pp. 26565–26577. Curran Associates, Inc., 2022.
- Bahjat Kawar, Gregory Vaksman, and Michael Elad. Snips: Solving noisy inverse problems stochastically. In *Advances in Neural Information Processing Systems 34 (NeurIPS)*, pp. 21757–21769. Curran Associates, Inc., 2021.
- Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising diffusion restoration models. In *Advances in Neural Information Processing Systems 35 (NeurIPS)*, pp. 23593–23606. Curran Associates, Inc., 2022.
- Haoying Li, Yifan Yang, Meng Chang, Shiqi Chen, Huajun Feng, Zhihai Xu, Qi Li, and Yueting Chen. SRDiff: Single image super-resolution with diffusion probabilistic models. *Neurocomputing*, 479:47–59, 2022.
- Julia Linhart, Gabriel Victorino Cardoso, Alexandre Gramfort, Sylvain Le Corff, and Pedro LC Rodrigues. Diffusion posterior sampling for simulation-based inference in tall data settings. *arXiv preprint arXiv:2404.07593*, 2024.
- Cheng Lu, Kaiwen Zheng, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Maximum likelihood training for score-based diffusion odes by high order denoising score matching. In *International Conference on Machine Learning (ICML)*, pp. 14429–14460. PMLR, 2022.
- David G Luenberger, Yinyu Ye, et al. *Linear and Nonlinear Programming*, volume 2. Springer, 1984.
- Ziwei Luo, Fredrik K Gustafsson, Zheng Zhao, Jens Sjölund, and Thomas B Schön. Taming diffusion models for image restoration: A review. *arXiv preprint arXiv:2409.10353*, 2024.

- Hila Manor and Tomer Michaeli. Zero-shot unsupervised and text-based audio editing using DDPM inversion. In *International Conference on Machine Learning (ICML)*, pp. 34603–34629. PMLR, 2024a.
- Hila Manor and Tomer Michaeli. On the posterior distribution in denoising: Application to uncertainty quantification. In *International Conference on Learning Representations (ICLR)*, 2024b.
- Morteza Mardani, Jiaming Song, Jan Kautz, and Arash Vahdat. A variational perspective on solving inverse problems with diffusion models. In *International Conference on Learning Representations (ICLR)*, 2024.
- Chenlin Meng, Yang Song, Wenzhe Li, and Stefano Ermon. Estimating high order gradients of the data distribution by denoising. In *Advances in Neural Information Processing Systems 34 (NeurIPS)*, pp. 25359–25369. Curran Associates, Inc., 2021.
- Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning (ICML)*, pp. 8162–8171. PMLR, 2021.
- Xinyu Peng, Ziyang Zheng, Wenrui Dai, Nuoqian Xiao, Chenglin Li, Junni Zou, and Hongkai Xiong. Improving diffusion models for inverse problems using optimal posterior covariance. In *International Conference on Machine Learning (ICML)*, pp. 40347–40370. PMLR, 2024.
- Litu Rout, Yujia Chen, Abhishek Kumar, Constantine Caramanis, Sanjay Shakkottai, and Wen-Sheng Chu. Beyond first-order tweedie: Solving inverse problems using latent diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9472–9481, 2024.
- François Rozet, G r me Andry, Fran ois Lanusse, and Gilles Louppe. Learning diffusion priors from observations by expectation maximization. *arXiv preprint arXiv:2405.13712*, 2024.
- Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(4):4713–4726, 2022.
- Pedro Sanchez, Xiao Liu, Alison Q O’Neil, and Sotirios A Tsaftaris. Diffusion models for causal discovery via topological ordering. In *International Conference on Learning Representations (ICLR)*, 2022.
- Simo S rkk  and Arno Solin. *Applied Stochastic Differential Equations*, volume 10. Cambridge University Press, 2019.
- Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning (ICML)*, pp. 2256–2265. PMLR, 2015.
- Jiaming Song, Arash Vahdat, Morteza Mardani, and Jan Kautz. Pseudoinverse-guided diffusion models for inverse problems. In *International Conference on Learning Representations (ICLR)*, 2023a.
- Jiaming Song, Qinsheng Zhang, Hongxu Yin, Morteza Mardani, Ming-Yu Liu, Jan Kautz, Yongxin Chen, and Arash Vahdat. Loss-guided diffusion models for plug-and-play controllable generation. In *International Conference on Machine Learning (ICML)*, pp. 32483–32498. PMLR, 2023b.
- Kaiyu Song and Hanjiang Lai. Fisher information improved training-free conditional diffusion model. *arXiv preprint arXiv:2404.18252*, 2024.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations (ICLR)*, 2021.
- Brian L Trippe, Jason Yim, Doug Tischer, David Baker, Tamara Broderick, Regina Barzilay, and Tommi S Jaakkola. Diffusion probabilistic modeling of protein backbones in 3D for the motif-scaffolding problem. In *International Conference on Learning Representations (ICLR)*, 2023.

- Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural Computation*, 23(7):1661–1674, 2011.
- Gregory K Wallace. The JPEG still picture compression standard. *Communications of the ACM*, 34(4):30–44, 1991.
- Fangyikang Wang, Hubery Yin, Huminhao Zhu, Shaobin Zhuang, Chao Zhang, Hanbin Zhao, Hui Qian, and Chen Li. Gradient-free analytical fisher information of diffused distributions. *OpenReview preprint*, 2024a.
- Hengkang Wang, Xu Zhang, Taihui Li, Yuxiang Wan, Tiancong Chen, and Ju Sun. Dmplug: A plug-in method for solving inverse problems with diffusion models. *arXiv preprint arXiv:2405.16749*, 2024b.
- Yinhui Wang, Jiwen Yu, and Jian Zhang. Zero-shot image restoration using denoising diffusion null-space model. In *International Conference on Learning Representations (ICLR)*, 2023.
- Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: From error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004.
- Jay Whang, Mauricio Delbracio, Hossein Talebi, Chitwan Saharia, Alexandros G Dimakis, and Peyman Milanfar. Deblurring via stochastic refinement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 16293–16303, 2022.
- Luhuan Wu, Brian Trippe, Christian Naeseth, David Blei, and John P Cunningham. Practical and asymptotically exact conditional sampling in diffusion models. In *Advances in Neural Information Processing Systems 36 (NeurIPS)*. Curran Associates, Inc., 2024.
- Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 586–595, 2018.
- Yuanzhi Zhu, Kai Zhang, Jingyun Liang, Jiezhang Cao, Bihan Wen, Radu Timofte, and Luc Van Gool. Denoising diffusion models for plug-and-play image restoration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1219–1229, 2023.

APPENDICES

A FULL DERIVATIONS FOR THE TIME UPDATE

For this section, we denote $p(\mathbf{x}_t) = \log p(\mathbf{x}, t)$ to explicitly separate the time variable from the spatial variable. This notation is useful for the derivations below, but in other parts of the paper we use the \mathbf{x}_t to separate it from \mathbf{x}_0 .

Recall the mean and covariance of the Gaussian approximation at location \mathbf{x} and time t :

$$\mathbf{m}(\mathbf{x}, t) = \mathbf{x} - \nabla_{\mathbf{x}}^2 \log p(\mathbf{x}, t)^{-1} \nabla_{\mathbf{x}} \log p(\mathbf{x}, t), \quad (26)$$

$$\mathbf{C}(\mathbf{x}, t) = -\nabla_{\mathbf{x}}^2 \log p(\mathbf{x}, t)^{-1}. \quad (27)$$

The evolution of the Gaussian has a closed form (Särkkä & Solin, 2019). In the variance-exploding case this results in

$$\mathbf{m}(\mathbf{x}, t + \Delta t) = \mathbf{m}(\mathbf{x}, t), \quad (28)$$

$$\mathbf{C}(\mathbf{x}, t + \Delta t) = \mathbf{C}(\mathbf{x}, t) + \Delta \sigma^2 \mathbf{I}, \quad (29)$$

where $\Delta \sigma^2 = \sigma^2(t + \Delta t) - \sigma^2(t)$. That is, the forward keeps the mean intact while increasing the covariance. Using the above equations we can derive:

$$\nabla_{\mathbf{x}}^2 \log p(\mathbf{x}, t + \Delta t) = (\nabla_{\mathbf{x}}^2 \log p(\mathbf{x}, t)^{-1} - \Delta \sigma^2 \mathbf{I})^{-1} \quad (30)$$

$$\nabla_{\mathbf{x}} \log p(\mathbf{x}, t + \Delta t) = \nabla_{\mathbf{x}}^2 \log p(\mathbf{x}, t + \Delta t) \nabla_{\mathbf{x}}^2 \log p(\mathbf{x}, t)^{-1} \nabla_{\mathbf{x}} \log p(\mathbf{x}, t) \quad (31)$$

When connected with Equation (7) and Equation (8), we can now derive the denoiser mean and covariance updates.

$$\begin{aligned} \boldsymbol{\mu}_{0|t+\Delta t}(\mathbf{x}) &= \mathbf{x} + \sigma^2(t + \Delta t) \nabla_{\mathbf{x}} \log p(\mathbf{x}, t + \Delta t) \\ &= \mathbf{x} + \sigma^2(t + \Delta t) \underbrace{(\nabla_{\mathbf{x}}^2 \log p(\mathbf{x}, t)^{-1} - \Delta \sigma^2 \mathbf{I})^{-1}}_{\text{Hessian projection}} \underbrace{\nabla_{\mathbf{x}}^2 \log p(\mathbf{x}, t)^{-1} \nabla_{\mathbf{x}} \log p(\mathbf{x}, t)}_{\text{Hessian-score product}} \end{aligned} \quad (32)$$

$$\begin{aligned} \boldsymbol{\mu}_{0|t+\Delta t}(\mathbf{x}) &= \mathbf{x} + \sigma(t + \Delta t)^2 \left((\sigma(t)^2 + \Delta \sigma^2) \mathbf{I} - \frac{\Delta \sigma^2}{\sigma(t)^2} \boldsymbol{\Sigma}_{0|t}(\mathbf{x}) \right)^{-1} (\boldsymbol{\mu}_{0|t}(\mathbf{x}) - \mathbf{x}), \\ &= \mathbf{x} + \sigma(t + \Delta t)^2 \left(\sigma(t + \Delta t)^2 \mathbf{I} - \frac{\Delta \sigma^2}{\sigma(t)^2} \boldsymbol{\Sigma}_{0|t}(\mathbf{x}) \right)^{-1} (\boldsymbol{\mu}_{0|t}(\mathbf{x}) - \mathbf{x}), \end{aligned} \quad (33)$$

$$\boldsymbol{\Sigma}_{0|t+\Delta t}(\mathbf{x}) = \left(\boldsymbol{\Sigma}_{0|t}(\mathbf{x})^{-1} + \Delta \sigma^{-2} \mathbf{I} \right)^{-1}. \quad (34)$$

This result is not entirely obvious, and next we will provide a detailed derivation.

Deriving the denoiser covariance update With the locally Gaussian approximation on $p(\mathbf{x}, t)$, we can represent the time evolution of the $p(\mathbf{x}, t)$ covariance as the following

$$\mathbf{C}(\mathbf{x}, t) = \mathbf{C}_0 + \sigma(t)^2 \mathbf{I}, \quad (35)$$

where \mathbf{C}_0 is the hypothetical covariance when extrapolating the Gaussian time evolution to $t = 0$. Then, moving back to the \mathbf{x}_t notation, we have the following connections (Eq. (11))

$$\nabla_{\mathbf{x}_t}^2 \log p(\mathbf{x}_t) = -(\mathbf{C}_0 + \sigma(t)^2 \mathbf{I})^{-1}, \quad (36)$$

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)^{-1} = -(\mathbf{C}_0 + \sigma(t)^2 \mathbf{I}). \quad (37)$$

On the other hand, the denoiser covariance is

$$\text{Cov}[\mathbf{x}_0 | \mathbf{x}_t] = (\nabla_{\mathbf{x}_t}^2 \log p(\mathbf{x}_t) \sigma(t)^2 + \mathbf{I}) \sigma(t)^2. \quad (38)$$

The inverse of the denoiser covariance is then (Sherman–Morrison–Woodbury formula):

$$\text{Cov}[\mathbf{x}_0 | \mathbf{x}_t]^{-1} = (\nabla_{\mathbf{x}_t}^2 \log p(\mathbf{x}_t) \sigma(t)^2 + \mathbf{I})^{-1} \sigma(t)^{-2} \quad (39)$$

$$= (\mathbf{I} - (\mathbf{I} + \nabla_{\mathbf{x}_t}^2 \log p(\mathbf{x}_t)^{-1} \sigma(t)^{-2})^{-1}) \sigma(t)^{-2} \quad (\text{Woodbury}) \quad (40)$$

$$= \left(\mathbf{I} - (\mathbf{I} - (\mathbf{C}_0 + \sigma(t)^2 \mathbf{I}) \sigma(t)^{-2})^{-1} \right) \sigma(t)^{-2} \quad (37) \quad (41)$$

$$= (\mathbf{I} + \mathbf{C}_0^{-1} \sigma(t)^2) \sigma(t)^{-2} \quad (42)$$

$$= \mathbf{C}_0^{-1} + \sigma(t)^{-2} \mathbf{I}. \quad (43)$$

So the inverse of the denoiser covariance is simply a constant term plus an identity scaled with $\sigma(t)^{-2}$. This means that

$$\text{Cov}[\mathbf{x}_0 | \mathbf{x}_{t+\Delta t}]^{-1} - \text{Cov}[\mathbf{x}_0 | \mathbf{x}_t]^{-1} = (\mathbf{C}_0^{-1} + \mathbf{I}\sigma(t + \Delta t)^{-2}) - (\mathbf{C}_0^{-1} + \mathbf{I}\sigma(t)^{-2}) \quad (44)$$

$$\text{Cov}[\mathbf{x}_0 | \mathbf{x}_{t+\Delta t}]^{-1} = \text{Cov}[\mathbf{x}_0 | \mathbf{x}_t]^{-1} + \mathbf{I}\sigma(t + \Delta t)^{-2} - \mathbf{I}\sigma(t)^{-2} \quad (45)$$

$$= \text{Cov}[\mathbf{x}_0 | \mathbf{x}_t]^{-1} + \Delta\sigma(t)^{-2}\mathbf{I}. \quad (46)$$

And thus

$$\text{Cov}[\mathbf{x}_0 | \mathbf{x}_{t+\Delta t}] = (\text{Cov}[\mathbf{x}_0 | \mathbf{x}_t]^{-1} + \Delta\sigma(t)^{-2}\mathbf{I})^{-1}. \quad (47)$$

Deriving the denoiser mean update We want to calculate the expression for the updated mean $\boldsymbol{\mu}_{0|t+\Delta t}(\mathbf{x}_t)$. This mean can be written as:

$$\boldsymbol{\mu}_{0|t+\Delta t}(\mathbf{x}_t) = \mathbf{x}_t + \sigma(t + \Delta t)^2 \cdot (\nabla_{\mathbf{x}_t}^2 \log p(\mathbf{x}_t)^{-1} - \Delta\sigma^2\mathbf{I})^{-1} \nabla_{\mathbf{x}_t}^2 \log p(\mathbf{x}_t)^{-1} \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t). \quad (48)$$

We aim to simplify the term inside the parentheses. Starting from the observation:

$$(\nabla_{\mathbf{x}_t}^2 \log p(\mathbf{x}_t)^{-1} - \Delta\sigma^2\mathbf{I})^{-1} \nabla_{\mathbf{x}_t}^2 \log p(\mathbf{x}_t)^{-1} = (\mathbf{I} - \Delta\sigma^2 \nabla_{\mathbf{x}_t}^2 \log p(\mathbf{x}_t))^{-1} \quad (49)$$

and using Eq. (7) and Eq. (8), we can express $\nabla_{\mathbf{x}_t}^2 \log p(\mathbf{x}_t)$ and $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$ as functions of $\boldsymbol{\Sigma}_{0|t}(\mathbf{x}_t)$ and $\sigma(t)$, yielding:

$$\boldsymbol{\mu}_{0|t+\Delta t}(\mathbf{x}_t) = \mathbf{x}_t + \sigma(t + \Delta t)^2 \cdot \left((\sigma(t)^2 + \Delta\sigma^2) \mathbf{I} - \frac{\Delta\sigma^2}{\sigma(t)^2} \boldsymbol{\Sigma}_{0|t}(\mathbf{x}_t) \right)^{-1} (\boldsymbol{\mu}_{0|t}(\mathbf{x}_t) - \mathbf{x}_t). \quad (50)$$

This is the final simplified expression for $\boldsymbol{\mu}_{0|t+\Delta t}(\mathbf{x}_t)$.

B INVERTING THE EFFICIENT MATRIX REPRESENTATION

Let's say we have a positive-definite matrix in represented in the format $\mathbf{C} = \mathbf{D} + \mathbf{U}\mathbf{U}^\top - \mathbf{V}\mathbf{V}^\top$, where \mathbf{D} is a diagonal matrix and \mathbf{U} and \mathbf{V} are $N \times k_1$ and $N \times k_2$ matrices, respectively. N is the data dimensionality and $k \ll N$. To invert it, we use the Woodbury identity twice, first for $\mathbf{A} = \mathbf{D} + \mathbf{U}\mathbf{U}^\top$, and second for $\mathbf{A} - \mathbf{V}\mathbf{V}^\top$. The first application of the identity is:

$$\mathbf{A}^{-1} = (\mathbf{D} + \mathbf{U}\mathbf{U}^\top)^{-1} = \mathbf{D}^{-1} - \mathbf{D}^{-1}\mathbf{U} \underbrace{(\mathbf{I} + \mathbf{U}^\top \mathbf{D}^{-1}\mathbf{U})^{-1}}_{=\mathbf{K}} \mathbf{U}^\top \mathbf{D}^{-1} \quad (51)$$

$$= \mathbf{D}^{-1} - \underbrace{\mathbf{D}^{-1}\mathbf{U} \text{sqrt}(\mathbf{K})}_{=\mathbf{V}'} \text{sqrt}(\mathbf{K})^\top \mathbf{U}^\top \mathbf{D}^{-1} \quad (52)$$

$$= \mathbf{D}^{-1} - \mathbf{V}'\mathbf{V}'^\top \quad (53)$$

that is, when we invert $\mathbf{A} = \mathbf{D} + \mathbf{U}\mathbf{U}^\top$, we get something in the form $\mathbf{D}^{-1} - \mathbf{V}'\mathbf{V}'^\top$. Note that $\mathbf{I} - \mathbf{U}^\top \mathbf{D}^{-1}\mathbf{U}$ is a $k_1 \times k_1$ matrix, instead of an $N \times N$ matrix and as such is much more efficient to invert than the full $N \times N$ matrix when $k_1 \ll N$. Now, invert $\mathbf{C} = \mathbf{A} - \mathbf{V}\mathbf{V}^\top$:

$$\mathbf{C}^{-1} = (\mathbf{A} - \mathbf{V}\mathbf{V}^\top)^{-1} = \mathbf{A}^{-1} + \mathbf{A}^{-1}\mathbf{V} \underbrace{(\mathbf{I} - \mathbf{V}^\top \mathbf{A}^{-1}\mathbf{V})^{-1}}_{=\mathbf{L}} \mathbf{V}^\top \mathbf{A}^{-1} \quad (54)$$

$$= \mathbf{A}^{-1} + \underbrace{\mathbf{A}^{-1}\mathbf{V} \text{sqrt}(\mathbf{L})}_{=\mathbf{U}'} \text{sqrt}(\mathbf{L})^\top \mathbf{V}^\top \mathbf{A}^{-1} \quad (55)$$

$$= \mathbf{A}^{-1} + \mathbf{U}'\mathbf{U}'^\top = \mathbf{D}^{-1} + \mathbf{U}'\mathbf{U}'^\top - \mathbf{V}'\mathbf{V}'^\top. \quad (56)$$

Note that $\mathbf{V}^\top \mathbf{A}^{-1}\mathbf{V}$ is again efficient to compute due to the low-rank structure of \mathbf{A}^{-1} :

$$\mathbf{V}^\top \mathbf{A}^{-1}\mathbf{V} = \mathbf{V}^\top (\mathbf{D}^{-1} - \mathbf{V}'\mathbf{V}'^\top) \mathbf{V} \quad (57)$$

$$= \mathbf{V}^\top \mathbf{D}^{-1}\mathbf{V} - \mathbf{V}^\top \mathbf{V}'\mathbf{V}'^\top \mathbf{V} \quad (58)$$

$$= \underbrace{\mathbf{V}^\top \mathbf{D}^{-1}\mathbf{V}}_{k_2 \times k_2} - \underbrace{(\mathbf{V}^\top \mathbf{V}')}_{k_2 \times k_2} \underbrace{(\mathbf{V}'^\top \mathbf{V})}_{k_2 \times k_2}. \quad (59)$$

This leads to a well-behaved $k_2 \times k_2$ matrix, and the inverse $(\mathbf{I} - \mathbf{V}^\top \mathbf{A}^{-1} \mathbf{V})^{-1}$ is also efficient to compute.

The matrix square root and complex numbers Note that in addition to the $k_1 \times k_1$ and $k_2 \times k_2$ inverses, the method requires matrix square root. One might imagine that $(\mathbf{I} + \mathbf{U}^\top \mathbf{D}^{-1} \mathbf{U})^{-1}$ is guaranteed to be positive-definite due to the original matrix being such, but this is not necessarily the case. $(\mathbf{I} + \mathbf{U}^\top \mathbf{U})^{-1}$ would be, but the diagonal term \mathbf{D}^{-1} can push the eigenvalues of the matrix to negative. This means that we can not use the Cholesky decomposition for the matrix square root operations, but instead we use the Schur decomposition as implemented in the scipy library. A side-effect is also that we have to use complex numbers to represent \mathbf{D} , \mathbf{U} , and \mathbf{V} in our implementation. This is not an issue, since in the calculation of the covariance $\mathbf{D} + \mathbf{U}\mathbf{U}^\top - \mathbf{V}\mathbf{V}^\top$, the imaginary components cancel out and we get a real matrix.

C EXTENDED ANALYSIS OF THE TOY EXAMPLE

As a recap, in the toy model, $\mathbf{A} = \mathbf{I}$ and all the pixels are perfectly correlated with $\text{Cov}[\mathbf{x}_0] = \mathbf{J}$, where \mathbf{J} is a matrix full of ones. The observation \mathbf{y} and the denoiser mean $\boldsymbol{\mu}_{0|t}(\mathbf{x}_t)$ are also vectors of ones $\vec{\mathbf{1}}$ scaled by a constant, so that $\mathbf{y} - \boldsymbol{\mu}_{0|t}(\mathbf{x}_t) = a\vec{\mathbf{1}}$.

PIGDM guidance without postprocessing

$$\begin{aligned} (\mathbf{y} - \boldsymbol{\mu}_{0|t}(\mathbf{x}_t))^\top (\mathbf{I} \frac{\sigma(t)^2}{1 + \sigma(t)^2} + \sigma_y^2 \mathbf{I})^{-1} \nabla_{\mathbf{x}_t} \boldsymbol{\mu}_{0|t}(\mathbf{x}_t) \\ \approx (\mathbf{y} - \boldsymbol{\mu}_{0|t}(\mathbf{x}_t))^\top (\mathbf{I} + \sigma_y^2 \mathbf{I})^{-1} \frac{\text{Cov}[\mathbf{x}_0 | \mathbf{x}_t]}{\sigma(t)^2} \end{aligned} \quad (60)$$

$$= \frac{a}{1 + \sigma_y^2} \vec{\mathbf{1}}^\top \frac{\mathbf{J}}{\sigma(t)^2} = \frac{aN}{(1 + \sigma_y^2)\sigma(t)^2} \vec{\mathbf{1}}^\top. \quad (61)$$

DPS guidance without postprocessing

$$\begin{aligned} (\mathbf{y} - \boldsymbol{\mu}_{0|t}(\mathbf{x}_t))^\top (\mathbf{I} \mathbf{0} + \sigma_y^2 \mathbf{I})^{-1} \nabla_{\mathbf{x}_t} \boldsymbol{\mu}_{0|t}(\mathbf{x}_t) \\ \approx (\mathbf{y} - \boldsymbol{\mu}_{0|t}(\mathbf{x}_t))^\top \sigma_y^{-2} \frac{\text{Cov}[\mathbf{x}_0 | \mathbf{x}_t]}{\sigma(t)^2} \end{aligned} \quad (62)$$

$$= \frac{a}{\sigma_y^2} \vec{\mathbf{1}}^\top \frac{\mathbf{J}}{\sigma(t)^2} = \frac{aN}{\sigma_y^2 \sigma(t)^2} \vec{\mathbf{1}}^\top. \quad (63)$$

Here N is the data dimensionality.

For PIGDM, the gradient is scaled by $\frac{\sigma(t)^2}{1 + \sigma(t)^2}$, but this does not change the result in high noise levels. Instead, the clipping of the denoiser mean to $[-1, 1]$ regularises the guidance such that the generation trajectory does not blow up. For DPS, the additional scaling results in

$$\sigma(t)^2 \nabla_{\mathbf{x}_t} p(\mathbf{y} | \mathbf{x}_t) \approx \frac{\xi \sigma_y^2}{\|\mathbf{y} - \mathbf{A}\mathbf{x}_0\|} \frac{N}{\sigma_y^2} \vec{\mathbf{1}} = \frac{\xi N}{\|\vec{\mathbf{1}}\|} \vec{\mathbf{1}} = \frac{\xi N}{\sqrt{N}} \vec{\mathbf{1}} = \xi \sqrt{N} \vec{\mathbf{1}} \quad (64)$$

which is less severe than PIGDM, but still requires additional clipping unless the scale ξ is set to very low values.

Solution with the correct covariance In the main text, we showed that the issue does not show up in the case $\sigma_y = 0$. This resulted in:

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t) = (\mathbf{y} - \mathbf{x}_0(\mathbf{x}_t))^\top \text{Cov}[\mathbf{x}_0 | \mathbf{x}_t]^{-1} \frac{\text{Cov}[\mathbf{x}_0 | \mathbf{x}_t]}{\sigma_t^2} = \frac{(\mathbf{y} - \mathbf{x}_0(\mathbf{x}_t))^\top}{\sigma_t^2} \quad (65)$$

The corresponding \mathbf{x}_0 estimate is:

$$\begin{aligned} \mathbf{x}_t + \sigma_t^2 \left(\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) + \frac{(\mathbf{y} - \mathbf{x}_0(\mathbf{x}_t))^\top}{\sigma_t^2} \right) \\ = \mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t] + \mathbf{y} - \mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t] = \mathbf{y} \end{aligned} \quad (66)$$

that is, the updated score points to the direction of the observation for all time steps t . For the case $t \rightarrow \infty$ and $\text{Cov}[\mathbf{x}_0 | \mathbf{x}_t] \approx \mathbf{J}$ and **not** assuming $\sigma_y = 0$, we can derive with the Sherman–Morrison formula that

$$(\text{Cov}[\mathbf{x}_0 | \mathbf{x}_t] + \sigma_y^2 \mathbf{I})^{-1} = (\mathbf{J} + \sigma_y^2 \mathbf{I})^{-1} = \frac{1}{\sigma_y^2} \mathbf{I} - \frac{1}{\sigma_y^4 + N\sigma_y^2} \mathbf{J}. \quad (67)$$

To simplify formulas, let us again assume that $\mathbf{y} - \mathbf{x}_0(\mathbf{x}_t) = a\bar{\mathbf{1}}$

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t) \approx a\bar{\mathbf{1}}^\top \left(\frac{1}{\sigma_y^2} \mathbf{I} - \frac{1}{\sigma_y^4 + N\sigma_y^2} \mathbf{J} \right) \frac{\mathbf{J}}{\sigma_t^2} \quad (68)$$

$$= a \left(\frac{1}{\sigma_y^2} - \frac{N}{\sigma_y^4 + N\sigma_y^2} \right) \bar{\mathbf{1}}^\top \frac{\mathbf{J}}{\sigma_t^2} \quad (69)$$

$$= a \frac{1}{\sigma_y^2 + N} \bar{\mathbf{1}}^\top \frac{\mathbf{J}}{\sigma_t^2} \quad (70)$$

$$= a \frac{1}{\sigma_y^2 + N} \frac{N}{\sigma_t^2} \bar{\mathbf{1}}^\top. \quad (71)$$

Here, again, since the inverse term is inversely dependent on N , the dependence of the last term on N is cancelled. In the case $\sigma_y^2 = 0$, we recover the exact same result as previously. With non-zero observation noise, the strength of the guidance becomes slightly smaller, reflecting the uncertainty about the underlying pure \mathbf{x}_0 value we have measured.

D FULL GUIDANCE ALGORITHMS FOR THE LINEAR-GAUSSIAN OBSERVATION MODEL

The algorithm [Alg. 3](#) details an implementation of the method with the Euler ODE solver. The algorithm [Alg. 4](#) is a more easily applicable implementation with any type of solver, including higher-order methods like the Heun method. In it, we instantiate a class at the beginning of sampling, and whenever a call to the denoiser / score model is made, it is passed to the class to calculate $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$ and update the covariance information. For image data, we only perform the space updates in $1 < \sigma(t) < 5$, as detailed in [App. J](#).

E ADDITIONAL TOY EXPERIMENT WITH CORRELATED DATA

To examine the effect mentioned in [Sec. 3.7](#), we constructed data $p(\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_0 | \mathbf{0}, (1-\rho)\mathbf{I} + \rho\mathbf{J})$, where \mathbf{J} is a matrix of ones and $\rho = 0.999$. We used an observation with noise $\sigma_y = 0.2$ and varied the dimension. We plot the variance of the generated samples in [Fig. 8](#). As expected, both IIGDM and DPS become overly confident as dimensionality increases. In contrast, our method, which explicitly accounts for data covariance, maintains correct uncertainty calibration across dimension counts. Note that the DPS results are obtained after tuning the guidance scale for this particular problem, making the comparisons somewhat favourable towards DPS.

F ADDITIONAL QUALITATIVE RESULTS

[Figure 9](#) shows the qualitative comparison with 30 Heun solver steps.

G QUANTIFYING THE ERROR IN THE COVARIANCE ESTIMATION

We study the case where we directly estimate the error in the denoiser covariance estimates for low-dimensional toy data, where this is directly feasible. We consider a 2D Gaussian mixture model with an likelihood function and posterior as shown in [Fig. 10](#).

We generate samples with four different methods, and compare the covariances true to the true covariance with the Frobenius norm. The methods are:

Algorithm 3: Free Hunch for Linear Inverse Problems (Euler solver, with diffusion parameters from (Karras et al., 2022))

Input: Linear operator \mathbf{A} , observation \mathbf{y} , noise σ_y

Input: Initial covariance Σ_{data} , score model s_θ

Input: Schedule params: $\sigma_{\min} = t_{\min} = 0.002$, $\sigma_{\max} = t_{\max} = 80$, $\rho = 7$, steps N

```

/* Define time discretization */
1  $t_i = (t_{\max}^\rho + \frac{i}{N-1}t_{\min}^\rho - t_{\max}^\rho)^\rho$  for  $i < N$ ,  $t_N = 0$  (Karras et al., 2022)
2 Initialize  $\mathbf{x}_t \sim \mathcal{N}(\mathbf{0}, \sigma_{\max}^2 \mathbf{I})$ 
3 Initialize  $\Sigma_{0|t}(\mathbf{x}_t) = \Sigma_{\text{data}}$ 
4 Initialize  $\mu_{0|t_i}^{\text{transferred}} = \text{null}$ 
5 Initialize  $\Delta \mathbf{x} = \text{null}$ 
6 for  $i = 1, \dots, N - 1$  do
7    $\sigma_{\text{curr}} = t_i$ ,  $\sigma_{\text{next}} = t_{i+1}$ 
8    $\Delta t = t_{i+1} - t_i$ 
9   /* New score and denoising mean evaluation */
10   $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) = s_\theta(\mathbf{x}_t, t_i)$ 
11   $\mu_{0|t}(\mathbf{x}_t) = \mathbf{x}_t + \sigma_{\text{curr}}^2 \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$ 
12  /* Space update for covariance */
13  if  $\mu_{0|t_i}^{\text{transferred}} \neq \text{null}$  And  $\Delta \mathbf{x} \neq \text{null}$  then
14     $\Delta \mathbf{e} = \sigma_{\text{curr}}^2 (\mu_{0|t}(\mathbf{x}_t) - \mu_{0|t_i}^{\text{transferred}})$ 
15     $\gamma = \frac{1}{\Delta \mathbf{e}^\top \Delta \mathbf{x}}$ 
16     $\Sigma_{0|t}(\mathbf{x}_{t_{\text{next}}}) = \Sigma_{0|t}(\mathbf{x}_t) - \frac{\Sigma_{0|t}(\mathbf{x}_t) \Delta \mathbf{x} \Delta \mathbf{x}^\top \Sigma_{0|t}(\mathbf{x}_t)}{\Delta \mathbf{x}^\top \Sigma_{0|t}(\mathbf{x}_t) \Delta \mathbf{x}} + \frac{\Delta \mathbf{e} \Delta \mathbf{e}^\top}{\Delta \mathbf{e}^\top \Delta \mathbf{x}}$ 
17  end
18  /* Reconstruction guidance */
19   $\nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) = (\mathbf{y} - \mathbf{A}\mu_{0|t}(\mathbf{x}_t))^\top (\mathbf{A}\Sigma_{0|t}(\mathbf{x}_t)\mathbf{A}^\top + \sigma_y^2 \mathbf{I})^{-1} \mathbf{A} \nabla_{\mathbf{x}_t} \mu_{0|t}(\mathbf{x}_t)$ 
20  /* Fall back to approximation if guidance too large */
21  if  $\|\sigma_{\text{curr}}^2 \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t)\| > 1$  then
22     $\nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) = (\mathbf{y} - \mathbf{A}\mu_{0|t}(\mathbf{x}_t))^\top (\mathbf{A}\Sigma_{0|t}(\mathbf{x}_t)\mathbf{A}^\top + \sigma_y^2 \mathbf{I})^{-1} \mathbf{A} \frac{\Sigma_{0|t}(\mathbf{x}_t)}{\sigma_{\text{curr}}^2}$ 
23  end
24  /* Update sample with Euler step */
25   $\Delta \mathbf{x} = -\sigma_{\text{curr}} (\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t)) \Delta t$ 
26   $\mathbf{x}_t = \mathbf{x}_t + \Delta \mathbf{x}$ 
27  /* Time update for mean */
28   $\Delta \sigma^2 = \sigma_{\text{next}}^2 - \sigma_{\text{curr}}^2$ 
29   $\mu_{0|t_{i+1}}^{\text{transferred}} = \mathbf{x}_t + \sigma_{\text{next}}^2 (\sigma_{\text{next}}^2 \mathbf{I} - \frac{\Delta \sigma^2}{\sigma_{\text{curr}}^2} \Sigma_{0|t}(\mathbf{x}_t))^{-1} (\mu_{0|t}(\mathbf{x}_t) - \mathbf{x}_t)$ 
30  /* Time update for covariance (moving to noise level  $\sigma_{\text{next}}$ ) */
31   $\Delta(\sigma^{-2}) = \sigma_{\text{next}}^{-2} - \sigma_{\text{curr}}^{-2}$ 
32   $\Sigma_{0|t}(\mathbf{x}_t)^{-1} = \Sigma_{0|t}(\mathbf{x}_t)^{-1} + \Delta(\sigma^{-2}) \mathbf{I}$ 
33   $\Sigma_{0|t}(\mathbf{x}_t) = (\Sigma_{0|t}(\mathbf{x}_t)^{-1})^{-1}$ 
34 end
35 return  $\mathbf{x}_t$ 

```

Algorithm 4: Free Hunch Guidance Class, applicable with any solver

```

1 class FreeHunchGuidance:
  /* Initialize with measurement model and data covariance */
2  constructor( $\mathbf{A}$ ,  $\mathbf{y}$ ,  $\sigma_y$ ,  $\Sigma_{\text{data}}$ ):
3    Store  $\mathbf{A}$ ,  $\mathbf{y}$ ,  $\sigma_y$ ,  $\Sigma_{0|t} = \Sigma_{\text{data}}$ 
4    Initialize  $\mu_{\text{prev}} = \text{null}$ ,  $\mathbf{x}_{\text{prev}} = \text{null}$ ,  $\sigma_{\text{prev}} = \text{null}$ 

  /* Process new denoiser evaluation and return guidance, to be
  used for updating  $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$  to  $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t)$ 
  before using it in the solver. */
5  function process_denoiser( $\mu_{\text{new}}$ ,  $\mathbf{x}_{\text{new}}$ ,  $\sigma_{\text{new}}$ ):
6    if  $\mu_{\text{prev}} \neq \text{null}$  then
7      /* Time update from previous step */
8       $\Delta(\sigma^{-2}) = \sigma_{\text{new}}^{-2} - \sigma_{\text{prev}}^{-2}$ 
9       $\Delta\sigma^2 = \sigma_{\text{new}}^2 - \sigma_{\text{prev}}^2$ 
10      $\Sigma_{0|t}^{-1} = \Sigma_{0|t}^{-1} + \Delta(\sigma^{-2})\mathbf{I}$ 
11      $\Sigma_{0|t} = (\Sigma_{0|t}^{-1})^{-1}$ 

     /* Transfer previous mu to new noise level */
12      $\mu_{\text{transferred}} = \mathbf{x}_{\text{prev}} + \sigma_{\text{new}}^2(\sigma_{\text{next}}^2\mathbf{I} - \frac{\Delta\sigma^2}{\sigma_{\text{prev}}^2}\Sigma_{0|t})^{-1}(\mu_{\text{prev}} - \mathbf{x}_{\text{prev}})$ 

     /* Space update */
13      $\Delta\mathbf{x} = \mathbf{x}_{\text{new}} - \mathbf{x}_{\text{prev}}$ 
14      $\Delta\mathbf{e} = \sigma_{\text{new}}^2(\mu_{\text{new}} - \mu_{\text{transferred}})$ 
15      $\gamma = \frac{1}{\Delta\mathbf{e}^\top\Delta\mathbf{x}}$ 
16      $\Sigma_{0|t} = \Sigma_{0|t} - \frac{\Sigma_{0|t}\Delta\mathbf{x}\Delta\mathbf{x}^\top\Sigma_{0|t}}{\Delta\mathbf{x}^\top\Sigma_{0|t}\Delta\mathbf{x}} + \frac{\Delta\mathbf{e}\Delta\mathbf{e}^\top}{\Delta\mathbf{e}^\top\Delta\mathbf{x}}$ 
17   end

   /* Calculate reconstruction guidance */
18    $\nabla_{\mathbf{x}} \log p(\mathbf{y} | \mathbf{x}_{\text{new}}) = (\mathbf{y} - \mathbf{A}\mu_{\text{new}})^\top (\mathbf{A}\Sigma_{0|t}\mathbf{A}^\top + \sigma_y^2\mathbf{I})^{-1} \mathbf{A}\nabla_{\mathbf{x}}\mu_{\text{new}}$ 

   /* Fall back to approximation if guidance too large */
19   if  $\|\sigma_{\text{new}}^2 \nabla_{\mathbf{x}} \log p(\mathbf{y} | \mathbf{x}_{\text{new}})\| > 1$  then
20      $\nabla_{\mathbf{x}} \log p(\mathbf{y} | \mathbf{x}_{\text{new}}) = (\mathbf{y} - \mathbf{A}\mu_{\text{new}})^\top (\mathbf{A}\Sigma_{0|t}\mathbf{A}^\top + \sigma_y^2\mathbf{I})^{-1} \mathbf{A}\frac{\Sigma_{0|t}}{t_{\text{new}}^2}$ 
21   end

   /* Update state variables */
22    $\mu_{\text{prev}} = \mu_{\text{new}}$ 
23    $\mathbf{x}_{\text{prev}} = \mathbf{x}_{\text{new}}$ 
24    $\sigma_{\text{prev}} = \sigma_{\text{new}}$ 
return  $\nabla_{\mathbf{x}} \log p(\mathbf{y} | \mathbf{x}_{\text{new}})$ 

```

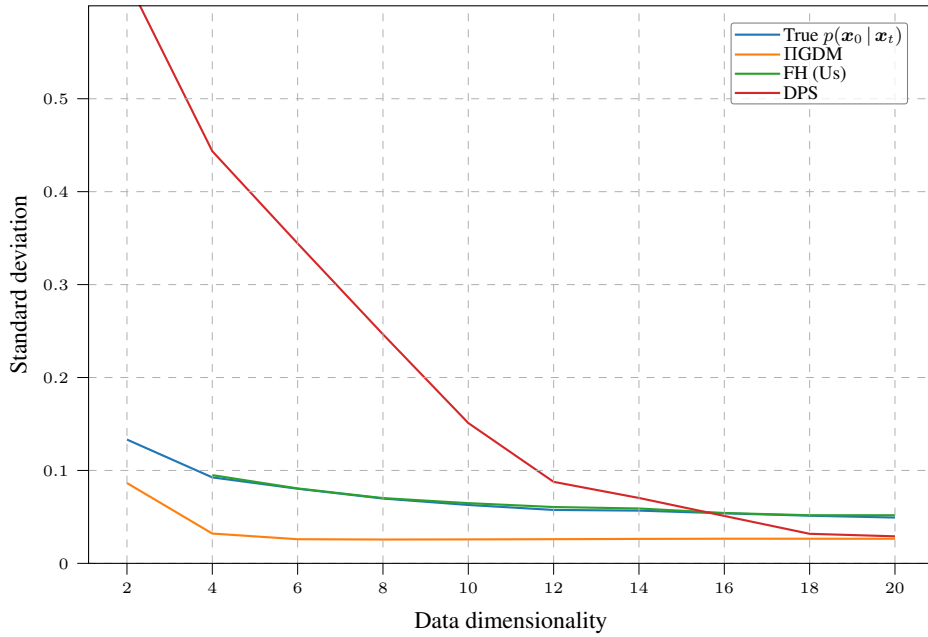


Figure 8: The standard deviation of posterior samples from different methods for the toy data discussed in App. E, showcasing the overconfidence problem caused by overestimated $\nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t)$.

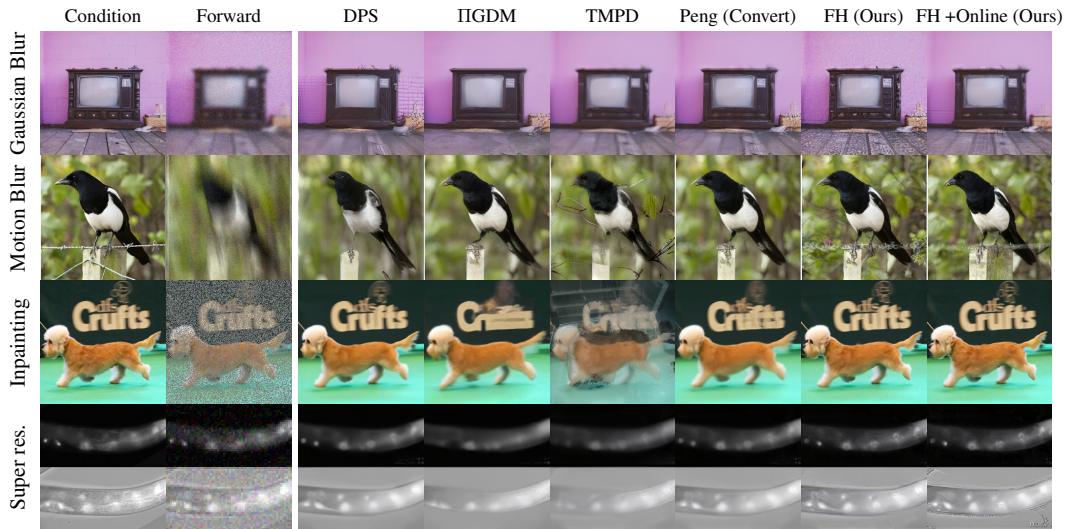


Figure 9: Qualitative results from 30-step Heun sampler.

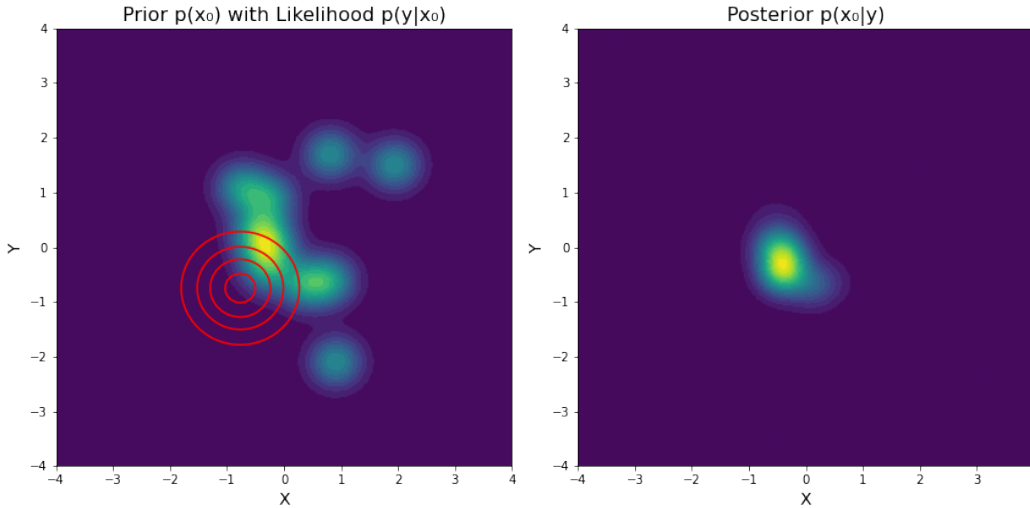


Figure 10: The posterior distribution of the toy data discussed in App. G and the prior distribution.

1. $\text{Cov}[\mathbf{x}_0 | \mathbf{x}_t] \approx \frac{\sigma(t)^2}{1+\sigma(t)^2} \mathbf{I}$, similarly to Π GDM.
2. $\text{Cov}[\mathbf{x}_0 | \mathbf{x}_t]$ approximated with our method by initialising at the data covariance, but not performing space updates.
3. $\text{Cov}[\mathbf{x}_0 | \mathbf{x}_t]$ approximated with our method by initialising at the data covariance, and performing space updates.
4. $\text{Cov}[\mathbf{x}_0 | \mathbf{x}_t]$ approximated with our method and with space updates, but estimating the BFGS updates by calculating $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t + \Delta \mathbf{x})$ and $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$ explicitly, requiring 2 denoiser calls per step, but without error from the time updates affecting the BFGS updates. This method is also discussed in App. J.

For these experiments, we used the Euler–Maruyama sampler. The results are shown in Fig. 11. The Π GDM covariance approximation is the furthest from the true value. Initialising at the data covariance helps, and adding the space updates decreases the error further. However, there is a clear gap in the standard method and using two score evaluations per step. We believe that this is due to the errors from the time updates affecting the BFGS updates.

We also perform an ablation comparing a deterministic Euler sampler and the stochastic Euler–Maryama sampler with the fourth method, and varying diffusion step count. The results are shown in Fig. 12. Whereas for the deterministic sampler, the covariance estimate does not significantly improve after a certain point, for the stochastic sampler, the error approaches zero across the sampling steps with more steps. This is because the reason the inherent curvature in the ODE path does not significantly change after increasing the step count above a certain limit. With a stochastic sampler, the generative path explores a slightly different direction at each step, and the BFGS updates get information from the curvature in all directions.

H ADDITIONAL RESULTS ON IMAGE MODELS

Here we list full additional results from:

1. ImageNet 256×256 results with Euler solver for 15, 30, 50, and 100 steps, in Table 2.
2. ImageNet 256×256 results with Heun solver for 15, 30, 50, and 100 steps, in Table 3.

For each task, we use 1000 samples from the ImageNet test set. We also evaluate DDNM+(Wang et al., 2023) and DiffPIR(Zhu et al., 2023) for a more thorough to different types of methods, although these can not directly be interpreted as reconstruction guidance with specific covariances. For DiffPIR, we use the implementation of (Peng et al., 2024), where they

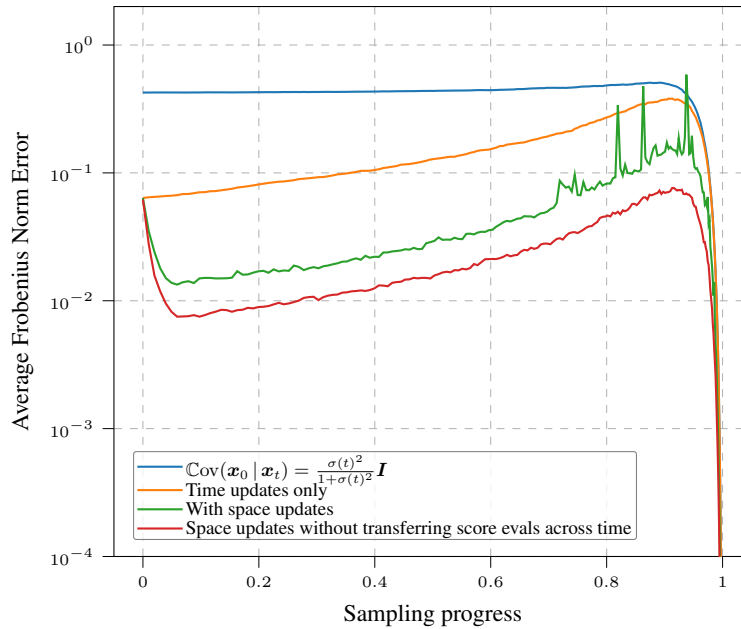


Figure 11: The Frobenius norm of the difference between the true covariance and the estimated covariance for different methods for different sampling steps (0 corresponds to the maximum noise level and 1 corresponds to zero noise level). As pointed out in App. J, using the time updates to transfer scores for use with the BFGS updates can cause inaccuracies with low-dimensional data, making the curve slightly rough, although still below the one with only time updates.

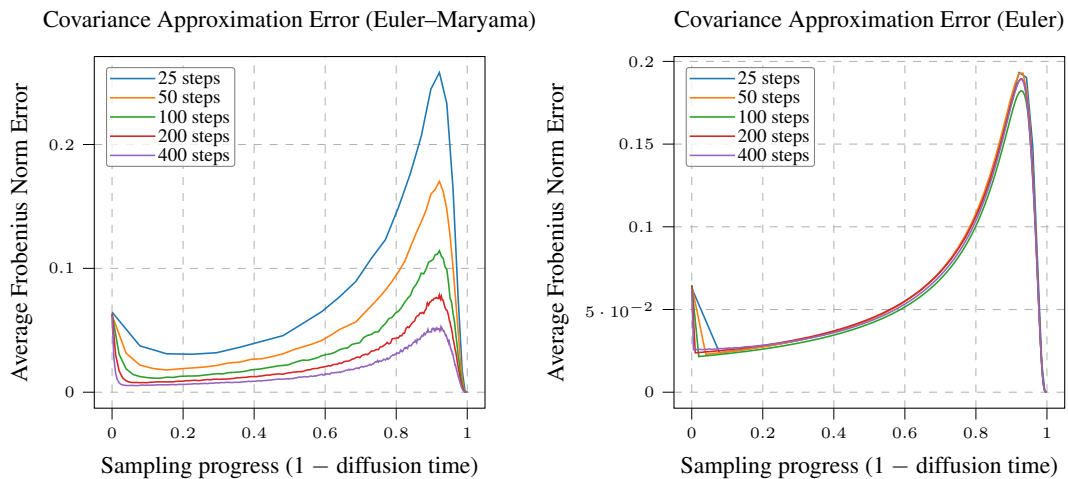


Figure 12: The Frobenius norm of the difference between the true covariance and the estimated covariance for different methods, with varying step count.

Notes on hyperparameters For DPS, Identity, Identity+Online updates, and DiffPIR, we tuned hyperparameters for each task using Gaussian blur as a baseline. We separately tuned the hyperparameters for the Euler and Heun solvers, and for each step count. While the optimal hyperparameters were similar for DiffPIR, Identity and Identity+Online updates, for DPS, the optimal values depended on the solver type and step count. We used 100 samples from the ImageNet validation set for tuning, and used these parameters for all experiments. The results are shown in Fig. 13, Fig. 14, Fig. 15 and Fig. 16.

I MOTIVATIONS FOR THE DCT BASIS AND THE BFGS UPDATE

The reason that we chose the DCT over, e.g., the DFT basis is that it is purely real-valued, does not assume periodic boundaries, and in practice needs less coefficients to efficiently represent natural images. This is also one of the reasons for its use in the JPEG compression standard (Wallace, 1991). The BFGS update has the attractive property of preserving positive-semidefiniteness (as opposed to, e.g., the symmetric rank-1 update). This combines well with performing the updates in the denoiser covariance, which is positive-definite (as opposed to the Hessian). Compared to Davidon-Fletcher-Powell (DFP), the difference is that the BFGS update minimizes a weighted Frobenius norm for the size of the update in inverse covariance (Dennis & Moré, 1977), instead of the covariance directly. In the update formula in Eq. (23), if we use $A=I$ and the observation noise is low, the inverse term is simply the inverse covariance. Thus, it could stabilise the updates across iterations, but this is more speculative, and DFP could work in practice as well.

J IMPLEMENTATION DETAILS

The solver We noticed that in large noise levels, it does not matter if the inversion in Eq. (23) is not exact, and we can set the tolerance quite high. We then defined a schedule for the tolerance such that it becomes lower towards the end. A lower tolerance towards the end of sampling is not an issue, since the covariance becomes closer to a diagonal, and the required matrix inverse becomes easier to calculate. In practice, we use the following schedule:

$$\begin{aligned} \sigma_{\max} &= 80 \\ \sigma_{\min} &= 1 \\ \text{rtol}_{\max} &= 1 \\ \text{rtol}_{\min} &= 1e - 14 \\ p &= 0.1 \\ \sigma_{\text{clipped}} &= \max(\min(\sigma, \sigma_{\max}), \sigma_{\min}) \\ \log_factor &= \left(\frac{\log_{10}(\sigma_{\text{clipped}}) - \log_{10}(\sigma_{\min})}{\log_{10}(\sigma_{\max}) - \log_{10}(\sigma_{\min})} \right)^p \end{aligned} \tag{72}$$

$$\log_rtol = \log_factor \cdot (\log_{10}(\text{rtol}_{\max}) - \log_{10}(\text{rtol}_{\min})) + \log_{10}(\text{rtol}_{\min}) \tag{73}$$

$$\text{rtol} = 10^{\log_rtol}, \tag{74}$$

where rtol_{\max} and rtol_{\min} control the maximum and minimum relative tolerances of the solver, and σ_{\max} and σ_{\min} control the noise levels outside of which the tolerances are rtol_{\max} and rtol_{\min} , respectively.

Note that scheduling the solver in this way does not improve image quality. Instead, it improves inference speed considerably, to the point where the solver is not a bottleneck anymore. Instead of using a standard off-the-shelf conjugate gradient implementation, we implemented one ourselves in PyTorch to utilize the speedup from the GPU.

Also for TMPD, we created a schedule for the conjugate gradient since a constant low tolerance slowed down the computation quite a bit. For TMPD, we use a standard scipy implementation that

Table 2: Results with the **Euler solver**. Our model performs especially well at small step sizes and remains competitive at larger step counts as well. DDNM+ is designed to enforce consistency with the measurement in cases where the measurement operator has a clearly defined nullspace, such as inpainting and super-resolution, potentially affecting the good PSNR and SSIM results there. In contrast, DDNM+ struggles with our Gaussian blur kernels. Motion blur results are not presented, as the code of DDNM+ assumes separable kernels.

Method	Deblur (Gaussian)			Inpainting (Random)			Deblur (Motion)			Super res. (4 \times)			
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	
15 steps	DPS	19.94	0.444	0.572	20.68	0.494	0.574	17.02	0.354	0.646	19.85	0.460	0.590
	IIGDM	20.29	0.474	0.574	19.87	0.468	0.598	19.21	0.429	0.602	20.17	0.474	0.582
	TMPD	22.56	0.572	0.486	17.70	0.447	0.589	20.40	0.481	0.567	21.15	0.517	0.541
	Peng Convert	22.53	0.563	0.490	22.23	0.579	0.489	20.46	0.475	0.556	21.92	0.541	0.517
	Peng Analytic	22.52	0.563	0.490	22.14	0.574	0.494	20.46	0.475	0.556	21.92	0.541	0.517
	DDNM+	7.21	0.029	0.822	23.95	0.667	0.352	–	–	–	24.30	0.669	0.398
	DiffPIR	22.77	0.575	0.403	16.10	0.284	0.661	19.75	0.381	0.527	21.76	0.540	0.436
	Identity	22.91	0.594	0.384	18.83	0.397	0.590	20.06	0.393	0.506	22.65	0.589	0.412
Identity+online	23.08	0.606	0.385	18.86	0.397	0.590	20.31	0.418	0.492	22.76	0.597	0.414	
FH	<u>23.41</u>	<u>0.625</u>	0.373	<u>24.76</u>	<u>0.702</u>	<u>0.327</u>	<u>21.69</u>	<u>0.534</u>	<u>0.447</u>	<u>23.39</u>	<u>0.632</u>	0.390	
FH+online	23.57	0.635	<u>0.378</u>	25.29	0.731	0.315	21.83	0.548	0.442	23.31	0.624	<u>0.393</u>	
30 steps	DPS	21.76	0.527	0.463	24.84	0.678	0.386	18.22	0.389	0.582	23.00	0.593	0.440
	IIGDM	22.27	0.559	0.468	21.24	0.518	0.517	21.15	0.508	0.503	22.11	0.552	0.479
	TMPD	22.92	0.591	0.451	18.27	0.465	0.563	20.71	0.495	0.538	21.59	0.536	0.507
	Peng Convert	<u>23.61</u>	0.627	0.405	23.74	0.648	0.403	21.99	0.553	0.463	23.21	0.608	0.430
	Peng Analytic	<u>23.61</u>	0.627	0.405	23.59	0.640	0.410	21.99	0.552	0.463	23.22	0.608	0.430
	DDNM+	7.51	0.033	0.814	26.66	0.769	0.272	–	–	–	24.09	0.657	0.418
	DiffPIR	22.34	0.552	0.404	15.94	0.262	0.667	19.38	0.368	0.523	21.25	0.512	0.443
	Identity	23.15	0.602	0.374	18.75	0.402	0.578	20.14	0.406	0.494	22.82	0.588	0.405
Identity+online	23.38	0.621	<u>0.359</u>	20.07	0.443	0.529	20.47	0.420	0.467	23.38	0.622	0.383	
FH	23.56	<u>0.630</u>	0.353	26.00	0.758	0.255	21.79	0.537	<u>0.410</u>	23.38	0.624	0.371	
FH+online	23.66	0.636	<u>0.359</u>	<u>26.17</u>	<u>0.766</u>	<u>0.268</u>	<u>21.89</u>	<u>0.548</u>	0.409	<u>23.46</u>	<u>0.629</u>	<u>0.375</u>	
50 steps	DPS	22.66	0.579	0.411	<u>26.31</u>	<u>0.761</u>	0.297	19.05	0.428	0.537	<u>23.79</u>	<u>0.642</u>	0.375
	IIGDM	22.64	0.577	0.434	21.67	0.536	0.484	21.56	0.526	0.468	22.48	0.571	0.442
	TMPD	23.09	0.600	0.434	18.50	0.472	0.551	20.83	0.500	0.524	21.76	0.543	0.492
	Peng Convert	23.81	0.638	0.377	24.75	0.698	0.346	22.30	0.567	<u>0.430</u>	23.44	0.622	0.400
	Peng Analytic	23.81	0.638	0.378	24.47	0.683	0.360	22.30	0.567	<u>0.430</u>	23.44	0.622	0.400
	DDNM+	7.83	0.038	0.806	27.15	0.771	0.301	–	–	–	24.05	0.655	0.424
	DiffPIR	22.10	0.539	0.407	15.82	0.251	0.670	19.17	0.358	0.525	21.00	0.498	0.448
	Identity	23.18	0.602	0.360	19.64	0.436	0.536	19.92	0.373	0.497	23.11	0.600	0.385
Identity+online	23.47	0.620	0.370	19.46	0.409	0.552	20.74	0.453	0.453	23.20	0.607	0.399	
FH	23.43	0.622	0.348	25.94	0.760	0.238	21.56	0.523	0.406	23.21	0.614	0.366	
FH+online	<u>23.59</u>	<u>0.631</u>	<u>0.353</u>	26.08	0.770	<u>0.252</u>	<u>21.71</u>	<u>0.534</u>	0.406	23.33	0.620	<u>0.369</u>	
100 steps	DPS	23.36	0.615	0.379	<u>26.61</u>	<u>0.800</u>	0.229	20.05	0.473	0.492	23.36	0.622	<u>0.366</u>
	IIGDM	22.76	0.585	0.408	21.97	0.551	0.459	<u>21.71</u>	<u>0.533</u>	0.440	22.63	0.580	0.416
	TMPD	23.18	0.604	0.423	18.67	0.477	0.543	20.90	0.502	0.514	21.88	0.548	0.480
	Peng Convert	23.74	0.638	0.358	24.89	0.706	0.332	22.36	0.570	0.407	<u>23.46</u>	<u>0.625</u>	0.379
	Peng Analytic	<u>23.73</u>	0.637	0.358	24.67	0.694	0.345	22.36	0.570	0.407	<u>23.46</u>	<u>0.625</u>	0.379
	DDNM+	8.77	0.054	0.786	28.28	0.809	0.278	–	–	–	23.55	0.630	0.450
	DiffPIR	21.88	0.527	0.410	15.69	0.241	0.672	18.95	0.345	0.529	20.78	0.485	0.451
	Identity	23.11	0.596	0.361	19.66	0.439	0.528	19.72	0.358	0.501	23.04	0.594	0.384
Identity+online	23.43	<u>0.616</u>	0.373	19.65	0.420	0.538	20.77	0.459	0.447	23.08	0.599	0.403	
FH	23.24	0.613	0.346	25.71	0.755	<u>0.233</u>	21.35	0.510	0.407	23.03	0.604	0.364	
FH+online	23.32	<u>0.616</u>	<u>0.356</u>	25.73	0.764	0.243	21.37	0.509	<u>0.417</u>	23.17	0.609	0.370	

Table 3: Results with the **Heun solver**. Our model performs especially well at small step sizes and remains competitive at larger step counts as well. DDNM+ is designed to enforce consistency with the measurement in cases where the measurement operator has a clearly defined nullspace, such as inpainting and super-resolution, potentially affecting the good PSNR and SSIM results there. In contrast, DDNM+ struggles with our Gaussian blur kernels. Motion blur results are not presented, as the code of DDNM+ assumes separable kernels.

Method	Deblur (Gaussian)			Inpainting (Random)			Deblur (Motion)			Super res. (4×)			
	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	
15 steps	DPS	19.94	0.444	0.572	20.68	0.494	0.574	17.02	0.354	0.646	19.85	0.460	0.590
	IIGDM	20.30	0.475	0.574	19.87	0.468	0.598	19.21	0.429	0.602	20.17	0.474	0.582
	TMPD	23.08	0.597	0.420	18.99	0.481	0.539	20.80	0.491	0.514	21.88	0.545	0.476
	Peng Convert	22.53	0.563	0.490	22.23	0.579	0.489	20.46	0.475	0.556	21.92	0.541	0.517
	Peng Analytic	22.53	0.563	0.490	22.14	0.574	0.494	20.46	0.475	0.556	21.92	0.541	0.517
	DDNM+	7.21	0.029	0.822	23.95	0.667	0.352	–	–	–	24.30	0.669	0.398
	DiffPIR	22.77	0.575	0.403	16.10	0.284	0.661	19.75	0.381	0.527	21.76	0.540	0.436
30 steps	Identity	22.91	0.594	0.384	18.83	0.397	0.590	20.06	0.393	0.506	22.65	0.589	0.412
	Identity+online	23.08	0.606	0.385	18.86	0.397	0.590	20.31	0.418	0.492	22.76	0.597	0.414
	FH	<u>23.39</u>	<u>0.624</u>	0.372	<u>24.73</u>	<u>0.701</u>	<u>0.327</u>	<u>21.69</u>	<u>0.534</u>	<u>0.446</u>	23.30	0.624	0.390
	FH+online	23.54	0.634	<u>0.378</u>	25.25	0.728	0.317	21.84	0.549	0.441	<u>23.39</u>	<u>0.632</u>	<u>0.394</u>
30 steps	DPS	21.76	0.527	0.463	24.84	0.678	0.387	18.22	0.389	0.582	23.00	0.593	0.440
	IIGDM	22.27	0.559	0.468	21.24	0.518	0.517	21.16	0.508	0.503	22.11	0.553	0.478
	TMPD	23.16	0.602	0.415	18.85	0.481	0.537	20.91	0.500	0.507	21.94	0.549	0.472
	Peng Convert	<u>23.61</u>	0.627	0.405	23.74	0.648	0.403	21.99	0.553	0.463	23.22	0.608	0.430
	Peng Analytic	<u>23.61</u>	0.626	0.405	23.59	0.640	0.411	21.99	<u>0.552</u>	0.463	23.21	0.608	0.430
	DDNM+	7.51	0.033	0.814	26.66	0.769	0.272	–	–	–	24.09	0.657	0.418
	DiffPIR	22.34	0.552	0.404	15.94	0.262	0.667	19.38	0.368	0.523	21.25	0.512	0.443
50 steps	Identity	23.15	0.602	0.374	18.75	0.402	0.578	20.14	0.406	0.494	22.82	0.588	0.405
	Identity+online	23.38	0.621	<u>0.359</u>	20.07	0.443	0.529	20.47	0.420	0.467	23.38	0.622	0.383
	FH	23.55	<u>0.630</u>	0.353	26.00	0.757	0.256	21.80	0.538	<u>0.411</u>	23.38	0.623	0.372
	FH+online	23.62	0.635	<u>0.358</u>	<u>26.18</u>	<u>0.767</u>	<u>0.268</u>	<u>21.88</u>	0.547	0.410	<u>23.44</u>	<u>0.628</u>	<u>0.375</u>
50 steps	DPS	22.66	0.578	0.412	<u>26.31</u>	<u>0.761</u>	0.296	19.05	0.428	0.537	<u>23.80</u>	<u>0.642</u>	0.375
	IIGDM	22.64	0.577	0.435	21.67	0.536	0.484	21.56	0.526	0.468	22.48	0.571	0.442
	TMPD	23.20	0.605	0.414	18.83	0.482	0.536	20.93	0.502	0.504	21.96	0.550	0.471
	Peng Convert	23.81	0.638	0.377	24.75	0.698	0.346	22.30	0.567	<u>0.429</u>	23.44	0.622	0.400
	Peng Analytic	<u>23.80</u>	0.638	0.378	24.47	0.683	0.360	22.30	0.567	0.430	23.44	0.622	0.400
	DDNM+	7.83	0.038	0.806	27.15	0.771	0.301	–	–	–	24.05	0.655	0.424
	DiffPIR	22.10	0.539	0.407	15.82	0.251	0.670	19.17	0.358	0.525	21.00	0.498	0.448
100 steps	Identity	23.18	0.602	0.360	19.64	0.436	0.536	19.92	0.373	0.497	23.11	0.600	0.385
	Identity+online	23.47	0.620	0.370	19.46	0.409	0.552	20.74	0.453	0.453	23.20	0.607	0.399
	FH	23.44	0.623	0.348	25.95	0.760	0.237	21.58	0.523	0.406	23.22	0.614	0.367
	FH+online	23.60	<u>0.631</u>	<u>0.353</u>	26.10	0.772	<u>0.250</u>	<u>21.73</u>	<u>0.535</u>	0.406	23.31	0.619	<u>0.369</u>
100 steps	DPS	23.36	0.615	0.378	<u>26.61</u>	<u>0.800</u>	0.228	20.05	0.473	0.492	23.36	0.622	<u>0.366</u>
	IIGDM	22.76	0.585	0.408	21.97	0.551	0.459	<u>21.71</u>	<u>0.533</u>	0.440	22.63	0.580	0.416
	TMPD	23.22	0.606	0.412	18.83	0.482	0.536	20.95	0.502	0.504	21.98	0.551	0.469
	Peng Convert	23.74	0.638	<u>0.358</u>	24.89	0.706	0.332	22.36	0.570	0.407	<u>23.46</u>	<u>0.625</u>	0.379
	Peng Analytic	<u>23.73</u>	0.637	<u>0.358</u>	24.67	0.694	0.345	22.36	0.570	0.407	<u>23.46</u>	<u>0.625</u>	0.379
	DDNM+	8.77	0.054	0.786	28.28	0.809	0.278	–	–	–	23.55	0.630	0.450
	DiffPIR	21.88	0.527	0.410	15.69	0.241	0.672	18.95	0.345	0.529	20.78	0.485	0.451
100 steps	Identity	23.11	0.596	0.361	19.66	0.439	0.528	19.72	0.358	0.501	23.04	0.594	0.384
	Identity+online	23.43	<u>0.616</u>	0.373	19.65	0.420	0.538	20.77	0.458	0.447	23.08	0.599	0.403
	FH	23.25	0.613	0.346	25.71	0.755	<u>0.233</u>	21.35	0.509	<u>0.408</u>	23.02	0.604	0.364
	FH+online	23.35	<u>0.616</u>	<u>0.357</u>	25.75	0.765	0.242	21.40	0.512	0.414	23.15	0.608	0.371

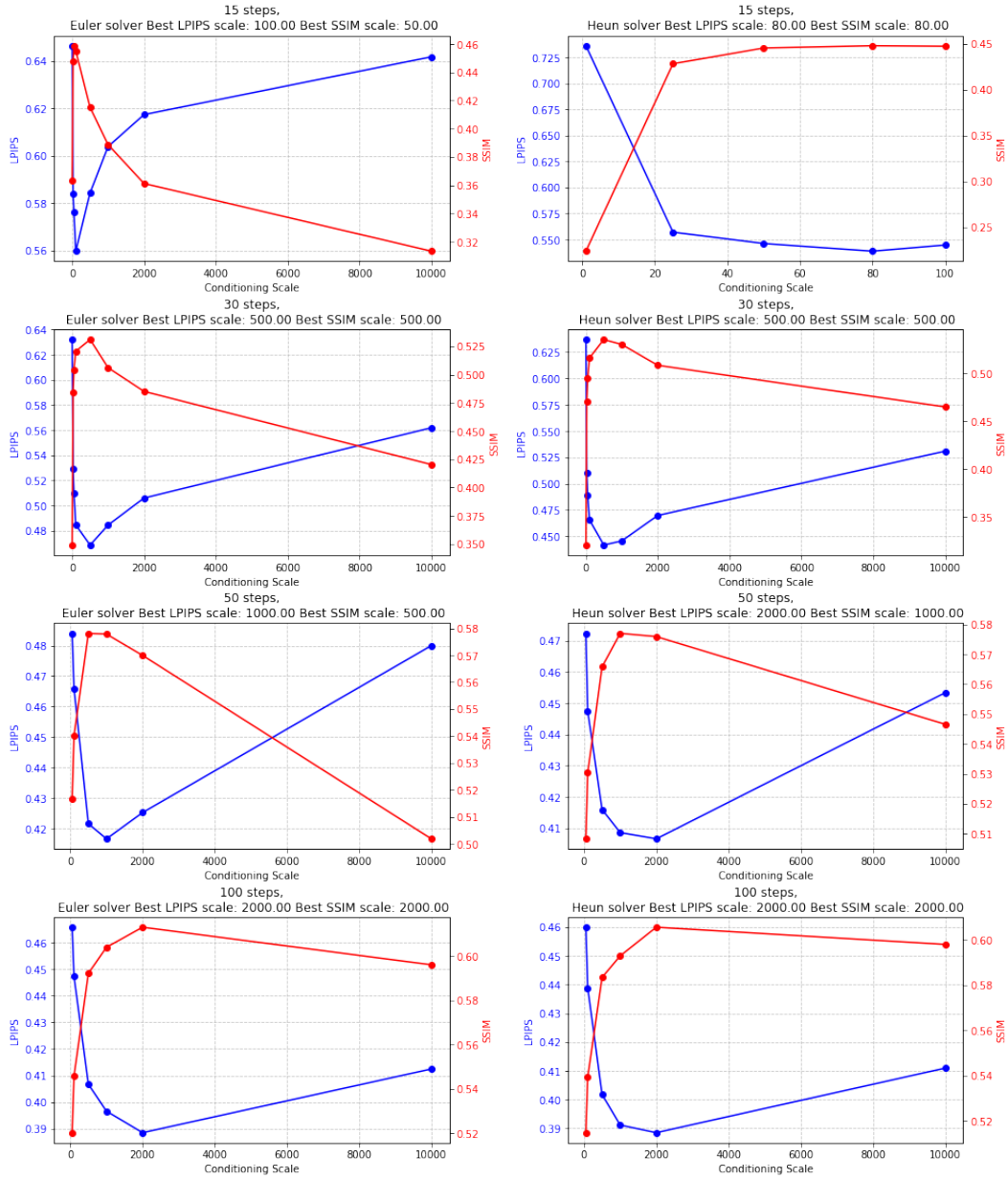


Figure 13: LPIPS and SSIM metrics across different solver steps and conditioning scales for DPS. The optimal LPIPS values are used in the experiments.

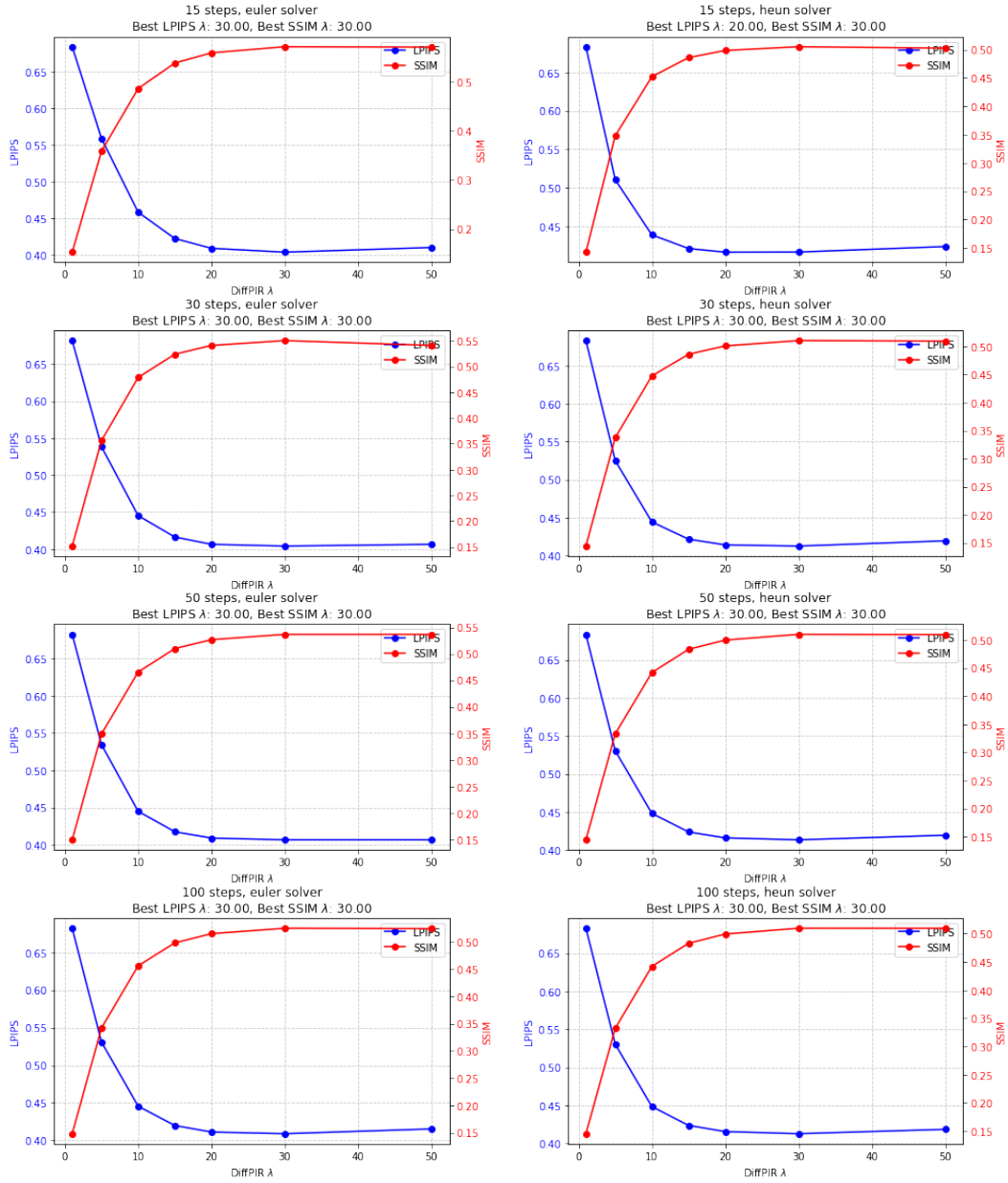


Figure 14: LPIPS and SSIM metrics across different solver steps and λ values for DiffPIR. The optimal LPIPS values are used in the experiments.

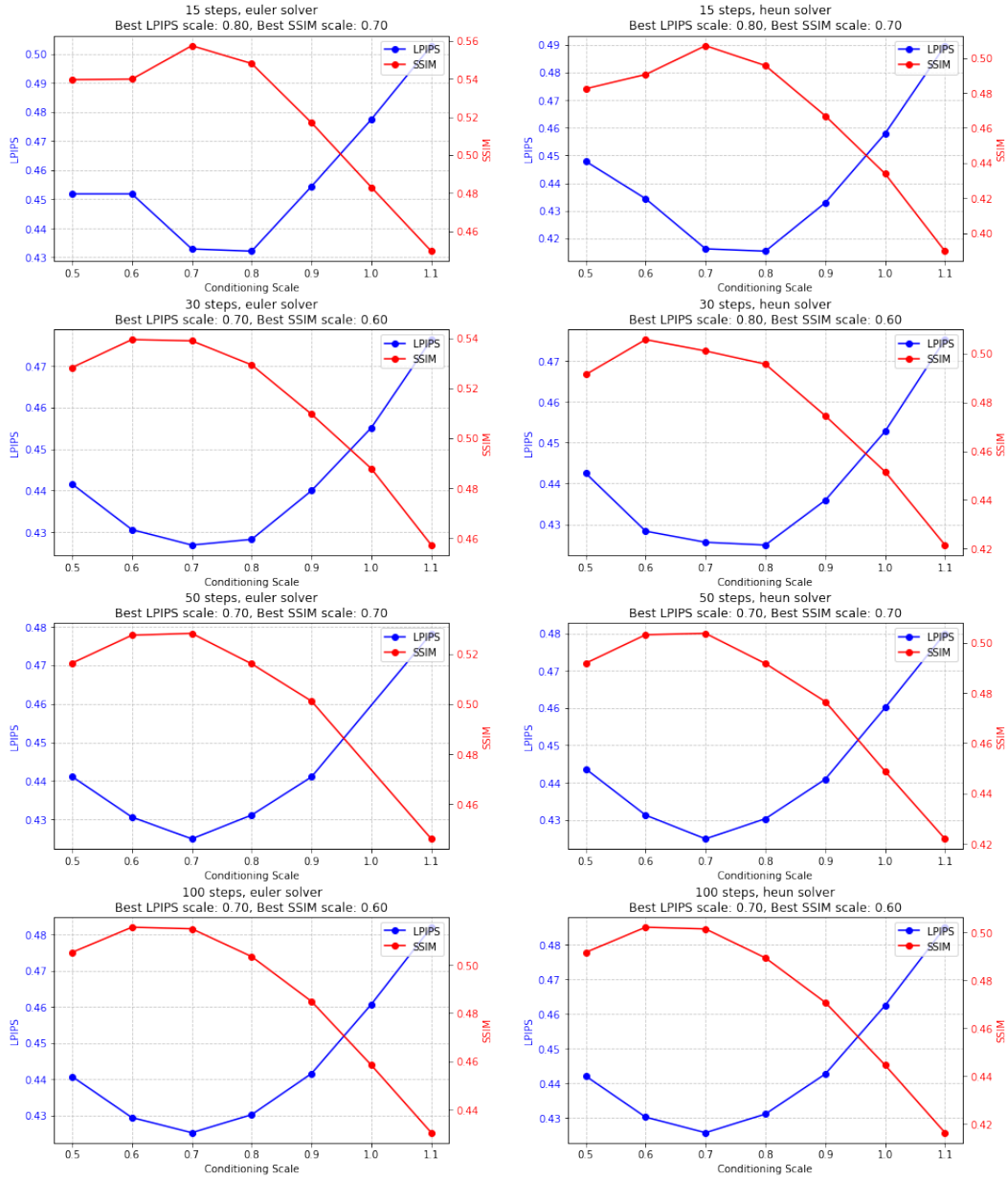


Figure 15: LPIPS and SSIM metrics across different solver steps and conditioning scales for our method with the identity base covariance. The optimal LPIPS values are used in the experiments.

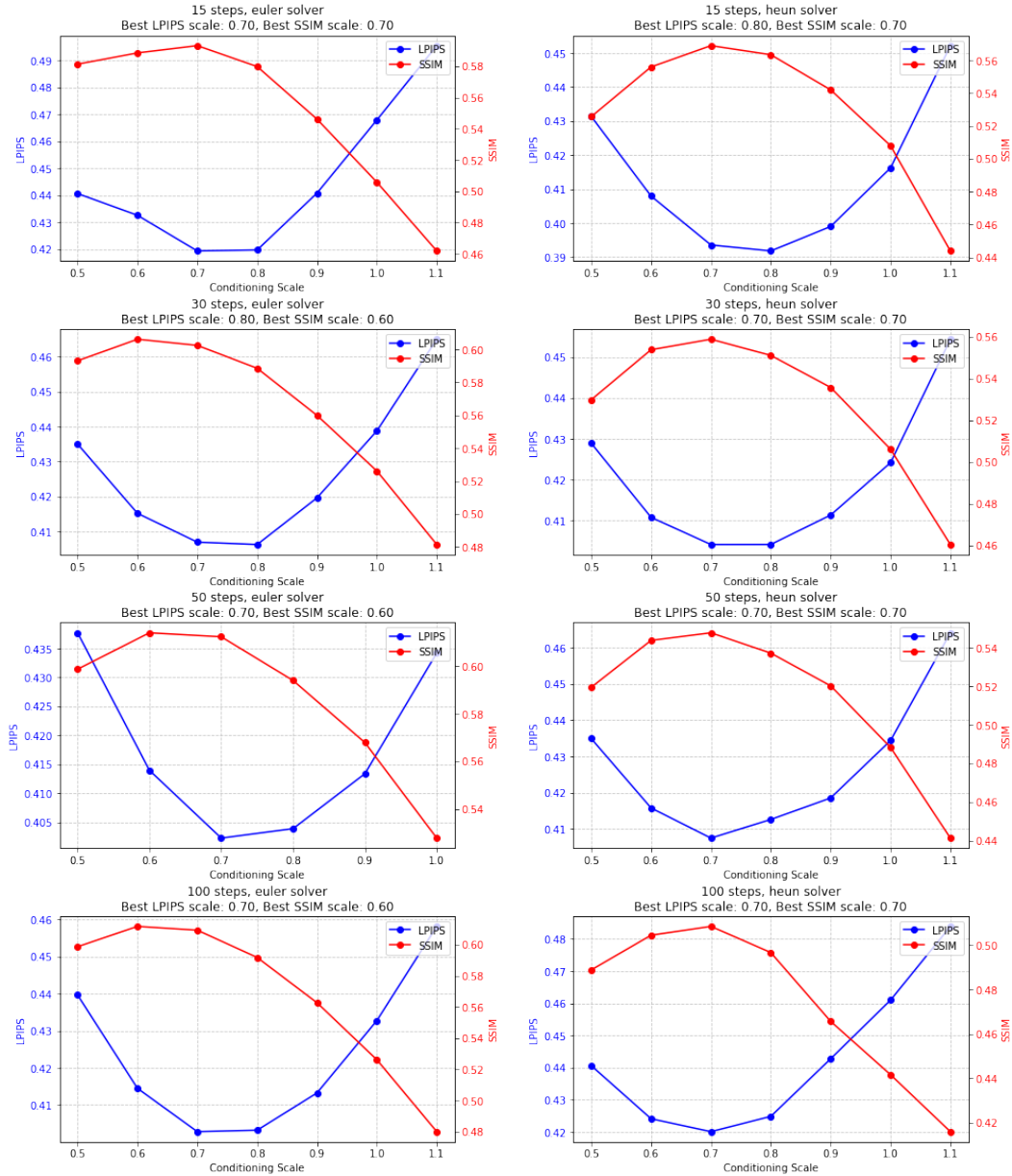


Figure 16: LPIPS and SSIM metrics across different solver steps and conditioning scales for our method with the identity base covariance and online updates. The optimal LPIPS values are used in the experiments.

is parameterized that is parameterized in terms of

$$\begin{aligned} \sigma_{\max} &= 80 \\ \sigma_{\min} &= 1 \\ \text{tol}_{\max} &= 1 \\ \text{tol}_{\min} &= 1e - 14 \\ p &= 0.05 \\ \sigma_{\text{clipped}} &= \max(\min(\sigma, \sigma_{\max}), \sigma_{\min}) \\ \log_factor &= \left(\frac{\log_{10}(\sigma_{\text{clipped}}) - \log_{10}(\sigma_{\min})}{\log_{10}(\sigma_{\max}) - \log_{10}(\sigma_{\min})} \right)^p \end{aligned} \quad (75)$$

$$\begin{aligned} \log_rtol &= \log_factor \cdot (\log_{10}(\text{rtol}_{\max}) - \log_{10}(\text{rtol}_{\min})) + \log_{10}(\text{rtol}_{\min}) \\ \text{rtol} &= 10^{\log_rtol} \end{aligned} \quad (76)$$

We tried to choose the schedule such that this does not degrade the performance noticeably, but also allows us to run experiments in reasonable time.

Range for the BFGS updates In practice, we do the space/BFGS updates for a range of $\sigma(t)$ values for image data, in particular $1 \leq \sigma(t) \leq 5$. A motivation for this choice is that we noticed the finite differences to not be numerically accurate at high noise levels, where the time updates Δt and the space updates $\Delta \mathbf{x}$ are large. For low noise levels, it is also unnecessary, given that the covariance approaches $\sigma(t)^2 I$ in any case. How to apply the space updates in the optimal way is an interesting direction for future research.

The solver We noticed that in large noise levels, it does not matter if the inversion in Eq. (23) is not exact, and we can set the tolerance quite high. We then defined a schedule for the tolerance such that it becomes lower towards the end. A lower tolerance towards the end of sampling is not an issue, since the covariance becomes closer to a diagonal, and the required matrix inverse becomes easier to calculate.

Details on the low-dimensional experiments In the Gaussian mixture experiments, we did not apply the time updates for $\mu_t(\mathbf{x})$, but instead evaluate $\mu_{0|t+\Delta t}(\mathbf{x})$ explicitly before applying the BFGS update. The reason is that on low-dimensional data, some of the prior samples are close to the actual data distribution. In that region, the time evolution is complex enough from the start that the denoiser mean time update coupled with the BFGS update in the next step sometimes causes numerical instability. To avoid additional score function evaluations, we could devise a schedule for when to apply the space updates.

Measurement operators. We obtained the measurement operator definitions from (Peng et al., 2024), which in turn are based on the operators in (Chung et al., 2023). We use a noise level $\sigma_y = 0.1$ for all measurement models (data scaled to $[-1,1]$).

K ADDITIONAL RESULTS ON OPTIMAL GUIDANCE STRENGTH

Figure 17 shows the PSNR, SSIM and LPIPS scores for the Gaussian deblurring task on ImageNet 256×256 with different post-hoc guidance scales on the initially calculated guidance term $\nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t)$.

L COMPUTATIONAL REQUIREMENTS

The sweep to obtain the results in Table 1 was done with multiple NVIDIA V100 GPUs in a few hours, and can be obtained with a single V100s in less than a day of compute. For some of the methods, the matrix inversion in Eq. (23) can slow down generation considerably, although this is not as significant an overhead with our tolerance-optimized conjugate gradient implementation.

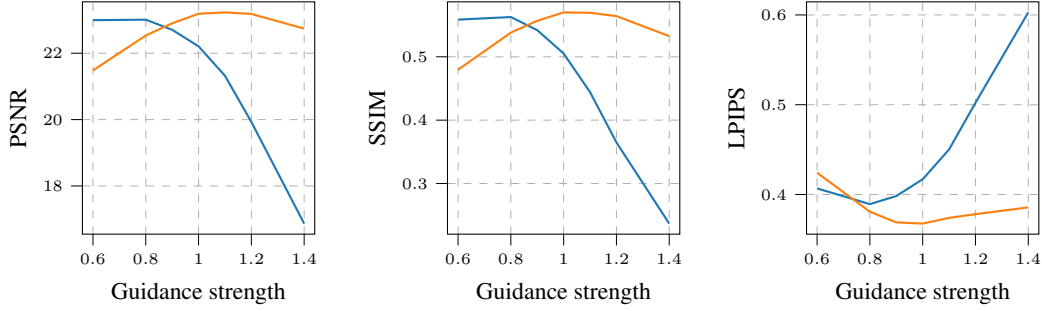


Figure 17: Different metrics with respect to guidance strength for 100 images from the Imagenet validation set, using an identity base covariance (blue) and a DCT-diagonal covariance (orange). With a better covariance approximation, the usefulness of adjusting the resulting guidance with post-hoc tricks becomes smaller.

M OPTIMAL $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y})$ FOR GAUSSIAN MIXTURE DATA AND GAUSSIAN OBSERVATION

In this section, we derive the optimal gradient $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y})$ for a situation with Gaussian mixture data and a Gaussian observation model. This is useful for performing toy experiments without having to retrain the model. Note that we can get the unconditional score from the end result by setting the observation noise Σ_y to infinity.

Model Definition We begin with the following components:

1. **Prior Distribution:** The prior on \mathbf{x}_0 is a Gaussian mixture model:

$$p(\mathbf{x}_0) = \sum_i w_i \mathcal{N}(\mathbf{x}_0 | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i), \quad (78)$$

where w_i are the mixture weights, $\boldsymbol{\mu}_i$ are the mean vectors, and $\boldsymbol{\Sigma}_i$ are the covariance matrices for each mixture component.

2. **Likelihood:** The observation model is Gaussian:

$$p(\mathbf{y} | \mathbf{x}_0) = \mathcal{N}(\mathbf{y} | \mathbf{x}_0, \boldsymbol{\Sigma}_y) \quad (79)$$

where $\boldsymbol{\Sigma}_y$ is the observation noise covariance.

3. **Transition Model:** The transition from \mathbf{x}_0 to \mathbf{x}_t is modeled as:

$$p(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t | \mathbf{x}_0, \sigma(t)^2 \mathbf{I}) \quad (80)$$

where $\sigma^2 \mathbf{I}$ is isotropic Gaussian noise with variance $\sigma(t)^2$.

M.1 POSTERIOR DISTRIBUTION

Given these components, the posterior distribution $p(\mathbf{x}_0 | \mathbf{x}_t, \mathbf{y})$ is also a Gaussian mixture:

$$p(\mathbf{x}_0 | \mathbf{x}_t, \mathbf{y}) = \sum_i w'_i \mathcal{N}(\mathbf{x}_0 | \boldsymbol{\mu}'_i, \boldsymbol{\Sigma}'_i), \quad (81)$$

where

$$\boldsymbol{\Sigma}'_i{}^{-1} = (\sigma(t)^2 \mathbf{I})^{-1} + \boldsymbol{\Sigma}_y^{-1} + \boldsymbol{\Sigma}_i^{-1} \quad (82)$$

$$\boldsymbol{\mu}'_i = \boldsymbol{\Sigma}'_i ((\sigma(t)^2 \mathbf{I})^{-1} \mathbf{x}_t + \boldsymbol{\Sigma}_y^{-1} \mathbf{y} + \boldsymbol{\Sigma}_i^{-1} \boldsymbol{\mu}_i) \quad (83)$$

$$w'_i \propto w_i \mathcal{N}(\mathbf{x}_t | \boldsymbol{\mu}_i, \sigma(t)^2 \mathbf{I} + \boldsymbol{\Sigma}_i) \mathcal{N}(\mathbf{y} | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_y + \boldsymbol{\Sigma}_i). \quad (84)$$

M.2 CONDITIONAL EXPECTATION

The conditional expectation of \mathbf{x}_0 given \mathbf{x}_t and \mathbf{y} is:

$$\mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t, \mathbf{y}] = \sum_i w'_i \boldsymbol{\mu}'_i \quad (85)$$

M.3 DERIVATION OF THE GRADIENT

Now, let's derive the gradient $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y})$:

1. We start with:

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}) = \nabla_{\mathbf{x}_t} \log \int p(\mathbf{x}_t | \mathbf{x}_0) p(\mathbf{x}_0 | \mathbf{y}) d\mathbf{x}_0 \quad (86)$$

2. Applying the chain rule and moving the gradient inside the integral:

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}) = \frac{\int \nabla_{\mathbf{x}_t} p(\mathbf{x}_t | \mathbf{x}_0) p(\mathbf{x}_0 | \mathbf{y}) d\mathbf{x}_0}{\int p(\mathbf{x}_t | \mathbf{x}_0) p(\mathbf{x}_0 | \mathbf{y}) d\mathbf{x}_0} \quad (87)$$

3. Given $p(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t | \mathbf{x}_0, \sigma(t)^2 \mathbf{I})$, we have:

$$\nabla_{\mathbf{x}_t} p(\mathbf{x}_t | \mathbf{x}_0) = -\frac{1}{\sigma(t)^2} (\mathbf{x}_t - \mathbf{x}_0) p(\mathbf{x}_t | \mathbf{x}_0) \quad (88)$$

4. Substituting this back:

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}) = -\frac{1}{\sigma(t)^2} \frac{\int (\mathbf{x}_t - \mathbf{x}_0) p(\mathbf{x}_t | \mathbf{x}_0) p(\mathbf{x}_0 | \mathbf{y}) d\mathbf{x}_0}{\int p(\mathbf{x}_t | \mathbf{x}_0) p(\mathbf{x}_0 | \mathbf{y}) d\mathbf{x}_0} \quad (89)$$

$$= -\frac{1}{\sigma(t)^2} \left(\mathbf{x}_t - \frac{\int \mathbf{x}_0 p(\mathbf{x}_t | \mathbf{x}_0) p(\mathbf{x}_0 | \mathbf{y}) d\mathbf{x}_0}{\int p(\mathbf{x}_t | \mathbf{x}_0) p(\mathbf{x}_0 | \mathbf{y}) d\mathbf{x}_0} \right) \quad (90)$$

$$= -\frac{1}{\sigma(t)^2} (\mathbf{x}_t - \mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t, \mathbf{y}]) \quad (91)$$

M.4 FINAL FORM OF THE GRADIENT

Therefore, the final form of the gradient is:

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}) = -\frac{1}{\sigma^2} (\mathbf{x}_t - \mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t, \mathbf{y}]) \quad (92)$$

M.5 DETAILED FORMULA FOR $\mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t, \mathbf{y}]$

Let's expand the formula for $\mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t, \mathbf{y}]$:

1. We start with the posterior distribution:

$$p(\mathbf{x}_0 | \mathbf{x}_t, \mathbf{y}) = \sum_i w'_i \mathcal{N}(\mathbf{x}_0 | \boldsymbol{\mu}'_i, \boldsymbol{\Sigma}'_i) \quad (93)$$

2. The expectation of this mixture is the weighted sum of the means:

$$\mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t, \mathbf{y}] = \sum_i w'_i \boldsymbol{\mu}'_i \quad (94)$$

3. Expanding $\boldsymbol{\mu}'_i$:

$$\boldsymbol{\mu}'_i = \boldsymbol{\Sigma}'_i ((\sigma(t)^2 \mathbf{I})^{-1} \mathbf{x}_t + \boldsymbol{\Sigma}_y^{-1} \mathbf{y} + \boldsymbol{\Sigma}_i^{-1} \boldsymbol{\mu}_i) \quad (95)$$

4. Substituting this into the expectation formula:

$$\mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t, \mathbf{y}] = \sum_i w'_i \boldsymbol{\Sigma}'_i ((\sigma(t)^2 \mathbf{I})^{-1} \mathbf{x}_t + \boldsymbol{\Sigma}_y^{-1} \mathbf{y} + \boldsymbol{\Sigma}_i^{-1} \boldsymbol{\mu}_i) \quad (96)$$

5. Rearranging:

$$\mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t, \mathbf{y}] = \left(\sum_i w'_i \boldsymbol{\Sigma}'_i (\sigma(t)^2 \mathbf{I})^{-1} \right) \mathbf{x}_t + \left(\sum_i w'_i \boldsymbol{\Sigma}'_i \boldsymbol{\Sigma}_y^{-1} \right) \mathbf{y} + \sum_i w'_i \boldsymbol{\Sigma}'_i \boldsymbol{\Sigma}_i^{-1} \boldsymbol{\mu}_i \quad (97)$$

Let's define:

$$\mathbf{A} = \sum_i w'_i \boldsymbol{\Sigma}'_i (\sigma(t)^2 \mathbf{I})^{-1}, \quad (98)$$

$$\mathbf{B} = \sum_i w'_i \boldsymbol{\Sigma}'_i \boldsymbol{\Sigma}_y^{-1}, \quad (99)$$

$$\mathbf{c} = \sum_i w'_i \boldsymbol{\Sigma}'_i \boldsymbol{\Sigma}_i^{-1} \boldsymbol{\mu}_i. \quad (100)$$

Then we can write the final formula as:

$$\mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t, \mathbf{y}] = \mathbf{A} \mathbf{x}_t + \mathbf{B} \mathbf{y} + \mathbf{c} \quad (101)$$

where:

$$w'_i \propto w_i \mathcal{N}(\mathbf{x}_t | \boldsymbol{\mu}_i, \sigma(t)^2 \mathbf{I} + \boldsymbol{\Sigma}_i) \mathcal{N}(\mathbf{y} | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_y + \boldsymbol{\Sigma}_i), \quad (102)$$

$$\boldsymbol{\Sigma}'_i = ((\sigma(t)^2 \mathbf{I})^{-1} + \boldsymbol{\Sigma}_y^{-1} + \boldsymbol{\Sigma}_i^{-1})^{-1}. \quad (103)$$

This formula shows that $\mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t, \mathbf{y}]$ is a linear combination of \mathbf{x}_t and \mathbf{y} , plus a constant term. The matrices \mathbf{A} and \mathbf{B} determine how much the expectation depends on \mathbf{x}_t and \mathbf{y} respectively, while \mathbf{c} represents a constant offset based on the prior distribution.