
Learning Diffusion Priors from Observations by Expectation Maximization

François Rozet
University of Liège
francois.rozet@uliege.be

Gérôme Andry
University of Liège
gandry@uliege.be

François Lanusse
Flatiron Institute
francois.lanusse@cnrs.fr

Gilles Louppe
University of Liège
g.louppe@uliege.be

Abstract

Diffusion models recently proved to be remarkable priors for Bayesian inverse problems. However, training these models typically requires access to large amounts of clean data, which could prove difficult in some settings. In this work, we present a novel method based on the expectation-maximization algorithm for training diffusion models from incomplete and noisy observations only. Unlike previous works, our method leads to proper diffusion models, which is crucial for downstream tasks. As part of our method, we propose and motivate a new posterior sampling scheme for unconditional diffusion models. We present empirical evidence supporting the effectiveness of our method.

1 Introduction

Many scientific applications can be formalized as Bayesian inference in latent variable models, where the target is the posterior distribution $p(x | y) \propto p(y | x)p(x)$ given an observation $y \in \mathbb{R}^M$ resulting from a forward process $p(y | x)$ and a prior distribution $p(x)$ over the latent variable $x \in \mathbb{R}^N$. Notable examples include gravitational lensing inversion [1–3], accelerated MRI [4–8], unfolding in particle physics [9, 10], and data assimilation [11–14]. In all of these examples, the observation y alone is either too incomplete or too noisy to recover the latent x . Additional knowledge in the form of an informative prior $p(x)$ is crucial for valuable inference.

Recently, diffusion models [15, 16] proved to be remarkable priors for Bayesian inference, demonstrating both quality and versatility [17–26]. However, to train a diffusion model for the latent x , one would typically need a large number of latent realizations, which by definition are not or rarely accessible. This is notably the case in earth and space sciences where the systems of interest can only be probed superficially.

Empirical Bayes (EB) methods [27–30] offer a solution to the problem of prior specification in latent variable models when only observations y are available. The objective of EB is to find the parameters θ of a prior model $q_\theta(x)$ for which the evidence distribution $q_\theta(y) = \int p(y | x) q_\theta(x) dx$ is closest to the empirical distribution of observations $p(y)$. Many EB methods have been proposed over the years, but they remain limited to low-dimensional settings [31–36] or simple parametric models [37, 38].

In this work, we use diffusion models for the prior $q_\theta(x)$, as they are best-in-class for modeling high-dimensional distributions and enable many downstream tasks, including Bayesian inference. This presents challenges for previous empirical Bayes methods which typically rely on models for which the density $q_\theta(x)$ or samples $x \sim q_\theta(x)$ are differentiable with respect to the parameters θ . Instead, we propose an adaptation of the expectation-maximization (EM) [39–43] algorithm where

we conduct the expectation step by sampling from the posterior $q_\theta(x | y)$, in the spirit of Monte Carlo EM [44–51]. As part of our method, we propose a new posterior sampling scheme for unconditional diffusion models, which we motivate theoretically and empirically.

2 Diffusion Models

The primary purpose of diffusion models (DMs) [15, 16], also known as score-based generative models [52, 53], is to generate plausible data from a distribution $p(x)$ of interest. Formally, adapting the continuous-time formulation of Song et al. [53], samples $x \in \mathbb{R}^N$ from $p(x)$ are progressively perturbed through a diffusion process expressed as a stochastic differential equation (SDE)

$$dx_t = f_t x_t dt + g_t dw_t \quad (1)$$

where $f_t \in \mathbb{R}$ is the drift coefficient, $g_t \in \mathbb{R}_+$ is the diffusion coefficient, $w_t \in \mathbb{R}^N$ denotes a standard Wiener process and $x_t \in \mathbb{R}^N$ is the perturbed sample at time $t \in [0, 1]$. Because the SDE is linear with respect to x_t , the perturbation kernel from x to x_t is Gaussian and takes the form

$$p(x_t | x) = \mathcal{N}(x_t | \alpha_t x, \Sigma_t) \quad (2)$$

where α_t and $\Sigma_t = \sigma_t^2 I$ are derived from f_t and g_t [53–56]. Crucially, the forward SDE (1) has an associated family of reverse SDEs [53–56]

$$dx_t = \left[f_t x_t - \frac{1 + \eta^2}{2} g_t^2 \nabla_{x_t} \log p(x_t) \right] dt + \eta g_t dw_t \quad (3)$$

where $\eta \geq 0$ is a parameter controlling stochasticity. In other words, we can draw noise samples $x_1 \sim p(x_1) \approx \mathcal{N}(0, \Sigma_1)$ and gradually remove the noise therein to obtain data samples $x_0 \sim p(x_0) \approx p(x)$ by simulating Eq. (3) from $t = 1$ to 0 using an appropriate discretization scheme [16, 52, 53, 56–59]. In this work, we adopt the variance exploding SDE [52] for which $f_t = 0$ and $\alpha_t = 1$.

In practice, the score function $\nabla_{x_t} \log p(x_t)$ in Eq. (3) is unknown, but can be approximated by a neural network trained via denoising score matching [60, 61]. Several equivalent parameterizations and objectives have been proposed for this task [16, 52, 53, 57–59]. In this work, we adopt the denoiser parameterization and its objective [58]

$$\arg \min_{\theta} \mathbb{E}_{p(x)p(t)p(x_t|x)} \left[\lambda_t \|d_\theta(x_t, t) - x\|_2^2 \right], \quad (4)$$

for which the optimal denoiser is the mean $\mathbb{E}[x | x_t]$ of $p(x | x_t)$. Importantly, $\mathbb{E}[x | x_t]$ is linked to the score function through Tweedie’s formula [62–65]

$$\mathbb{E}[x | x_t] = x_t + \Sigma_t \nabla_{x_t} \log p(x_t), \quad (5)$$

which allows to use $s_\theta(x_t) = \Sigma_t^{-1}(d_\theta(x_t, t) - x_t)$ a score estimate in Eq. (3).

3 Expectation-Maximization

The objective of the expectation-maximization (EM) algorithm [39–43] is to find the parameters θ of a latent variable model $q_\theta(x, y)$ that maximize the log-evidence $\log q_\theta(y)$ of an observation y . For a distribution of observations $p(y)$, the objective is to maximize the expected log-evidence [42, 43] or, equivalently, to minimize the Kullback-Leibler (KL) divergence between $p(y)$ and $q_\theta(y)$. That is,

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{p(y)} [\log q_\theta(y)] \quad (6)$$

$$= \arg \min_{\theta} \text{KL}(p(y) \| q_\theta(y)). \quad (7)$$

The key idea behind the EM algorithm is that for any two sets of parameters θ_1 and θ_2 , we have

$$\log q_{\theta_2}(y) - \log q_{\theta_1}(y) = \log \mathbb{E}_{q_{\theta_1}(x|y)} \left[\frac{q_{\theta_2}(x, y)}{q_{\theta_1}(x, y)} \right] \quad (8)$$

$$\geq \mathbb{E}_{q_{\theta_1}(x|y)} [\log q_{\theta_2}(x, y) - \log q_{\theta_1}(x, y)] \quad (9)$$

by Jensen’s inequality. This inequality also holds in expectation over $p(y)$. Therefore, starting from arbitrary parameters θ_1 , the EM update

$$\theta_{k+1} = \arg \max_{\theta} \mathbb{E}_{p(y)} \mathbb{E}_{q_{\theta_k}(x|y)} [\log q_{\theta}(x, y) - \log q_{\theta_k}(x, y)] \quad (10)$$

$$= \arg \max_{\theta} \mathbb{E}_{p(y)} \mathbb{E}_{q_{\theta_k}(x|y)} [\log q_{\theta}(x, y)] \quad (11)$$

leads to a sequence of parameters θ_k for which the expected log-evidence $\mathbb{E}_{p(y)} [\log q_{\theta_k}(y)]$ is monotonically increasing and converges to a local optimum [41–43].

When the expectation in Eq. (11) is intractable, many have proposed to use Monte Carlo approximations instead [44–51]. Previous approaches include Markov chain Monte Carlo (MCMC) sampling, importance sampling, rejection sampling and their variations [48–51]. A perhaps surprising advantage of Monte Carlo EM (MCEM) algorithms is that they may be able to overcome local optimum traps [45, 46]. We refer the reader to Ruth [51] for a recent review of MCEM algorithms.

4 Methods

Although rarely mentioned in the literature, the expectation-maximization algorithm is a possible solution to the empirical Bayes problem. Indeed, both have the same objective: minimizing the KL between the empirical distribution of observations $p(y)$ and the evidence $q_{\theta}(y)$. In the empirical Bayes setting, the forward model $p(y | x)$ is known and only the parameters of the prior $q_{\theta}(x)$ should be optimized. In this case, Eq. (11) becomes

$$\theta_{k+1} = \arg \max_{\theta} \mathbb{E}_{p(y)} \mathbb{E}_{q_{\theta_k}(x|y)} [\log q_{\theta}(x) + \log p(y | x)] \quad (12)$$

$$= \arg \max_{\theta} \mathbb{E}_{p(y)} \mathbb{E}_{q_{\theta_k}(x|y)} [\log q_{\theta}(x)] \quad (13)$$

$$= \arg \min_{\theta} \text{KL}(\pi_k(x) \| q_{\theta}(x)) \quad (14)$$

where $\pi_k(x) = \int q_{\theta_k}(x | y) p(y) dy$. Intuitively, $\pi_k(x)$ assigns less density to regions which are inconsistent with observations. In this work, we consider a special case of the empirical Bayes problem where each observation y has an associated measurement matrix $A \sim p(A)$ and the forward process takes a linear Gaussian form $p(y | x, A) = \mathcal{N}(y | Ax, \Sigma_y)$. This formulation allows the forward process to be potentially different for each observation y . As a result, $\pi_k(x)$ in Eq. (14) becomes $\pi_k(x) = \int q_{\theta_k}(x | y, A) p(y, A) dy$.

4.1 Pipeline

Now that our goals and assumptions are set, we present our method to learn a diffusion model $q_{\theta}(x)$ for the latent x from observations y by expectation-maximization. The idea is to decompose Eq. (14) into (i) generating a dataset of i.i.d. samples from $\pi_k(x)$ and (ii) training $q_{\theta}(x)$ to reproduce the generated dataset. We summarize the pipeline in Algorithms 1, 2 and 3, provided in Appendix A due to space constraints.

Expectation To draw from $\pi_k(x)$, we first sample a pair $(y, A) \sim p(y, A)$ and then generate $x \sim q_{\theta_k}(x | y, A)$ from the posterior. Unlike previous MCEM algorithms that rely on expensive and hard to tune sampling strategies [48–51], the use of a diffusion model enables efficient and embarrassingly parallelizable posterior sampling [21–23]. However, the quality of posterior samples is critical for the EM algorithm to converge properly [48–51] and, in this regard, we find previous posterior sampling methods [21–23] to be unsatisfactory (see Section 5). Therefore, we propose a new posterior sampling scheme, named moment matching posterior sampling (MMPS), which we present and motivate in Section 4.2.

Maximization We parameterize our diffusion model $q_{\theta}(x)$ by a denoiser network $d_{\theta}(x_t, t)$ and train it via denoising score matching [60, 61], as presented in Section 2. To accelerate the training, we start each iteration with the previous parameters θ_k .

Initialization An important part of our pipeline is the initial model $q_0(x)$. Any initial model leads to a local optimum [41–43], but an informed initial model can reduce the number of iterations until

convergence. In our approach, we take a Gaussian distribution $\mathcal{N}(x \mid \mu_x, \Sigma_x)$ as initial model and fit its mean and covariance by – you guessed it! – expectation-maximization. Fitting a Gaussian model by EM is very fast as the maximization step can be conducted in closed-form, especially for low-rank covariance approximations [66].

4.2 Moment Matching Posterior Sampling

In the diffusion model paradigm, to sample from a posterior distribution $p(x \mid y)$, we have to estimate the posterior score $\nabla_{x_t} \log p(x_t \mid y)$ and plug it into the reverse SDE (3). Thanks to Bayes’ rule, the posterior score can be decomposed into two terms [17, 18, 21–25, 53]

$$\nabla_{x_t} \log p(x_t \mid y) = \nabla_{x_t} \log p(x_t) + \nabla_{x_t} \log p(y \mid x_t). \quad (15)$$

As denoising score matching [60, 61] already provides a way to estimate the prior score $\nabla_{x_t} \log p(x_t)$, the remaining task is to estimate the likelihood score $\nabla_{x_t} \log p(y \mid x_t)$.

Diffusion posterior sampling Assuming a differentiable measurement function \mathcal{A} and a Gaussian forward process $p(y \mid x) = \mathcal{N}(y \mid \mathcal{A}(x), \Sigma_y)$, Chung et al. [21] propose the approximation

$$p(y \mid x_t) = \int p(y \mid x) p(x \mid x_t) dx \approx \mathcal{N}(y \mid \mathcal{A}(\mathbb{E}[x \mid x_t]), \Sigma_y) \quad (16)$$

which allows to estimate the likelihood score $\nabla_{x_t} \log p(y \mid x_t)$ without training any other network than $d_\theta(x_t, t) \approx \mathbb{E}[x \mid x_t]$. The motivation behind Eq. (16) is that, when σ_t is small, assuming that $p(x \mid x_t)$ is narrowly concentrated around its mean $\mathbb{E}[x \mid x_t]$ is reasonable. However, this approximation is very poor when σ_t is not negligible. Consequently, DPS [21] is unstable, does not properly cover the support of the posterior $p(x \mid y)$ and often leads to samples x which are inconsistent with the observation y [22–25].

Moment matching To address these flaws, following studies [22–25] take the covariance $\mathbb{V}[x \mid x_t]$ into account when estimating the likelihood score $\nabla_{x_t} \log p(y \mid x_t)$. Specifically, they consider the Gaussian approximation

$$q(x \mid x_t) = \mathcal{N}(x \mid \mathbb{E}[x \mid x_t], \mathbb{V}[x \mid x_t]) \quad (17)$$

which is closest to $p(x \mid x_t)$ in Kullback-Leibler (KL) divergence [67]. Then, assuming a linear Gaussian forward process $p(y \mid x) = \mathcal{N}(y \mid Ax, \Sigma_y)$, we obtain [67]

$$q(y \mid x_t) = \int p(y \mid x) q(x \mid x_t) dx = \mathcal{N}(y \mid A\mathbb{E}[x \mid x_t], \Sigma_y + A\mathbb{V}[x \mid x_t]A^\top) \quad (18)$$

which allows to estimate the likelihood score $\nabla_{x_t} \log p(y \mid x_t)$ as

$$\nabla_{x_t} \log q(y \mid x_t) = \nabla_{x_t} \mathbb{E}[x \mid x_t]^\top A^\top (\Sigma_y + A\mathbb{V}[x \mid x_t]A^\top)^{-1} (y - A\mathbb{E}[x \mid x_t]) \quad (19)$$

under the assumption that $\mathbb{V}[x \mid x_t]$ is constant with respect to x_t [24, 25]. Even with simple heuristics for $\mathbb{V}[x \mid x_t]$, such as Σ_t [20] or $(\Sigma_t^{-1} + \Sigma_x^{-1})^{-1}$ [22, 23], this adaptation leads to significantly more stable sampling and better coverage of the posterior $p(x \mid y)$ than DPS [21]. However, we find that heuristics lead to overly dispersed posteriors $q(x_t \mid y) \propto p(x_t) q(y \mid x_t)$ in the presence of local covariances – *i.e.* covariances in the neighborhood of x_t .

We illustrate this behavior in Figure 1 and measure its impact as the Sinkhorn divergence [68, 69] between the posteriors $p(x_t \mid y)$ and $q(x_t \mid y)$ when the prior $p(x)$ lies on randomly generated 1-dimensional manifolds [70] embedded in \mathbb{R}^3 . The prior $p(x)$ is modeled as a mixture of isotropic Gaussians centered around points of the manifold, which gives access to $p(x_t)$, $\mathbb{E}[x \mid x_t]$ and $\mathbb{V}[x \mid x_t]$ analytically. The results, presented in Figure 2, indicate that using $\mathbb{V}[x \mid x_t]$ instead of heuristics leads to orders of magnitude more accurate posteriors $q(x_t \mid y)$. We expect this gap to further increase with real high-dimensional data as the latter often lies along low-dimensional manifolds and, therefore, presents strong local covariances.

Because the MCEM algorithm is sensitive to the accuracy of posterior samples [48–51], we choose to estimate $\mathbb{V}[x \mid x_t]$ using Tweedie’s covariance formula [62–65]

$$\mathbb{V}[x \mid x_t] = \Sigma_t + \Sigma_t \nabla_{x_t}^2 \log p(x_t) \Sigma_t \quad (20)$$

$$= \Sigma_t \nabla_{x_t}^\top \mathbb{E}[x \mid x_t] \approx \Sigma_t \nabla_{x_t}^\top d_\theta(x_t, t). \quad (21)$$

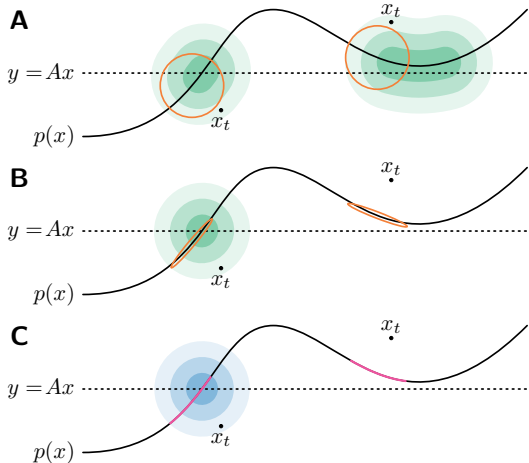


Figure 1. Illustration of the posterior $q(x_t | y)$ for the Gaussian approximation $q(x | x_t)$ when the prior $p(x)$ lies on a manifold. Ellipses represent 95% credible regions of $q(x | x_t)$. **(A)** With Σ_t as heuristic for $\mathbb{V}[x | x_t]$, any x_t whose mean $\mathbb{E}[x | x_t]$ is close to the plane $y = Ax$ is considered likely. **(B)** With $\mathbb{V}[x | x_t]$, more regions are correctly pruned. **(C)** Ground-truth $p(x_t | y)$ and $p(x | x_t)$ for reference.

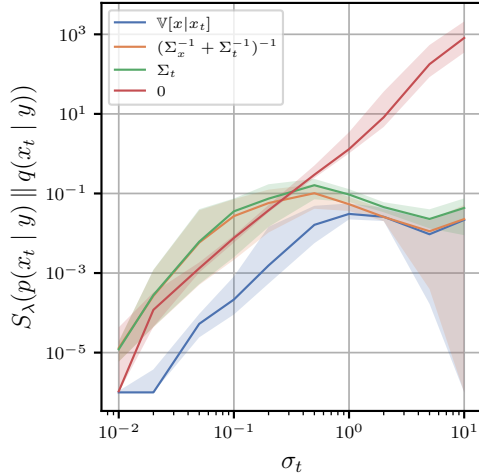


Figure 2. Sinkhorn divergence [68] between the posteriors $p(x_t | y)$ and $q(x_t | y)$ for different heuristics of $\mathbb{V}[x | x_t]$ when the prior $p(x)$ lies on 1-d manifolds embedded in \mathbb{R}^3 . Lines and shades represent the 25-50-75 percentiles for 64 randomly generated manifolds [70] and measurement matrices $A \in \mathbb{R}^{1 \times 3}$. Using $\mathbb{V}[x | x_t]$ instead of heuristics leads to orders of magnitude more accurate posteriors $q(x_t | y)$.

Conjugate gradient method As noted by Finzi et al. [24], explicitly computing and materializing the Jacobian $\nabla_{x_t}^\top d_\theta(x_t, t) \in \mathbb{R}^{N \times N}$ is intractable in high dimension. Furthermore, even if we had access to $\mathbb{V}[x | x_t]$, naively computing the inverse of the matrix $\Sigma_y + A\mathbb{V}[x | x_t]A^\top$ in Eq. (19) would still be intractable. Fortunately, we observe that the matrix $\Sigma_y + A\mathbb{V}[x | x_t]A^\top$ is symmetric positive definite (SPD) and, therefore, compatible with the conjugate gradient (CG) method [71]. The CG method is an iterative algorithm to solve linear systems of the form $Mv = b$ where the SPD matrix M and the vector b are known. Importantly, the CG method only requires implicit access to M through an operator that performs the matrix-vector product Mv given a vector v . In our case, the linear system to solve is

$$y - A\mathbb{E}[x | x_t] = (\Sigma_y + A\mathbb{V}[x | x_t]A^\top)v \quad (22)$$

$$\approx \Sigma_y v + A \underbrace{\Sigma_t \left(v^\top A \nabla_{x_t}^\top d_\theta(x_t, t) \right)^\top}_{\text{vector-Jacobian product}}. \quad (23)$$

Within automatic differentiation frameworks [72, 73], the vector-Jacobian product in the right-hand side can be cheaply evaluated. In practice, due to numerical errors and imperfect training, the Jacobian $\nabla_{x_t}^\top d_\theta(x_t, t)$ is not always perfectly SPD. Consequently, the CG method becomes unstable after a large number of iterations and fails to reach an exact solution. Fortunately, we find that using very few iterations (1 to 3) of the CG method as part of the computation of Eq. (19) already leads to significant improvements over using heuristics for the covariance $\mathbb{V}[x | x_t]$.

5 Results

We conduct three experiments to demonstrate the effectiveness of our method. We design the first experiment around a low-dimensional latent variable x whose ground-truth distribution $p(x)$ is known. In this setting, we can use asymptotically exact sampling schemes such as predictor-corrector sampling [23, 53] or twisted diffusion sampling [74] without worrying about their computational cost. This allows us to validate our expectation-maximization pipeline (see Algorithm 1) in the limit of (almost) exact posterior sampling. The remaining experiments target two benchmarks from previous studies: corrupted CIFAR-10 and accelerated MRI. These tasks provide a good understanding of how our method would perform in typical empirical Bayes applications with limited data and compute.

5.1 Low-dimensional manifold

In this experiment, the latent variable $x \in \mathbb{R}^5 \sim p(x)$ lies on a random 1-dimensional manifold embedded in \mathbb{R}^5 represented in Figure 7. Each observation $y \in \mathbb{R}^2 \sim \mathcal{N}(y | Ax, \Sigma_y)$ is the result of a random linear projection of a latent x plus isotropic Gaussian noise ($\Sigma_y = 10^{-4}I$). The rows of the measurement matrix $A \in \mathbb{R}^{2 \times 5}$ are drawn uniformly from the unit sphere \mathbb{S}^4 . We note that observing all push-forward distributions $p(u^\top x)$ with $u \in \mathbb{S}^{N-1}$ of a distribution $p(x)$ in \mathbb{R}^N is sufficient to recover $p(x)$ in theory [75, 76]. In practice, we generate a finite training set of 2^{16} pairs (y, A) .

We train a DM $q_\theta(x)$ parameterized by a multi-layer perceptron $d_\theta(x_t, t)$ for $K = 32$ EM iterations. We apply Algorithm 3 to estimate the posterior score $\nabla_{x_t} \log q_\theta(x_t | y, A)$, but rely on the predictor-corrector [23, 53] sampling scheme with a large number (4096) of correction steps to sample from the posterior $q_\theta(x | y, A)$. We provide additional details such as noise schedule, network architectures, learning rate, etc. in Appendix C.

As expected, the model $q_{\theta_k}(x)$ converges towards a stationary distribution whose marginals are close to the marginals of the ground-truth $p(x)$, as illustrated in Figure 3. We attribute the remaining artifacts to finite data and inaccuracies in our sampling scheme.

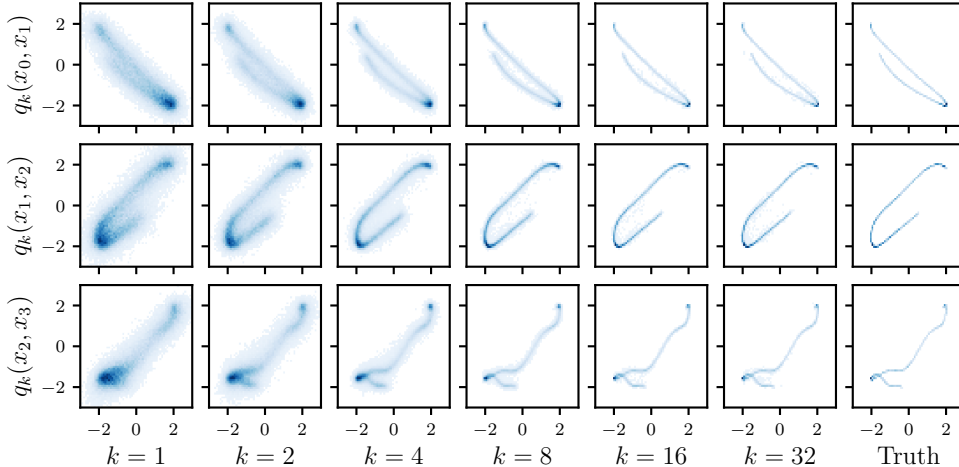


Figure 3. Illustration of 2-d marginals of the model $q_{\theta_k}(x)$ along the EM iterations. The Gaussian initial prior $q_0(x)$ leads to a very dispersed first model $q_{\theta_1}(x)$. The EM algorithm gradually prunes the density regions which are inconsistent with observations, until it reaches a stationary distribution. The marginals of the final distribution are close to the marginals of the ground-truth distribution.

5.2 Corrupted CIFAR-10

Following Daras et al. [77], we take the 50 000 training images of the CIFAR-10 [78] dataset as latent variables x . A single observation y is generated for each image x by randomly deleting pixels with 75 % probability. The measurement matrix A is therefore a binary diagonal matrix. We add negligible isotropic Gaussian noise ($\Sigma_y = 10^{-6}I$) for fair comparison with AmbientDiffusion [77] which cannot handle noisy observations.

We train a DM $q_\theta(x)$ parameterized by a U-Net [79] inspired network $d_\theta(x_t, t)$ for $K = 32$ EM iterations. We apply Algorithm 2 with $T = 256$ discretization steps to approximately sample from the posterior $q_\theta(x | y, A)$. We apply Algorithm 3 with several heuristics for $\mathbb{V}[x | x_t]$ to compare their results against Tweedie’s covariance formula. For the latter, we truncate the conjugate gradient method in Algorithm 4 to a single iteration.

For each model $q_{\theta_k}(x)$, we generate a set of 50 000 images and evaluate its Inception score (IS) [80] and Fréchet Inception distance (FID) [81] against the uncorrupted training set of CIFAR-10. We report the results in Table 1 and Figures 4 and 5. At 75 % of corruption, our method performs similarly to AmbientDiffusion [77] at only 40 % of corruption. On the contrary, using heuristics for $\mathbb{V}[x | x_t]$ leads to poor sample quality.

Method	Deleted	FID ↓	IS ↑
AmbientDiffusion [77]	0.4	18.85	7.45
	0.6	28.88	6.88
	0.8	46.27	6.14
Ours w/ Tweedie	0.75	19.56	7.60
Ours w/ $(I + \Sigma_t^{-1})^{-1}$	0.75	49.45	6.99
Ours w/ Σ_t	0.75	125.58	3.98

Table 1. Evaluation of final models trained on corrupted CIFAR-10. Our method outperforms AmbientDiffusion [77] at similar corruption level. Using heuristics for $\mathbb{V}[x | x_t]$ instead of Tweedie’s formula greatly decreases the sample quality.

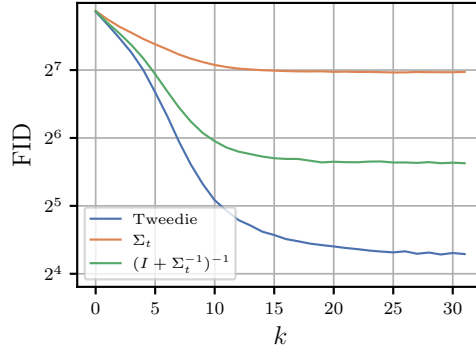


Figure 4. FID of $q_{\theta_k}(x)$ along the EM iterations for the corrupted CIFAR-10 experiment.



Figure 5. Example of samples from the model $q_{\theta_k}(x)$ along the EM iterations for the corrupted CIFAR-10 experiment. We use the deterministic DDIM [57] sampling scheme for easier comparison. Generated images become gradually more detailed and less noisy.

5.3 Accelerated MRI

Magnetic resonance imaging (MRI) is a non-invasive medical imaging technique used in radiology to inspect the internal anatomy and physiology of the body. MRI measurements of an object are obtained in the frequency domain, also called k -space, using strong magnetic fields. However, measuring the entire k -space can be time-consuming and expensive. Accelerated MRI [4–8] consists in inferring the scanned object based on partial, possibly randomized and noisy, frequency measurements.

In this experiment, following Kawar et al. [82], we consider the single-coil knee MRI scans from the fastMRI [7, 8] dataset. We treat each slice between the 10th and 40th of each scan as an independent latent variable x , represented as a 320×320 gray-scale image. Scans are min-max normalized such that pixel values range between -2 and 2 . A single observation y is generated for each slice x by first applying the discrete Fourier transform and then a random horizontal frequency sub-sampling with acceleration factor $R = 6$ [82, 83], meaning that a proportion $1/R$ of all frequencies are observed on average. Eventually, we obtain 24 853 k -space observations to which we add isotropic Gaussian noise ($\Sigma_y = 10^{-4}I$) to match Kawar et al. [82].

Once again, we train a DM $q_{\theta}(x)$ parameterized by a U-Net [79] inspired network $d_{\theta}(x_t, t)$ for $K = 16$ EM iterations. We apply Algorithm 2 with $T = 64$ discretization steps to approximately sample from the posterior $q_{\theta}(x | y, A)$ and truncate the conjugate gradient method in Algorithm 4 to 3 iterations. After training, we employ our final model $q_{\theta_{16}}(x)$ as a diffusion prior for accelerated MRI. Thanks to our moment matching posterior sampling, we are able to infer plausible scans at acceleration factors up to $R = 32$, as shown in Figure 6. Our samples are noticeably more detailed than the ones of Kawar et al. [82]. We choose not to report the PSNR/SSIM of our samples as these metrics only make sense for regression tasks and unfairly penalize proper generative models [84, 85]. We provide prior samples in Figure 13 and posterior samples for another kind of forward process in Figure 14.

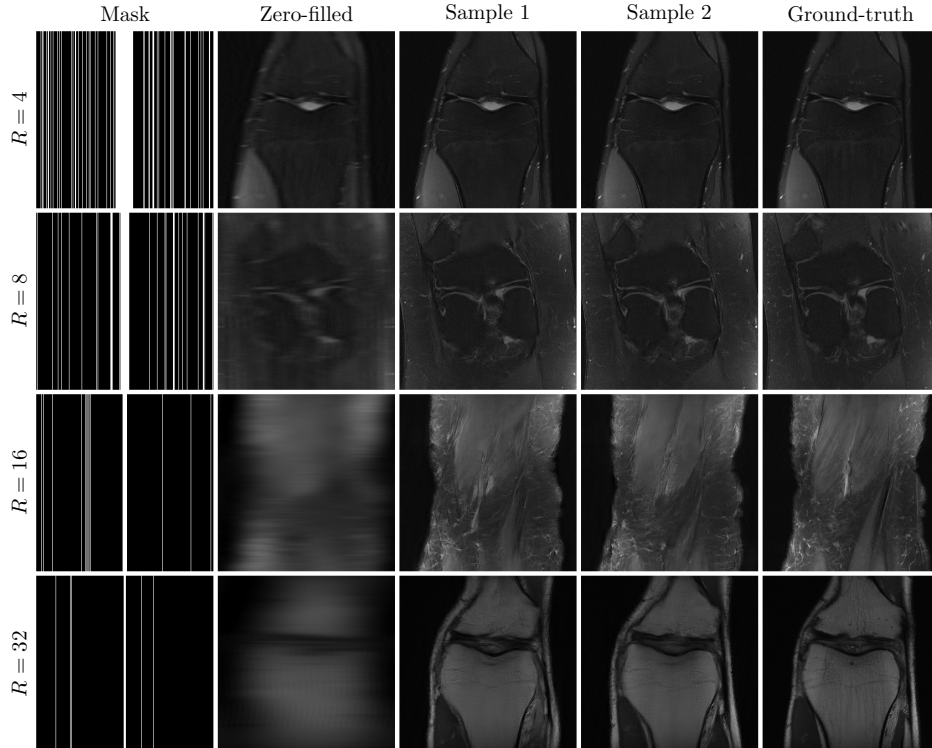


Figure 6. Examples of posterior samples for accelerated MRI using a diffusion prior trained from k -space observations only. Posterior samples are detailed and present plausible variations, while remaining consistent with the observation. We provide the zero-filled inverse, where missing frequencies are set to zero, as baseline.

6 Related Work

Empirical Bayes A number of previous studies have investigated the use of deep learning to solve the empirical Bayes problem. Louppe et al. [34] use adversarial training for learning a prior distribution that reproduces the empirical distribution of observations when pushed through a non-differentiable black-box forward process. Dockhorn et al. [32] use normalizing flows [86, 87] to estimate the prior density by variational inference when the forward process consists of additive noise. Vandegar et al. [35] also use normalizing flows but consider black-box forward processes for which the likelihood $p(y | x)$ is intractable. They note that empirical Bayes is an ill-posed problem in that distinct prior distributions may result in the same distribution over observations. Vetter et al. [36] address this issue by targeting the prior distribution of maximum entropy while minimizing the sliced-Wasserstein distance [75, 76] with the empirical distribution of observations. These methods rely on generative models $q_\theta(x)$ for which the density $q_\theta(x)$ or samples $x \sim q_\theta(x)$ are differentiable with respect to the parameters θ , which is not or hardly the case for diffusion models.

Closer to this work, Daras et al. [77] and Kawar et al. [82] attempt to train DMs from linear observations only. Daras et al. [77] consider noiseless observations of the form $y = Ax$ and train a network $d_\theta(Ax_t, A, t)$ to approximate $\mathbb{E}[x | Ax_t]$ under the assumption that $\mathbb{E}[A^\top A]$ is full-rank. The authors argue that $\mathbb{E}[x | Ax_t]$ can act as a surrogate for $\mathbb{E}[x | x_t]$. Similarly, Kawar et al. [82] assume Gaussian observations $y \sim \mathcal{N}(y | Ax, \Sigma_y)$ and train a network $d_\theta(Px_t, t)$ to approximate $\mathbb{E}[x | Px_t]$ under the assumption that $\mathbb{E}[P]$ is full-rank where $P = A^+A$ and A^+ is the Moore-Penrose pseudo-inverse of A . The authors assume that $d_\theta(Px_t, t)$ can generalize to $P = I$, even if the training data does not contain $P = I$. In both cases, the trained network is not a proper denoiser approximating $\mathbb{E}[x | x_t]$ and, therefore, cannot reliably parameterize a standard diffusion model. This is problematic for downstream applications such as Bayesian inference. Conversely, in this work, we do not make modifications to the denoising score matching objective [60, 61], which guarantees a proper DM at every iteration. In addition, we find that our method leads to quantitatively and qualitatively better diffusion priors.

Posterior sampling Recently, there has been much work on conditional generation using unconditional diffusion models, most of which adopt the posterior score decomposition in Eq. (15). As covered in Section 4.2, Chung et al. [21] propose an analytical approximation for the likelihood score $\nabla_{x_t} \log p(y | x_t)$ when the forward process $p(y | x)$ is Gaussian. Song et al. [22] and Rozet et al. [23] improve this approximation by taking the covariance $\mathbb{V}[x | x_t]$ into account in the form of simple heuristics. We build upon this idea and replace heuristics with a proper estimate of the covariance $\mathbb{V}[x | x_t]$ based on Tweedie’s covariance formula [62–65]. Finzi et al. [24] take the same approach, but materialize the matrix $A\mathbb{V}[x | x_t]A^\top$ which is intractable in high dimension. Boys et al. [25] replace the term $A\mathbb{V}[x | x_t]A^\top$ with a row-sum approximation $\text{diag}(Ae^\top \mathbb{V}[x | x_t]A^\top)$ where e is the all-ones vector. This approximation is only valid when A and $\mathbb{V}[x | x_t]$ are diagonal, which greatly limits its applicability. Instead, we take advantage of the conjugate gradient method [71] to avoid materializing $A\mathbb{V}[x | x_t]A^\top$.

A parallel line of work [88–90] extends the conditioning of diffusion models to arbitrary loss terms $\ell(x, y) \propto -\log p(y | x)$, for which no reliable analytical approximation of the likelihood score $\nabla_{x_t} \log p(y | x_t)$ exists. Instead, Song et al. [88] rely on Monte Carlo approximations of the likelihood $p(y | x_t) = \int p(y | x)p(x | x_t) dx$ by sampling from a Gaussian approximation of $p(x | x_t)$. Conversely, He et al. [90] use the mean $\mathbb{E}[x | x_t]$ as a point estimate for $p(x | x_t)$, but leverage a pre-trained encoder-decoder pair to project the updates of x_t within its manifold. We note that our use of the covariance $\mathbb{V}[x | x_t]$ similarly leads to updates tangent to the manifold of x_t .

Finally, Wu et al. [74] propose a particle-based posterior sampling scheme that is asymptotically exact in the limit of infinitely many particles as long as the likelihood approximation $q(y | x_t)$ – here named the *twisting* function – converges to $p(y | x)$ as t approaches 0. Using TDS [74] as part of our expectation-maximization pipeline could lead to better results and/or faster convergence, at the cost of computational resources. In addition, the authors note that the efficiency of TDS [74] depends on how closely the twisting function approximates the exact likelihood. In this regard, our moment matching Gaussian approximation in Eq. (18) could be a good twisting candidate.

7 Discussion

To the best of our knowledge, we are the first to formalize the training of diffusion models from corrupted observations as an empirical Bayes [27–30] problem. In this work, we limit our analysis to linear Gaussian forward processes to take advantage of accurate and efficient high-dimensional posterior sampling schemes. This contrasts with typical empirical Bayes methods which target low-dimensional latent spaces and highly non-linear forward processes [32–36]. As such, we choose to benchmark our work against similar methods in the diffusion model literature, but stress that a proper comparison with previous empirical Bayes methods would be valuable for both communities. We also note that Monte Carlo EM [44–51] can handle arbitrary forward processes $p(y | x)$ as long as one can sample from the posterior $q_\theta(x | y)$. Therefore, our method could be adapted to any kind of forward processes in the future. We believe that the works of Dhariwal et al. [91] and Ho et al. [92] on diffusion guidance are good avenues for adapting our method to non-linear, non-differentiable, or even black-box forward processes.

From a computational perspective, the iterative nature of our method is a drawback compared to previous works [77, 82]. Notably, generating enough samples from the posterior can be expensive, although embarrassingly parallelizable. In addition, even though the architecture and training of the model $q_\theta(x)$ itself are simpler than in previous works [77, 82], the sampling step adds a significant amount of complexity, especially as the convergence of our method is sensitive to the quality of posterior samples. Furthermore, the use of a vector-Jacobian product in the linear system solved as part of Eq. (19) significantly increases the cost of each posterior score evaluation compared to heuristics. This increase scales linearly with the number of conjugate gradient iterations. However, as mentioned in Section 4.2, we find that truncating the CG method to very few iterations (1 to 3) already leads to significant improvements over heuristics for the covariance $\mathbb{V}[x | x_t]$. One of these improvements is the ability to greatly reduce the number of discretization steps without apparent loss of quality, which mitigates the higher per-step cost of moment matching posterior sampling.

Finally, as mentioned in Section 6, empirical Bayes is an ill-posed problem in that distinct prior distributions may result in the same distribution over observations. Following the maximum entropy principle, as advocated by Vetter et al. [36], is left to future work.

Acknowledgments and Disclosure of Funding

François Rozet and G r me Andry are research fellows of the F.R.S.-FNRS (Belgium) and acknowledge its financial support.

The present research benefited from computational resources made available on Lucia, the Tier-1 supercomputer of the Walloon Region, infrastructure funded by the Walloon Region under the grant n 1910247. The computational resources have been provided by the Consortium des  quipements de Calcul Intensif (C ICI), funded by the Fonds de la Recherche Scientifique de Belgique (F.R.S.-FNRS) under the grant n 2.5020.11 and by the Walloon Region.

MRI data used in the preparation of this article were obtained from the NYU fastMRI Initiative database [7, 8]. As such, NYU fastMRI investigators provided data but did not participate in analysis or writing of this report. A listing of NYU fastMRI investigators, subject to updates, can be found at <https://fastmri.med.nyu.edu/>. The primary goal of fastMRI is to test whether machine learning can aid in the reconstruction of medical images.

References

- [1] S. J. Warren and S. Dye. “Semilinear Gravitational Lens Inversion”. In *The Astrophysical Journal* (2003).
- [2] Warren R. Morningstar et al. “Data-driven Reconstruction of Gravitationally Lensed Galaxies Using Recurrent Inference Machines”. In *The Astrophysical Journal* (2019).
- [3] Siddharth Mishra-Sharma and Ge Yang. “Strong Lensing Source Reconstruction Using Continuous Neural Fields”. 2022.
- [4] Shanshan Wang et al. “Accelerating magnetic resonance imaging via deep learning”. In *International Symposium on Biomedical Imaging*. 2016.
- [5] Kerstin Hammernik et al. “Learning a variational network for reconstruction of accelerated MRI data”. In *Magnetic Resonance in Medicine* (2018).
- [6] Yoseo Han et al. “k-Space Deep Learning for Accelerated MRI”. In *Transactions on Medical Imaging* (2020).
- [7] Jure Zbontar et al. “fastMRI: An Open Dataset and Benchmarks for Accelerated MRI”. 2018.
- [8] Florian Knoll et al. “fastMRI: A Publicly Available Raw k-Space and DICOM Dataset of Knee Images for Accelerated MR Image Reconstruction Using Machine Learning”. In *Radiology: Artificial Intelligence* (2020).
- [9] G. Cowan. “A survey of unfolding methods for particle physics”. In *Conf. Proc. C* (2002).
- [10] Volker Blobel. “Unfolding Methods in Particle Physics”. In *PHYSTAT*. CERN, 2011.
- [11] Fran ois-Xavier Le Dimet and Olivier Talagrand. “Variational algorithms for analysis and assimilation of meteorological observations: theoretical aspects”. In *Tellus A: Dynamic Meteorology and Oceanography* (1986).
- [12] Yannick Tr molet. “Accounting for an imperfect model in 4D-Var”. In *Quarterly Journal of the Royal Meteorological Society* (2006).
- [13] Thomas M. Hamill. “Ensemble-based atmospheric data assimilation”. In *Predictability of Weather and Climate*. 2006.
- [14] Alberto Carrassi et al. “Data assimilation in the geosciences: An overview of methods, issues, and perspectives”. In *WIREs Climate Change* (2018).
- [15] Jascha Sohl-Dickstein et al. “Deep Unsupervised Learning using Nonequilibrium Thermodynamics”. In *Proceedings of the 32nd International Conference on Machine Learning*. 2015.
- [16] Jonathan Ho et al. “Denoising Diffusion Probabilistic Models”. In *Advances in Neural Information Processing Systems*. 2020.
- [17] Yang Song et al. “Solving Inverse Problems in Medical Imaging with Score-Based Generative Models”. In *International Conference on Learning Representations*. 2022.
- [18] Bahjat Kawar et al. “SNIPS: Solving Noisy Inverse Problems Stochastically”. In *Advances in Neural Information Processing Systems*. 2021.

- [19] Bahjat Kawar et al. “Denoising Diffusion Restoration Models”. In *Advances in Neural Information Processing Systems*. 2022.
- [20] Alexandre Adam et al. “Posterior samples of source galaxies in strong gravitational lenses with score-based priors”. 2022.
- [21] Hyungjin Chung et al. “Diffusion Posterior Sampling for General Noisy Inverse Problems”. In *International Conference on Learning Representations*. 2023.
- [22] Jiaming Song et al. “Pseudoinverse-Guided Diffusion Models for Inverse Problems”. In *International Conference on Learning Representations*. 2023.
- [23] François Rozet and Gilles Louppe. “Score-based Data Assimilation”. In *Advances in Neural Information Processing Systems*. 2023.
- [24] Marc Anton Finzi et al. “User-defined Event Sampling and Uncertainty Quantification in Diffusion Models for Physical Dynamical Systems”. In *Proceedings of the 40th International Conference on Machine Learning*. 2023.
- [25] Benjamin Boys et al. “Tweedie Moment Projected Diffusions For Inverse Problems”. 2023.
- [26] Noe Dia et al. “Bayesian Imaging for Radio Interferometry with Score-based Priors”. 2023.
- [27] Herbert E. Robbins. “An Empirical Bayes Approach to Statistics”. In *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability*. 1956.
- [28] George Casella. “An Introduction to Empirical Bayes Data Analysis”. In *The American Statistician* (1985).
- [29] Bradley P. Carlin and Thomas A. Louis. “Empirical Bayes: Past, Present and Future”. In *Journal of the American Statistical Association* (2000).
- [30] Bradley Efron. “Two Modeling Strategies for Empirical Bayes Estimation”. In *Statistical Science* (2014).
- [31] G. D’Agostini. “A multidimensional unfolding method based on Bayes’ theorem”. In *Nuclear Instruments and Methods in Physics Research* (1995).
- [32] Tim Dockhorn et al. “Density Deconvolution with Normalizing Flows”. 2020.
- [33] Anders Andreassen et al. “OmniFold: A Method to Simultaneously Unfold All Observables”. In *Physical Review Letters* (2020).
- [34] Gilles Louppe et al. “Adversarial Variational Optimization of Non-Differentiable Simulators”. In *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*. 2019.
- [35] Maxime Vandegar et al. “Neural Empirical Bayes: Source Distribution Estimation and its Applications to Simulation-Based Inference”. In *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*. 2021.
- [36] Julius Vetter et al. “Sourcerer: Sample-based Maximum Entropy Source Distribution Estimation”. 2024.
- [37] Bradley Efron. “Empirical Bayes deconvolution estimates”. In *Biometrika* (2016).
- [38] Balasubramanian Narasimhan and Bradley Efron. “deconvolveR: A G-Modeling Program for Deconvolution and Empirical Bayes Estimation”. In *Journal of Statistical Software* (2020).
- [39] H. O. Hartley. “Maximum Likelihood Estimation from Incomplete Data”. In *Biometrics* (1958).
- [40] A. P. Dempster et al. “Maximum Likelihood from Incomplete Data Via the EM Algorithm”. In *Journal of the Royal Statistical Society* (1977).
- [41] C. F. Jeff Wu. “On the Convergence Properties of the EM Algorithm”. In *The Annals of Statistics* (1983).
- [42] Geoffrey J McLachlan and Thriyambakam Krishnan. “The EM algorithm and extensions”. John Wiley & Sons, 2007.
- [43] Sivaraman Balakrishnan et al. “Statistical guarantees for the EM algorithm: From population to sample-based analysis”. In *The Annals of Statistics* (2017).
- [44] Greg C. G. Wei and Martin A. Tanner. “A Monte Carlo Implementation of the EM Algorithm and the Poor Man’s Data Augmentation Algorithms”. In *Journal of the American Statistical Association* (1990).
- [45] Gilles Celeux and Jean Diebolt. “A stochastic approximation type EM algorithm for the mixture problem”. In *Stochastics and Stochastic Reports* (1992).

- [46] Bernard Delyon et al. “Convergence of a stochastic approximation version of the EM algorithm”. In *The Annals of Statistics* (1999).
- [47] James G. Booth and James P. Hobert. “Maximizing Generalized Linear Mixed Model Likelihoods with an Automated Monte Carlo EM Algorithm”. In *Journal of the Royal Statistical Society* (1999).
- [48] Richard A. Levine and George Casella. “Implementations of the Monte Carlo EM Algorithm”. In *Journal of Computational and Graphical Statistics* (2001).
- [49] Brian S. Caffo et al. “Ascent-Based Monte Carlo Expectation-Maximization”. In *Journal of the Royal Statistical Society* (2005).
- [50] Wolfgang Jank. “The EM Algorithm, Its Randomized Implementation and Global Optimization”. In *Perspectives in Operations Research*. 2006.
- [51] William Ruth. “A review of Monte Carlo-based versions of the EM algorithm”. 2024.
- [52] Yang Song and Stefano Ermon. “Generative Modeling by Estimating Gradients of the Data Distribution”. In *Advances in Neural Information Processing Systems*. 2019.
- [53] Yang Song et al. “Score-Based Generative Modeling through Stochastic Differential Equations”. In *International Conference on Learning Representations*. 2021.
- [54] Brian D. O. Anderson. “Reverse-time diffusion equation models”. In *Stochastic Processes and their Applications* (1982).
- [55] Simo Särkkä and Arno Solin. “Applied Stochastic Differential Equations”. Institute of Mathematical Statistics Textbooks. Cambridge University Press, 2019.
- [56] Qinsheng Zhang and Yongxin Chen. “Fast Sampling of Diffusion Models with Exponential Integrator”. In *International Conference on Learning Representations*. 2023.
- [57] Jiaming Song et al. “Denoising Diffusion Implicit Models”. In *International Conference on Learning Representations*. 2021.
- [58] Tero Karras et al. “Elucidating the Design Space of Diffusion-Based Generative Models”. In *Advances in Neural Information Processing Systems*. 2022.
- [59] Yaron Lipman et al. “Flow Matching for Generative Modeling”. In *International Conference on Learning Representations*. 2023.
- [60] Aapo Hyvärinen. “Estimation of Non-Normalized Statistical Models by Score Matching”. In *Journal of Machine Learning Research* (2005).
- [61] Pascal Vincent. “A Connection Between Score Matching and Denoising Autoencoders”. In *Neural Computation* (2011).
- [62] M. C. K. Tweedie. “Functions of a statistical variate with given means, with special reference to Laplacian distributions”. In *Mathematical Proceedings of the Cambridge Philosophical Society* (1947).
- [63] Bradley Efron. “Tweedie’s Formula and Selection Bias”. In *Journal of the American Statistical Association* (2011).
- [64] Kwanyoung Kim and Jong Chul Ye. “Noise2Score: Tweedie’s Approach to Self-Supervised Image Denoising without Clean Images”. In *Advances in Neural Information Processing Systems*. 2021.
- [65] Chenlin Meng et al. “Estimating High Order Gradients of the Data Distribution by Denoising”. In *Advances in Neural Information Processing Systems*. 2021.
- [66] Michael E. Tipping and Christopher M. Bishop. “Mixtures of Probabilistic Principal Component Analyzers”. In *Neural Computation* (1999).
- [67] Christopher M. Bishop. “Pattern Recognition and Machine Learning”. Information Science and Statistics. Springer, 2006.
- [68] Lénaïc Chizat et al. “Faster Wasserstein Distance Estimation with the Sinkhorn Divergence”. In *Advances in Neural Information Processing Systems*. 2020.
- [69] Rémi Flamary et al. “POT: Python Optimal Transport”. In *Journal of Machine Learning Research* (2021).
- [70] Friedemann Zenke and Tim P. Vogels. “The Remarkable Robustness of Surrogate Gradient Learning for Instilling Complex Function in Spiking Neural Networks”. In *Neural Computation* (2021).

- [71] Magnus R. Hestenes and Eduard Stiefel. “Methods of Conjugate Gradients for Solving Linear Systems”. In *Journal of Research of the National Bureau of Standards* (1952).
- [72] James Bradbury et al. “JAX: Composable transformations of Python + NumPy programs”. 2018.
- [73] Adam Paszke et al. “PyTorch: An Imperative Style, High-Performance Deep Learning Library”. In *Advances in Neural Information Processing Systems*. 2019.
- [74] Luhuan Wu et al. “Practical and Asymptotically Exact Conditional Sampling in Diffusion Models”. In *Advances in Neural Information Processing Systems*. 2023.
- [75] Nicolas Bonneel et al. “Sliced and Radon Wasserstein Barycenters of Measures”. In *Journal of Mathematical Imaging and Vision* (2015).
- [76] Kimia Nadjahi et al. “Statistical and Topological Properties of Sliced Probability Divergences”. In *Advances in Neural Information Processing Systems*. 2020.
- [77] Giannis Daras et al. “Ambient Diffusion: Learning Clean Distributions from Corrupted Data”. In *Advances in Neural Information Processing Systems*. 2023.
- [78] Alex Krizhevsky, Geoffrey Hinton, et al. “Learning Multiple Layers of Features from Tiny Images”. 2009.
- [79] Olaf Ronneberger et al. “U-Net: Convolutional Networks for Biomedical Image Segmentation”. In *Medical Image Computing and Computer-Assisted Intervention*. 2015.
- [80] Tim Salimans et al. “Improved Techniques for Training GANs”. In *Advances in Neural Information Processing Systems*. 2016.
- [81] Martin Heusel et al. “GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium”. In *Advances in Neural Information Processing Systems*. 2017.
- [82] Bahjat Kawar et al. “GSURE-Based Diffusion Model Training with Corrupted Data”. In *Transactions on Machine Learning Research* (2024).
- [83] Ajil Jalal et al. “Robust Compressed Sensing MRI with Deep Generative Priors”. In *Advances in Neural Information Processing Systems*. 2021.
- [84] Yochai Blau and Tomer Michaeli. “The Perception-Distortion Tradeoff”. In *Conference on Computer Vision and Pattern Recognition*. 2018.
- [85] Mauricio Delbracio and Peyman Milanfar. “Inversion by Direct Iteration: An Alternative to Denoising Diffusion for Image Restoration”. In *Transactions on Machine Learning Research* (2023).
- [86] E. G. Tabak and Cristina V. Turner. “A family of nonparametric density estimation algorithms”. In *Communications on Pure and Applied Mathematics* (2013).
- [87] Danilo Rezende and Shakir Mohamed. “Variational Inference with Normalizing Flows”. In *Proceedings of the 32nd International Conference on Machine Learning*. 2015.
- [88] Jiaming Song et al. “Loss-Guided Diffusion Models for Plug-and-Play Controllable Generation”. In *Proceedings of the 40th International Conference on Machine Learning*. 2023.
- [89] Arpit Bansal et al. “Universal Guidance for Diffusion Models”. In *International Conference on Learning Representations*. 2024.
- [90] Yutong He et al. “Manifold Preserving Guided Diffusion”. In *International Conference on Learning Representations*. 2024.
- [91] Prafulla Dhariwal and Alexander Quinn Nichol. “Diffusion Models Beat GANs on Image Synthesis”. In *Advances in Neural Information Processing Systems*. 2021.
- [92] Jonathan Ho and Tim Salimans. “Classifier-Free Diffusion Guidance”. 2022.
- [93] Stefan Elfving et al. “Sigmoid-weighted linear units for neural network function approximation in reinforcement learning”. In *Neural Networks* (2018).
- [94] Jimmy Lei Ba et al. “Layer Normalization”. 2016.
- [95] Diederik P. Kingma and Jimmy Ba. “Adam: A Method for Stochastic Optimization”. In *International Conference on Learning Representations*. 2015.
- [96] Kaiming He et al. “Deep Residual Learning for Image Recognition”. In *Conference on Computer Vision and Pattern Recognition*. 2016.
- [97] William Peebles and Saining Xie. “Scalable Diffusion Models with Transformers”. In *International Conference on Computer Vision*. 2023.

- [98] Ashish Vaswani et al. "Attention is All you Need". In *Advances in Neural Information Processing Systems*. 2017.
- [99] Anton Obukhov et al. "High-fidelity performance metrics for generative models in PyTorch". 2020.
- [100] Wenzhe Shi et al. "Real-Time Single Image and Video Super-Resolution Using an Efficient Sub-Pixel Convolutional Neural Network". In *Conference on Computer Vision and Pattern Recognition*. 2016.

A Algorithms

Algorithm 1 Expectation-maximization pipeline

```

1 Choose an initial prior  $q_0(x)$ 
2 Initialize the parameters  $\theta$  of the denoiser  $d_\theta(x_t, t)$ 
3 for  $k = 1$  to  $K$  do
4   for  $i = 1$  to  $S$  do
5      $y_i, A_i \sim p(y, A)$ 
6      $x_i \sim q_k(x | y_i, A_i)$  # Posterior sampling
7   repeat
8      $i \sim \mathcal{U}(\{1, \dots, S\})$ 
9      $t \sim \mathcal{U}(0, 1)$ 
10     $z \sim \mathcal{N}(0, I)$ 
11     $x_t \leftarrow x_i + \sigma_t z$ 
12     $\ell \leftarrow \lambda_t \|d_\theta(x_t, t) - x_i\|^2$  # Denoising score matching
13     $\theta \leftarrow \text{GRADIENTDESCENT}(\theta, \nabla_\theta \ell)$ 
14  until convergence
15   $\theta_k \leftarrow \theta$ 
16   $q_k := q_{\theta_k}$ 
17 return  $\theta_K$ 

```

Algorithm 2 DDPM-style posterior sampling

```

1  $x_T \sim \mathcal{N}(0, \Sigma_T)$ 
2 for  $t = T$  to 1 do
3    $\hat{x} \leftarrow x_t + \Sigma_t s_\theta(x_t | y, A)$ 
4    $z \sim \mathcal{N}(0, I)$ 
5    $x_{t-1} \leftarrow \frac{\sigma_{t-1}^2}{\sigma_t^2} x_t + \left(1 - \frac{\sigma_{t-1}^2}{\sigma_t^2}\right) \hat{x} + \sqrt{\sigma_{t-1}^2 - \frac{\sigma_{t-1}^4}{\sigma_t^2}} z$ 
6 return  $x_0$ 

```

Algorithm 3 Moment matching posterior score

```

1 function  $s_\theta(x_t | y, A)$  # Estimate  $\nabla_{x_t} \log q_\theta(x_t | y, A)$ 
2    $\hat{x} \leftarrow d_\theta(x_t, t)$ 
3   if Tweedie then
4      $\Sigma_{x|x_t} \leftarrow \Sigma_t \nabla_{x_t} d_\theta(x_t, t)^\top$ 
5   else
6      $\Sigma_{x|x_t} \leftarrow \Sigma_t$  or  $(I + \Sigma_t^{-1})^{-1}$  or  $(\Sigma_x^{-1} + \Sigma_t^{-1})^{-1}$ 
7    $b \leftarrow y - A\hat{x}$ 
8    $u \leftarrow \text{CONJUGATEGRADIENT}(\Sigma_y + A\Sigma_{x|x_t}A^\top, b)$ 
9    $s_{y|x} \leftarrow \nabla_{x_t} d_\theta(x_t, t)^\top A^\top u$  # Estimate  $\nabla_{x_t} \log q_\theta(y | x_t, A)$ 
10   $s_x \leftarrow \Sigma_t^{-1}(\hat{x} - x_t)$  # Estimate  $\nabla_{x_t} \log q_\theta(x_t)$ 
11  return  $s_x + s_{y|x}$ 

```

Algorithm 4 Conjugate gradient method

```
1 function CONJUGATEGRADIENT( $A, b, x_0$ )
2    $r_0 \leftarrow b - Ax_0$ 
3    $p_0 \leftarrow r_0$ 
4   for  $i = 0$  to  $N - 1$  do
5     if  $\|r_i\| \leq \epsilon$  then
6       return  $x_i$ 
7      $\alpha_i \leftarrow \frac{r_i^\top r_i}{p_i^\top Ap_i}$ 
8      $x_{i+1} \leftarrow x_i + \alpha_i p_i$ 
9      $r_{i+1} \leftarrow r_i - \alpha_i Ap_i$ 
10     $\beta_i \leftarrow \frac{r_{i+1}^\top r_{i+1}}{r_i^\top r_i}$ 
11     $p_{i+1} \leftarrow r_{i+1} + \beta_i p_i$ 
12  return  $x_N$ 
```

B Tweedie's formulae

Theorem 1. For any distribution $p(x)$ and $p(x_t | x) = \mathcal{N}(x_t | x, \Sigma_t)$, the first and second moments of the distribution $p(x | x_t)$ are linked to the score function $\nabla_{x_t} \log p(x_t)$ through

$$\mathbb{E}[x | x_t] = x_t + \Sigma_t \nabla_{x_t} \log p(x_t) \quad (24)$$

$$\mathbb{V}[x | x_t] = \Sigma_t + \Sigma_t \nabla_{x_t}^2 \log p(x_t) \Sigma_t \quad (25)$$

Proof.

$$\begin{aligned} \nabla_{x_t} \log p(x_t) &= \frac{1}{p(x_t)} \nabla_{x_t} p(x_t) \\ &= \frac{1}{p(x_t)} \int \nabla_{x_t} p(x, x_t) dx \\ &= \frac{1}{p(x_t)} \int p(x, x_t) \nabla_{x_t} \log p(x, x_t) dx \\ &= \int p(x | x_t) \nabla_{x_t} \log p(x_t | x) dx \\ &= \int p(x | x_t) \Sigma_t^{-1} (x - x_t) dx \\ &= \Sigma_t^{-1} \mathbb{E}[x | x_t] - \Sigma_t^{-1} x_t \quad \square \end{aligned}$$

Proof.

$$\begin{aligned} \nabla_{x_t}^2 \log p(x_t) &= \nabla_{x_t} \nabla_{x_t}^\top \log p(x_t) \\ &= \nabla_{x_t} \mathbb{E}[x | x_t]^\top \Sigma_t^{-1} - \Sigma_t^{-1} \\ &= \left(\int \nabla_{x_t} p(x | x_t) x^\top dx \right) \Sigma_t^{-1} - \Sigma_t^{-1} \\ &= \left(\int p(x | x_t) \nabla_{x_t} \log \frac{p(x_t | x)}{p(x_t)} x^\top dx \right) \Sigma_t^{-1} - \Sigma_t^{-1} \\ &= \left(\int p(x | x_t) \Sigma_t^{-1} (x - \mathbb{E}[x | x_t]) x^\top dx \right) \Sigma_t^{-1} - \Sigma_t^{-1} \\ &= \Sigma_t^{-1} \left(\mathbb{E}[xx^\top | x_t] - \mathbb{E}[x | x_t] \mathbb{E}[x | x_t]^\top \right) \Sigma_t^{-1} - \Sigma_t^{-1} \\ &= \Sigma_t^{-1} \mathbb{V}[x | x_t] \Sigma_t^{-1} - \Sigma_t^{-1} \quad \square \end{aligned}$$

C Experiment details

All experiments are implemented within the JAX [72] automatic differentiation framework. The code for all experiments is made available at <https://github.com/anonymous/anonymous>.

Diffusion models As mentioned in Section 2, in this work, we adopt the variance exploding SDE [52] and the denoiser parameterization [58]. Following Karras et al. [58], we precondition our denoiser $d_\theta(x_t, t)$ as

$$d_\theta(x_t, t) = \frac{1}{\sigma_t^2 + 1} x_t + \frac{\sigma_t}{\sqrt{\sigma_t^2 + 1}} h_\theta(x_t, \sigma_t) \quad (26)$$

where $h_\theta(x_t, \sigma_t)$ is an arbitrary noise-conditional network. The scalar noise σ_t is embedded as a vector using a smooth one-hot encoding. In our experiments, we use an exponential noise schedule

$$\sigma_t = \exp((1 - t) \log 10^{-3} + t \log 10^2), \quad (27)$$

loss weights $\lambda_t = \frac{1}{\sigma_t^2} + 1$ and sample t from a Beta distribution $\mathcal{B}(\alpha = 3, \beta = 3)$ during training.

Low-dimensional manifold The noise-conditional network $h_\theta(x_t, \sigma_t)$ is a multi-layer perceptron with 3 hidden layers of 256 neurons and SiLU [93] activation functions. A layer normalization [94] function is inserted after each activation. The input of the network is the concatenation of x_t and the noise embedding vector. We train the network with Algorithm 1 for $K = 32$ EM iterations. Each iteration consists of 16 384 optimization steps of the Adam [95] optimizer. The optimizer and learning rate are re-initialized after each EM iteration. Other hyperparameters are provided in Table 2.

Table 2. Hyperparameters for the low-dimensional manifold experiment.

Architecture	MLP
Input shape	(5)
Hidden features	(256, 256, 256)
Activation	SiLU
Normalization	LayerNorm
Optimizer	Adam
Weight decay	0.0
Scheduler	linear
Initial learning rate	1×10^{-3}
Final learning rate	1×10^{-6}
Gradient norm clipping	1.0
Batch size	1024
Steps per EM iteration	16 384
EM iterations	32

We apply Algorithm 3 to estimate the posterior score $\nabla_{x_t} \log p(x_t)$ and truncate Algorithm 4 to 3 iterations. We rely on the predictor-corrector [23, 53] sampling scheme to sample from the posterior $q_\theta(x | y, A)$. Following Rozet et al. [23], the predictor is a deterministic DDIM [57] step and the corrector is a Langevin Monte Carlo step. We perform 4096 prediction steps, each followed by 1 correction step. At each EM iteration, we generate a single latent x for each pair (y, A) .

We generate smooth random manifolds according to a procedure described by Zenke et al. [70]. We evaluate the Sinkhorn divergences using the POT [69] package with an entropic regularization factor $\lambda = 1e - 3$.

Corrupted CIFAR-10 The noise-conditional network $h_\theta(x_t, \sigma_t)$ is a U-Net [79] with residual blocks [96], SiLU [93] activation functions and layer normalization [94]. Each residual block is modulated with respect to the noise σ_t in the style of diffusion transformers [97]. A multi-head self-attention block [98] is inserted after each residual block at the last level of the U-Net. We train the network with Algorithm 1 for $K = 32$ EM iterations. Each iteration consists of 256 epochs over the training set (50 000 images). To prevent overfitting, images are augmented with horizontal flips

and “pad & crop” (reflection padding + random crop). The optimizer is re-initialized after each EM iteration. Other hyperparameters are provided in Table 3.

Table 3. Hyperparameters for the corrupted CIFAR-10 and accelerated MRI experiments.

Experiment	corrupted CIFAR-10	accelerated MRI
Architecture	U-Net	U-Net
Input shape	(32, 32, 3)	(80, 80, 16)
Residual blocks per level	(5, 5, 5)	(3, 3, 3, 3)
Channels per level	(128, 256, 384)	(128, 256, 384, 512)
Attention heads per level	(0, 0, 4)	(0, 0, 0, 4)
Kernel size	3	3
Activation	SiLU	SiLU
Normalization	LayerNorm	LayerNorm
Optimizer	Adam	Adam
Weight decay	0.0	0.0
Learning rate	2×10^{-4}	10^{-4}
Gradient norm clipping	1.0	1.0
EMA decay	0.999	0.999
Dropout	0.1	0.1
Augmentation	h-flip, pad & crop	h-flip, pad & crop
Batch size	256	256
Epochs per EM iteration	256	64
EM iterations	32	16

We apply Algorithm 2 with $T = 256$ discretization steps to sample from the posterior $q_\theta(x | y, A)$. We apply Algorithm 3 with several heuristics for $\mathbb{V}[x | x_t]$ to compare their results against Tweedie’s covariance formula. For the latter, we truncate the conjugate gradient method in Algorithm 4 to a single iteration. At each EM iteration, we generate a single latent x for each pair (y, A) . Each EM iteration (including sampling and training) takes around 3 h on 4 A100 (40GB) GPUs.

We evaluate the Inception score (IS) [80] and Fréchet Inception distance (FID) [81] of generated images using the torch-fidelity [99] package.

Accelerated MRI The noise-conditional network architecture is the same as for the corrupted CIFAR-10 experiment. The $320 \times 320 \times 1$ tensor x_t is reshaped into a $80 \times 80 \times 16$ tensor using pixel shuffling [100] before entering the network. We train the network with Algorithm 1 for $K = 16$ EM iterations. Each iteration consists of 64 epochs over the training set (24 853 images). The optimizer is re-initialized after each EM iteration. Other hyperparameters are provided in Table 3.

We apply Algorithm 2 with $T = 64$ discretization steps to sample from the posterior $q_\theta(x | y, A)$. We truncate the conjugate gradient method in Algorithm 4 to 3 iterations. At each EM iteration, we generate a single latent x for each pair (y, A) . Each EM iteration (including sampling and training) takes around 3 h on 4 A100 (40GB) GPUs.

D Additional figures

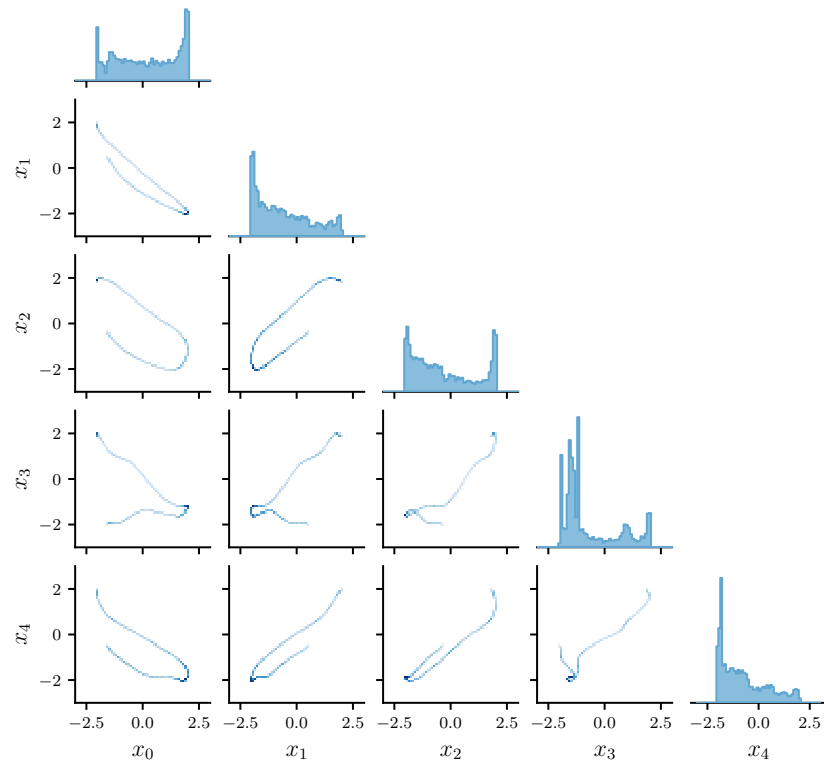


Figure 7. 1-d and 2-d marginals of the ground-truth distribution $p(x)$ used in the low-dimensional manifold experiment. The distribution lies on a random 1-dimensional manifold embedded in \mathbb{R}^5 .

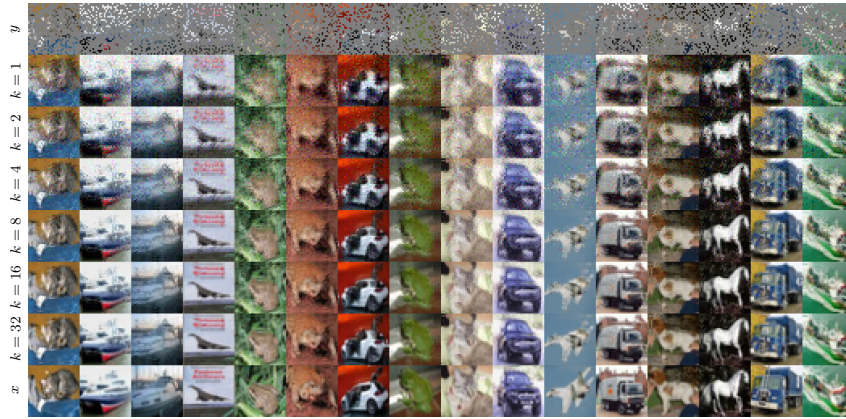


Figure 8. Example of samples from the posterior $q_{\theta_k}(x | y)$ along the EM iterations for the CIFAR-10 experiment. The generated images become gradually more detailed and less noisy.

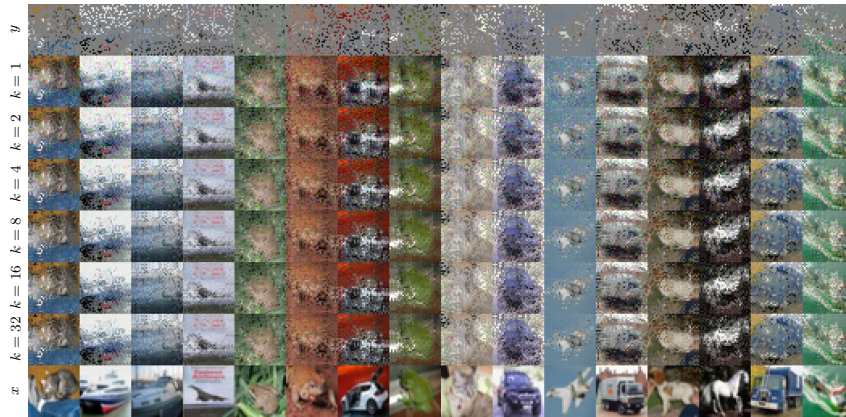


Figure 9. Example of samples from the posterior $q_{\theta_k}(x | y)$ along the EM iterations for the CIFAR-10 experiment when the heuristic Σ_t is used for $\mathbb{V}[x | x_t]$. The generated images remain very noisy.

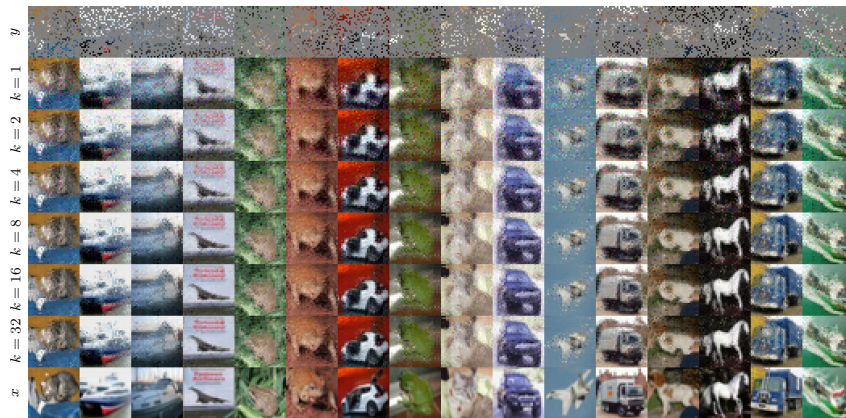


Figure 10. Example of samples from the posterior $q_{\theta_k}(x | y)$ along the EM iterations for the CIFAR-10 experiment when the heuristic $(I + \Sigma_t^{-1})^{-1}$ is used for $\mathbb{V}[x | x_t]$. The generated images become gradually more detailed but some noise remains.

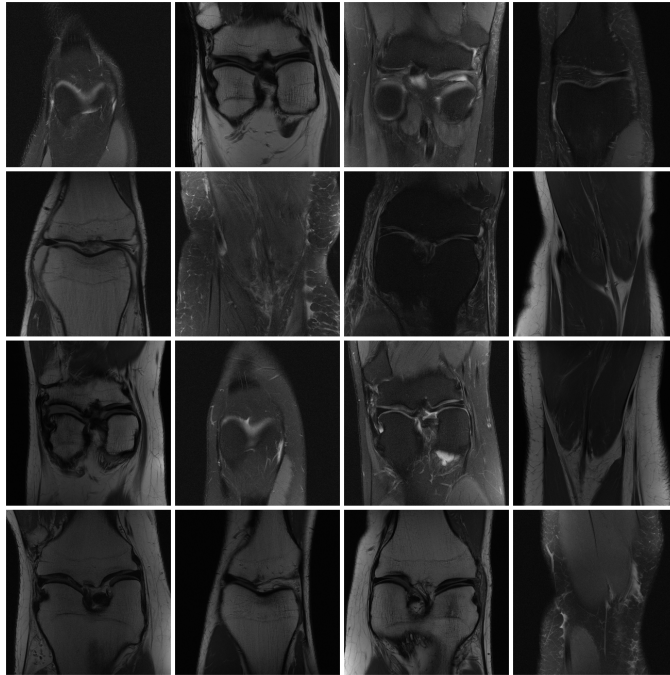


Figure 11. Example of scan slices from the fastMRI [7, 8] dataset.

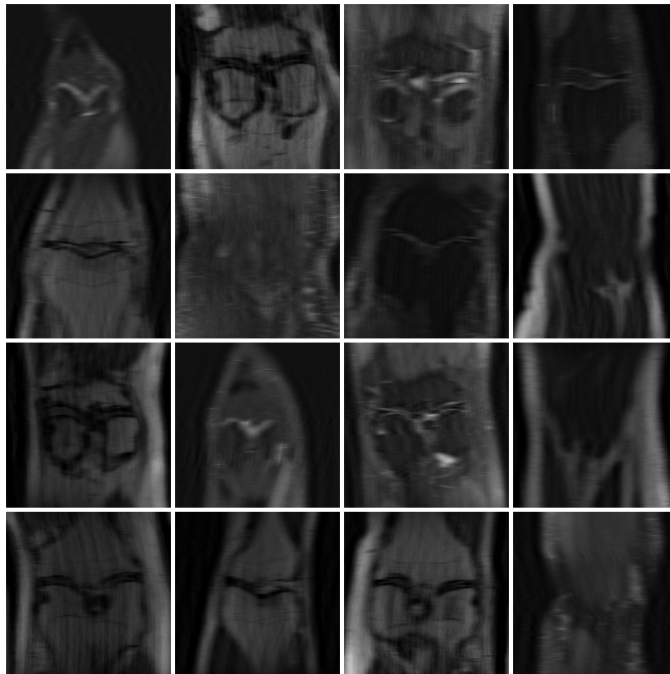


Figure 12. Example of k -space sub-sampling observations with acceleration factor $R = 6$ for the accelerated MRI experiment. We represent each observation by its zero-filled inverse, where missing frequencies are set to zero before taking the inverse discrete Fourier transform.

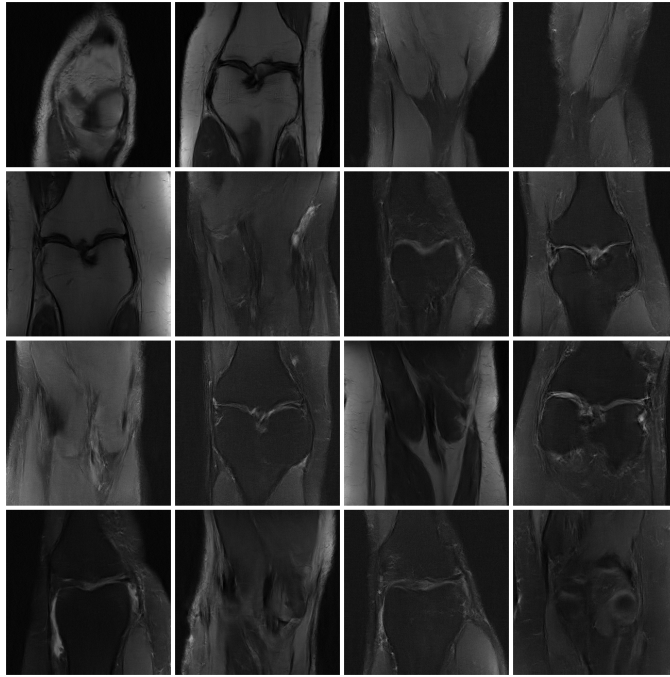


Figure 13. Example of samples from the final model $q_{\theta_k}(x)$ for the accelerated MRI experiment. The samples present varied and coherent global structures. Samples seem slightly less sharp than real scans (see Figure 11), but do not present artifacts typical to unresolved frequencies (see Figure 12).

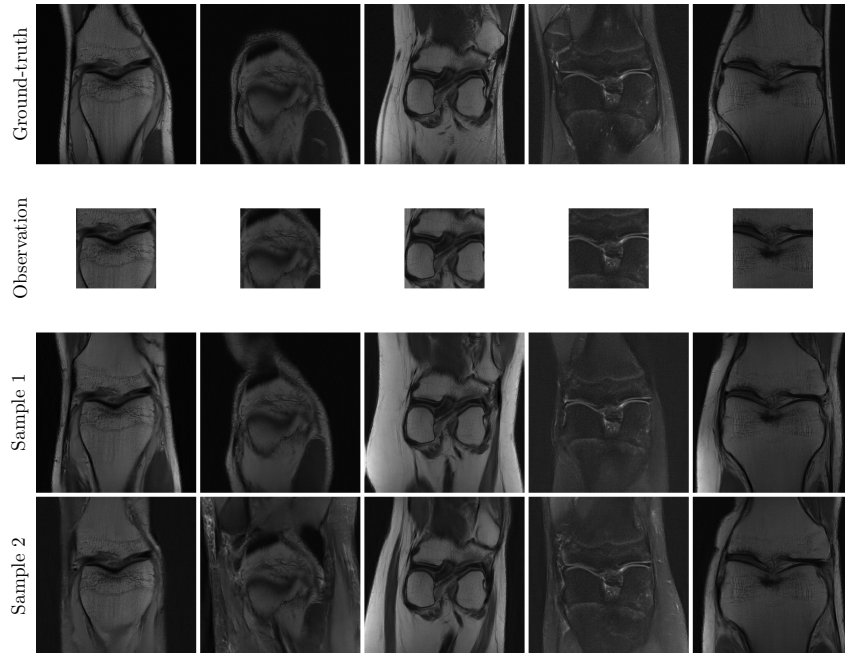


Figure 14. Examples of posterior samples using a diffusion prior trained from k -space observations only. The forward process crops the latent x to a centered 160×160 window. Moment matching posterior sampling is used to sample from the posterior. Samples are consistent with the ground-truth where observed, but present plausible variations elsewhere.

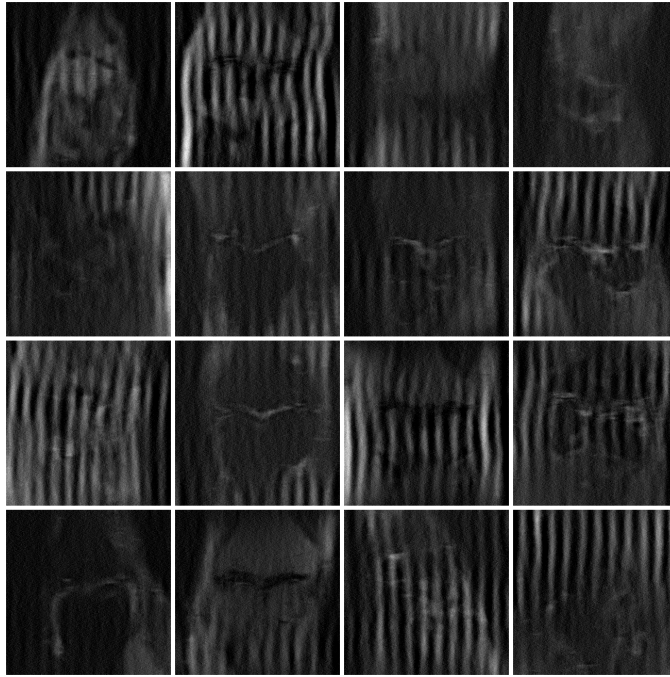


Figure 15. Example of samples from the model $q_{\theta_k}(x)$ after $k = 2$ EM iterations for the accelerated MRI experiment when the heuristic $(I + \Sigma_t^{-1})^{-1}$ is used for $\mathbb{V}[x | x_t]$. The samples start to present vertical artifacts due to poor sampling.

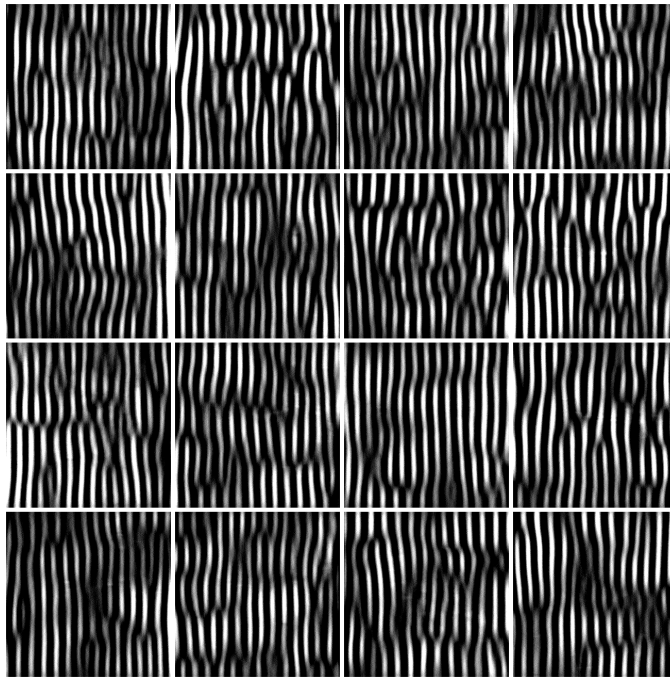


Figure 16. Example of samples from the model $q_{\theta_k}(x)$ after $k = 4$ EM iterations for the accelerated MRI experiment when the heuristic $(I + \Sigma_t^{-1})^{-1}$ is used for $\mathbb{V}[x | x_t]$. The artifacts introduced by the poor sampling get amplified at each iteration, leading to a total collapse after few iterations.