

END-TO-END LEARNING OF GAUSSIAN MIXTURE PRIORS FOR DIFFUSION SAMPLER

Anonymous authors

Paper under double-blind review

ABSTRACT

Diffusion models optimized via variational inference (VI) have emerged as a promising tool for generating samples from unnormalized target densities. These models create samples by simulating a stochastic differential equation, starting from a simple, tractable prior, typically a Gaussian distribution. However, when the support of this prior differs greatly from that of the target distribution, diffusion models often struggle to explore effectively or suffer from large discretization errors. Moreover, learning the prior distribution can lead to mode-collapse, exacerbated by the mode-seeking nature of reverse Kullback-Leibler divergence commonly used in VI. To address these challenges, we propose end-to-end learnable Gaussian mixture priors (GMPs). GMPs offer improved control over exploration, adaptability to target support, and increased expressiveness to counteract mode collapse. We further leverage the structure of mixture models by proposing a strategy to iteratively refine the model through the addition of mixture components during training. Our experimental results demonstrate significant performance improvements across a diverse range of real-world and synthetic benchmark problems when using GMPs without requiring additional target evaluations.

1 INTRODUCTION

Sampling methods are designed to address the challenge of generating approximate samples or estimating the intractable normalization constant Z for a probability density π on \mathbb{R}^d of the form

$$\pi(\mathbf{x}) = \frac{\rho(\mathbf{x})}{Z}, \quad Z = \int_{\mathbb{R}^d} \rho(\mathbf{x}) d\mathbf{x}, \quad (1)$$

where $\rho : \mathbb{R}^d \rightarrow (0, \infty)$ can be evaluated pointwise. This formulation has broad applications in fields such as Bayesian statistics, the natural sciences (Liu & Liu, 2001; Stoltz et al., 2010; Frenkel & Smit, 2023).

Monte Carlo (MC) methods (Hammersley, 2013), Annealed Importance Sampling (AIS) (Neal, 2001), and their Sequential Monte Carlo (SMC) extensions (Del Moral et al., 2006; Arbel et al., 2021; Matthews et al., 2022; Midgley et al., 2022) have long been regarded as the gold standard for tackling complex sampling problems. An alternative approach is variational inference (VI) (Blei et al., 2017), which approximates an intractable target distribution by parameterizing a family of tractable distributions. Recently, there has been growing interest in diffusion models (Zhang & Chen, 2021; Berner et al., 2022; Richter et al., 2023; Vargas et al., 2023a;b), which employ stochastic processes to transport samples from a simple, tractable prior distribution to the target distribution. While diffusion models have shown great success in generative modeling (Ho et al., 2020; Song et al., 2020), their application to sampling tasks introduces unique challenges.

We identify these challenges as follows (C1–C3): Unlike generative modeling, where the support of the target distribution is often known, in sampling tasks, the target’s support is usually unknown. This makes it difficult to set the prior appropriately and requires the model to explore the relevant regions of the space—an exploration that becomes exponentially harder as dimensionality increases (C1). Additionally, large discrepancies between the support of the prior and the target distribution can lead to highly non-linear dynamics, necessitating many diffusion steps to mitigate discretization errors (C2). Finally, while joint optimization of the prior and diffusion process is possible, using

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

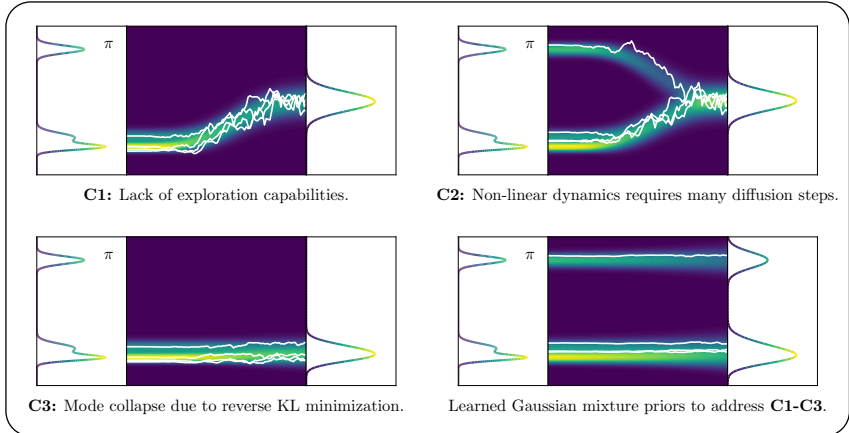


Figure 1: Illustration of challenges (C1-C3) associated with diffusion-based sampling methods and how learned Gaussian mixture priors address them (bottom right). Here, π denotes the target distribution.

simple priors like Gaussians can result in mode collapse due to the mode-seeking behavior of the reverse Kullback-Leibler (KL) divergence commonly used in VI (C3). These challenges are illustrated in Figure 1.

Outline. In Section 3, we present an overview of diffusion-based sampling methods within the framework of variational inference. Next, we discuss the necessary adaptations for supporting the learning of arbitrary prior distributions, illustrated through specific examples of diffusion models (Section 4). We then provide a rationale for our choice of Gaussian Mixture Priors (GMPs) and introduce a novel training scheme designed to iteratively refine diffusion models during training (Section 5). Finally, in Section 6, we assess our method through experiments on a range of real-world and synthetic benchmark problems, demonstrating consistent improvements in performance.

2 RELATED WORK

Sampling and Variational Inference. Numerous works have studied the problem of sampling from unnormalized densities to estimate the partition function Z , including Monte Carlo (MC) methods such as Markov Chain Monte Carlo (MCMC) (Kass et al., 1998) and Sequential Importance Sampling (Liu et al., 2001). Seminal works include Annealed Importance Sampling (Neal, 2001) and its Sequential Monte Carlo extensions (Del Moral et al., 2006; Arbel et al., 2021; Wu et al., 2020; Matthews et al., 2022; Midgley et al., 2022). Another line of work approaches the sampling problem by utilizing tools from optimization to fit a parametric family of distributions to the target density π , known as Variational Inference (VI) (Blei et al., 2017). To that end, one typically uses the reverse Kullback-Leibler divergence, although other discrepancies have been studied (Li & Turner, 2016; Midgley et al., 2022; Dieng et al., 2017; Richter et al., 2020; Wan et al., 2020; Naesseth et al., 2020).

Diffusion-based Sampling Methods. Recently, there has been growing interest in combining Monte Carlo methods with variational techniques by constructing a sequence of variational distributions through the parameterization of Markov chains (Naesseth et al., 2018; Geffner & Domke, 2021; Thin et al., 2021; Zhang et al., 2021; Chen et al., 2024). In the limit of infinitely many steps, these Markov chains converge to stochastic differential equations (SDEs) (Särkkä & Solin, 2019), which has led to further research on diffusion-based models for sampling, particularly in light of advances in generative modeling (Ho et al., 2020; Song et al., 2020). One line of work considered parameterized drift functions to improve annealed Langevin diffusions in the overdamped (Doucet et al., 2022a) or underdamped (Geffner & Domke, 2022) regime. Another line of work casts diffusion-based sampling as a stochastic optimal control problem (Dai Pra, 1991) including denoising diffusion models (Bernier et al., 2022; Vargas et al., 2023a), and Föllmer sampling (Föllmer, 2005; Zhang & Chen, 2021; Vargas et al., 2023b). A unifying view was later provided by Vargas et al. (2024); Richter et al. (2023). Further extensions to diffusion-based sampling methods have been proposed such as improved learning objectives (Zhang et al., 2023; Akhound-Sadegh et al., 2024) or combinations with sequential importance sampling (Phillips et al., 2024). Another study

leverages physics-informed neural networks (PINNs, (Raissi et al., 2019)) to learn the Fokker-Planck equation governing the density evolution of the diffusion process (Sun et al., 2024).

3 PRELIMINARIES

In this section, we offer a concise overview of diffusion models within the context of variational inference. Our discussion draws primarily from the works of Richter et al. (2023); Vargas et al. (2024). While these studies emphasize the continuous-time perspective, we adopt an approach that largely emphasizes discrete time, aiming to make the topic more accessible to readers without a background in stochastic calculus.

3.1 CONTROLLED DIFFUSIONS, DISCRETIZATION, AND COUPLINGS

We consider two \mathbb{R}^d -valued stochastic processes on the time-interval $[0, T]$: $\vec{\mathbf{X}}$ starts from a prior distribution p_0 and runs forward in time whereas $\overleftarrow{\mathbf{X}}$ starts from the target distribution $p_T = \pi$ and runs backward in time. These processes are governed by the stochastic differential equations (SDEs) given by controlled diffusions, that is,

$$d\mathbf{X}_t = [f(\mathbf{X}_t, t) + \sigma u^\theta(\mathbf{X}_t, t)] dt + \sqrt{2}\sigma d\mathbf{B}_t, \quad \mathbf{X}_0 \sim p_0, \quad (2a)$$

$$d\mathbf{X}_t = [f(\mathbf{X}_t, t) - \sigma v^\gamma(\mathbf{X}_t, t)] dt + \sqrt{2}\sigma d\mathbf{B}_t, \quad \mathbf{X}_T \sim p_T = \pi, \quad (2b)$$

with drift, and parameterized control functions $f, u^\theta, v^\gamma : \mathbb{R}^d \times [0, T] \rightarrow \mathbb{R}^d$, respectively. Further, $(\mathbf{B}_t)_{t \in [0, T]}$ is a d -dimensional Brownian motion and $\sigma \in \mathbb{R}^+$ a diffusion coefficient. For integration, we consider the Euler-Maruyama (EM) method with constant discretization step size $\delta t \geq 0$ such that $N = T/\delta t$ is an integer. To simplify notation, we write \mathbf{x}_n , instead of $\mathbf{x}_{n\delta t}$. Integrating Eq. 2a yields

$$\mathbf{x}_{n+1} = \mathbf{x}_n + [f(\mathbf{x}_n, n) + \sigma u^\theta(\mathbf{x}_n, n)] \delta t + \sigma \sqrt{2\delta t} \epsilon_n, \quad \mathbf{x}_0 \sim p_0, \quad (3)$$

where $\epsilon_n \sim \mathcal{N}(0, \mathbf{I})$. The EM discretizations of $\vec{\mathbf{X}}$ and $\overleftarrow{\mathbf{X}}$ admit the following Markov Processes

$$\mathcal{P}^\theta(\mathbf{x}_{0:N}) = p_0(\mathbf{x}_0) \prod_{n=1}^N F_n^\theta(\mathbf{x}_n | \mathbf{x}_{n-1}), \quad \text{and} \quad (4)$$

$$\mathcal{Q}^\gamma(\mathbf{x}_{0:N}) = p_T(\mathbf{x}_N) \prod_{n=1}^N B_{n-1}^\gamma(\mathbf{x}_{n-1} | \mathbf{x}_n), \quad (5)$$

in a sense that \mathcal{P}^θ and \mathcal{Q}^γ converge to the law of $\vec{\mathbf{X}}$ and $\overleftarrow{\mathbf{X}}$, respectively, as $\delta t \rightarrow 0$. Here,

$$F_n^\theta(\mathbf{x}_{n+1} | \mathbf{x}_n) = \mathcal{N}(\mathbf{x}_{n+1} | \mathbf{x}_n + [f(\mathbf{x}_n, n) + \sigma u^\theta(\mathbf{x}_n, n)] \delta t, 2\sigma^2 \delta t \mathbf{I}), \quad \text{and} \quad (6)$$

$$B_{n-1}^\gamma(\mathbf{x}_{n-1} | \mathbf{x}_n) = \mathcal{N}(\mathbf{x}_{n-1} | \mathbf{x}_n - [f(\mathbf{x}_n, n) - \sigma v^\gamma(\mathbf{x}_n, n)] \delta t, 2\sigma^2 \delta t \mathbf{I}). \quad (7)$$

The goal of diffusion-based sampling methods is to obtain a coupling/bridge between p_0 and $p_T = \pi$ by learning control functions u^θ and v^γ such that

$$\mathcal{P}^\theta(\mathbf{x}_{0:N}) = \mathcal{Q}^\gamma(\mathbf{x}_{0:N}). \quad (8)$$

Assuming Eq. 8 holds, we have $\int \mathcal{P}^\theta(\mathbf{x}_{0:N}) d\mathbf{x}_{0:N-1} = \int \mathcal{Q}^\gamma(\mathbf{x}_{0:N}) d\mathbf{x}_{0:N-1} = \pi(\mathbf{x}_N)$, meaning that we can sample $\mathbf{x}_0 \sim p_0$ and integrate the ‘forward’ diffusion process $\vec{\mathbf{X}}$ to obtain samples from π . In contrast, the ‘backward’ process $\overleftarrow{\mathbf{X}}$ is not needed for generating samples from π , but is required for estimating the normalization constant Z , and obtaining a tractable optimization objective which is discussed in the next section. Lastly, we want to highlight that this formulation of diffusion-based sampling is very generic and that most instances of samplers, such as denoising diffusion samplers (Berner et al., 2022; Vargas et al., 2023a), can be recovered by choosing the drift and/or control functions in Eq. 2 appropriately. We refer the interested reader to Richter et al. (2023) for further details.

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

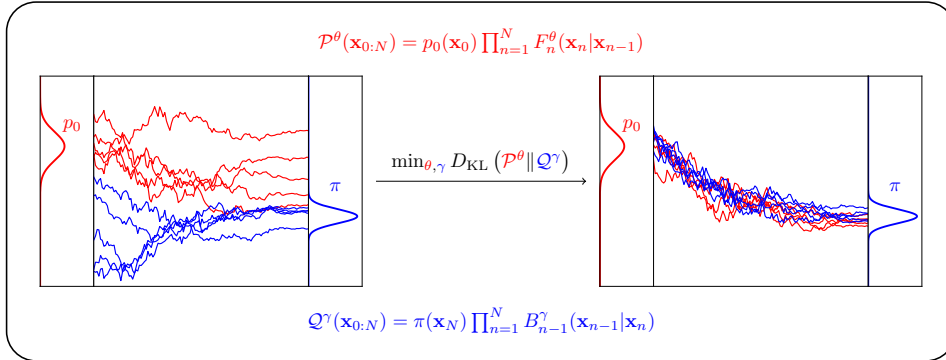


Figure 2: Diffusion-Based Sampling: The goal is to align two parameterized Markov Processes \mathcal{P}^θ and \mathcal{Q}^γ . The former starts at the prior p_0 and runs forward in time while the latter starts at the target π and runs backward.

3.2 VARIATIONAL INFERENCE FOR DIFFUSION MODELS

Variational Inference (Blei et al., 2017) uses a parameterized tractable distribution p^θ and minimizes a divergence to the target distribution π with respect to its parameters θ , typically the Kullback-Leibler divergence, i.e.,

$$D_{\text{KL}}(p^\theta(\mathbf{x}) \| \pi(\mathbf{x})) = \mathbb{E}_{\mathbf{x} \sim p^\theta} \left[\log \frac{p^\theta(\mathbf{x})}{\rho(\mathbf{x})} \right] + \log Z = -\text{ELBO}(\theta) + \log Z, \quad (9)$$

It directly follows that minimizing D_{KL} , or equivalently, maximizing the ELBO¹ does not require access to the true normalization constant Z as it is independent of θ . Moreover, using the fact that $D_{\text{KL}} \geq 0$, it is straightforward to see that $\text{ELBO}(\theta) \leq \log Z$.

In the case of diffusion models, we are interested in learning the parameters θ, γ of the control functions u^θ, v^γ . Directly minimizing D_{KL} between $p_T^\theta(\mathbf{x}_N) = \int \mathcal{P}^\theta(\mathbf{x}_{0:N}) d\mathbf{x}_{0:N-1}$ and π is challenging. However, the data-processing inequality (Cover, 1999), that is,

$$D_{\text{KL}}(p_T^\theta(\mathbf{x}_T) \| \pi(\mathbf{x}_T)) \leq D_{\text{KL}}(\mathcal{P}^\theta(\mathbf{x}_{0:N}) \| \mathcal{Q}^\gamma(\mathbf{x}_{0:N})), \quad (10)$$

provides an auxiliary, tractable, objective for optimizing (θ, γ) , that is,

$$D_{\text{KL}}(\mathcal{P}^\theta(\mathbf{x}_{0:N}) \| \mathcal{Q}^\gamma(\mathbf{x}_{0:N})) = \underbrace{\mathbb{E}_{\mathbf{x}_{0:N} \sim \mathcal{P}^\theta} \left[\log \frac{p_0(\mathbf{x}_0)}{\rho(\mathbf{x}_N)} + \sum_{n=1}^N \log \frac{F_n^\theta(\mathbf{x}_n | \mathbf{x}_{n-1})}{B_{n-1}^\gamma(\mathbf{x}_{n-1} | \mathbf{x}_n)} \right]}_{-\mathcal{L}(\theta, \gamma)} + \log Z, \quad (11)$$

where $\mathcal{L}(\theta, \gamma)$ is often referred to as augmented or extended evidence lower bound, as it has additional looseness due to the latent variables $\mathbf{x}_{0:N-1}$ (Geffner & Domke, 2021). Note that the VI setting for optimizing diffusion models is different from techniques used when samples from the target, i.e., $\mathbf{x}_N \sim \pi$ are available. The former requires simulations $\mathbf{x}_{0:N} \sim \mathcal{P}^\theta$ for optimization, while the latter minimizes the forward KL $D_{\text{KL}}(\mathcal{Q}^\gamma(\mathbf{x}_{0:N}) \| \mathcal{P}^\theta(\mathbf{x}_{0:N}))$, allowing for simulation-free optimization techniques such as denoising score-matching (Vincent, 2011; Song & Ermon, 2019) or bridge matching (Liu et al., 2022; Shi et al., 2024). Moreover, recent works consider minimizing other loss functions than the KL divergence in Eq. 11. A recent overview of possible alternatives can be found in Domingo-Enrich (2024). For further details, the interested reader is referred to Berner et al. (2022); Vargas et al. (2024).

4 END-TO-END LEARNING OF PRIOR DISTRIBUTIONS

We aim to learn a parametric prior p_0^ϕ with parameters ϕ end-to-end when maximizing the extended ELBO \mathcal{L} (Eq. 11). To that end, we consider two requirements:

¹Evidence Lower Bound. The terminology stems from Bayesian inference, where $\log Z$ is equivalent to the evidence of the data.

1. We can compute gradients of \mathcal{L} with respect to ϕ .
2. There exists a ϕ and γ such that $p_0^\phi(\mathbf{x}_0) = \int \mathcal{Q}^\gamma(\mathbf{x}_{0:N})d\mathbf{x}_{1:N}$.

For the former, we assume that $p_0^\phi(\mathbf{x}_0)$ is amendable to the reparameterization trick², i.e., we can express a sample \mathbf{x}_0 from p_0^ϕ as a deterministic function of a random variable ξ with some fixed distribution and the parameters ϕ , i.e., $\mathbf{x}_0 = g(\xi, \phi)$. We can then obtain gradients of

$$\mathcal{L}(\theta, \gamma, \phi) = \mathbb{E}_{\mathbf{x}_{0:N} \sim \mathcal{P}^{\theta, \phi}} \left[\log \frac{\rho(\mathbf{x}_N)}{p_0^\phi(\mathbf{x}_0)} + \sum_{n=1}^N \log \frac{B_{n-1}^\gamma(\mathbf{x}_{n-1}|\mathbf{x}_n)}{F_n^\theta(\mathbf{x}_n|\mathbf{x}_{n-1})} \right], \quad (12)$$

with $\mathcal{P}^{\theta, \phi}(\mathbf{x}_{0:N}) = p_0^\phi(\mathbf{x}_0) \prod_{n=1}^N F_n^\theta(\mathbf{x}_n|\mathbf{x}_{n-1})$, with respect to ϕ , by differentiating through the stochastic process

$$\mathbf{x}_{n+1} = \mathbf{x}_n + [f(\mathbf{x}_n, n) + \sigma u^\theta(\mathbf{x}_n, n)] \delta t + \sigma \sqrt{2\delta t} \epsilon_n, \quad \mathbf{x}_0 = g(\xi, \phi). \quad (13)$$

The second requirement is necessary to obtain a coupling between p_0^ϕ and π , i.e., to satisfy Eq. 8. This requirement is trivially fulfilled for a controlled process $\tilde{\mathbf{X}}$, where we can learn a v^γ such that \mathcal{Q}^γ transports π back to p_0^ϕ . However, this requirement can be more intricate for other processes and will be discussed in the next sections. In particular, we look at two instances of Eq. 2, namely denoising diffusion models (Berner et al., 2022; Vargas et al., 2023a) and annealed Langevin diffusions (Doucet et al., 2022a; Vargas et al., 2024).

4.1 DENOISING DIFFUSION MODELS

Denoising diffusion models use an Ornstein Uhlenbeck (OU) process³ for $\tilde{\mathbf{X}}$, that is,

$$d\mathbf{X}_t = -\sigma^2 \mathbf{X}_t dt + \sqrt{2}\sigma d\mathbf{B}_t, \quad \mathbf{X}_T \sim p_T = \pi, \quad (14)$$

and, hence, a special case of Eq. 2 with $f(\mathbf{X}_t, t) = -\sigma^2 \mathbf{X}_t$ and $v^\gamma = 0$. Assuming a sufficiently large σ (or T), it holds that $p_0(\mathbf{x}_0) = \int \mathcal{Q}(\mathbf{x}_{0:N})d\mathbf{x}_{1:N} \approx \mathcal{N}(0, \mathbf{I})$. In other words, the OU process transports the target π to a Gaussian distribution. We extend denoising diffusion models to support learning arbitrary priors based on Proposition 1, whose proof can be found in Appendix A.1.

Proposition 1. *Let $\tilde{\mathbf{X}}$ be a (uncontrolled) stochastic process as defined in Eq. 2 with $v^\gamma = 0$, starting from $p_T = \pi$. For a time-independent drift, i.e., $f(\mathbf{x}, t) = f(\mathbf{x})$, the stationary distribution $p_s(\mathbf{x})$ for which $\frac{\partial p_s(\mathbf{x}_t)}{\partial t} = 0$ holds, is given by*

$$p_s(\mathbf{x}) = \frac{1}{Z_s} \exp \left(-\frac{1}{\sigma^2} \int f(\mathbf{x}) d\mathbf{x} \right), \quad (15)$$

with normalization constant Z_s .

Rewriting Eq. 15, yields $f = \sigma^2 \nabla_{\mathbf{x}} \log p_s$, resulting in the SDE

$$d\mathbf{X}_t = \sigma^2 \nabla_{\mathbf{x}} \log p_s(\mathbf{X}_t) dt + \sqrt{2}\sigma d\mathbf{B}_t, \quad \mathbf{X}_T \sim p_T = \pi, \quad (16)$$

with stationary distribution $p_s(\mathbf{x})$. Note that denoising diffusion models leverage this result by setting $p_s = \mathcal{N}(0, \mathbf{I})$, resulting in the OU process (Eq. 14) since $\nabla_{\mathbf{x}} \log p_s(\mathbf{x}) = -\mathbf{x}$. Hence, we can adapt existing denoising diffusion sampling methods (Vargas et al., 2023a; Berner et al., 2022) to arbitrary priors p^ϕ using

$$d\mathbf{X}_t = \sigma^2 \nabla_{\mathbf{x}} \log p^\phi(\mathbf{X}_t) dt + \sqrt{2}\sigma d\mathbf{B}_t, \quad \mathbf{X}_T \sim p_T = \pi. \quad (17)$$

However, contrary to the OU process, where the relaxation time, i.e., the time scale over which the system loses memory of its initial conditions and approaches its stationary distribution, can be estimated analytically, it is unknown for general p^ϕ and is only guaranteed as $T \rightarrow \infty$ (Roberts & Tweedie, 1996).

²Note that this requirement is not necessary when minimizing loss function where the expectation is not computed with respect to samples from \mathcal{P}^θ . For further details see e.g. (Richter et al., 2023).

³Often referred to as Variance Preserving (VP) SDE, a term coined by Song et al. (2020).

We address this by additionally learning $T = N\delta t$ by treating the discretization step size δt as a learnable parameter. As such, the parameters ϕ of the stationary distribution, i.e., the prior distribution and the discretization step size δt are optimized jointly by maximizing the extended ELBO

$$\mathcal{L}(\theta, \phi, \delta t) = \mathbb{E}_{\mathbf{x}_{0:N} \sim \mathcal{P}^{\theta, \phi, \delta t}} \left[\log \frac{\rho(\mathbf{x}_N)}{p_0^\phi(\mathbf{x}_0)} + \sum_{n=1}^N \log \frac{B_{n-1}^{\phi, \delta t}(\mathbf{x}_{n-1} | \mathbf{x}_n)}{F_n^{\theta, \phi, \delta t}(\mathbf{x}_n | \mathbf{x}_{n-1})} \right], \quad (18)$$

with additional parameters $\phi, \delta t$. Proposition 1 thus suggests, that for any ϕ , there exists a δt such that $p_0^\phi(\mathbf{x}_0) = \int \mathcal{Q}^{\phi, \delta t}(\mathbf{x}_{0:N}) d\mathbf{x}_{1:N}$ as $N \rightarrow \infty$. Empirically, we observe substantial improvements for finite values of N , as demonstrated in Section 6

4.2 ANNEALED LANGEVIN DIFFUSIONS

Annealed Langevin Diffusions use an annealed version of the (overdamped) Langevin diffusion equation by constructing a sequence of distributions $(\pi_t)_{t \in [0, T]}$ that anneal smoothly from the prior distribution $\pi_0 = p_0$ to the target distribution $\pi_T = \pi$. One typically uses the geometric average, that is, $\pi_t(\mathbf{x}) = p_0(\mathbf{x})^{\beta_t} \pi(\mathbf{x})^{1-\beta_t}$, for β_t monotonically increasing in t with $\beta_0 = 0$ and $\beta_T = 1$. When learning the prior, we can use a parametric annealing, i.e., $\pi_t^\phi(\mathbf{x}) = p_0^\phi(\mathbf{x})^{\beta_t} \pi(\mathbf{x})^{1-\beta_t}$. The corresponding stochastic processes $\vec{\mathbf{X}}$ and $\vec{\mathbf{X}}$ can be described as an instance of Eq. 2 given by

$$d\mathbf{X}_t = \left[\sigma^2 \nabla_{\mathbf{x}} \log \pi_t^\phi(\mathbf{X}_t) + \sigma u^\theta(\mathbf{X}_t, t) \right] dt + \sqrt{2} \sigma d\mathbf{B}_t, \quad \mathbf{X}_0 \sim p_0 = p^\phi, \quad (19)$$

$$d\mathbf{X}_t = \left[\sigma^2 \nabla_{\mathbf{x}} \log \pi_t^\phi(\mathbf{X}_t) - \sigma v^\gamma(\mathbf{X}_t, t) \right] dt + \sqrt{2} \sigma d\mathbf{B}_t, \quad \mathbf{X}_T \sim p_T = \pi, \quad (20)$$

when setting $f = \nabla_{\mathbf{x}} \log \pi_t^\phi$. Note that $\nabla_{\mathbf{x}} \log \pi_t^\phi$ can be computed without knowing the normalization constant Z of π . Different variants can be derived from using either controlled or uncontrolled processes: Monte Carlo Diffusions (MCD) (Doucet et al., 2022b) uses a controlled process $\vec{\mathbf{X}}$ but uncontrolled $\vec{\mathbf{X}}$ and Controlled Monte Carlo Diffusions (CMCD) (Vargas et al., 2024) control both processes. Since both methods use controlled backward processes $\vec{\mathbf{X}}$, the second requirement is satisfied. Finally, while this work focuses on overdamped approaches, we want to highlight that there exist methods that are based on the underdamped Langevin equation (Geffner & Domke, 2021; Geffner & Domke, 2022), however, the idea of learning a prior end-to-end straightforwardly transfers to these approaches.

5 GAUSSIAN MIXTURE PRIORS AND ITERATIVE MODEL REFINEMENT

In this work, we focus on end-to-end learned Gaussian mixture priors (GMPs), that is,

$$p_0^\phi(\mathbf{x}_0) = \sum_{k=1}^K \alpha_k p_0^{\phi_k}(\mathbf{x}_0) = \sum_{k=1}^K \alpha_k \mathcal{N}(\mathbf{x}_0 | \mu_k, \Sigma_k), \quad \alpha_k \geq 0, \quad \sum_{k=1}^K \alpha_k = 1, \quad (21)$$

with mixture weights α_k , Gaussian components $p_0^{\phi_k}(\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_0 | \mu_k, \Sigma_k)$ and parameters $\phi = \bigcup_{k=1}^K \{\alpha_k, \phi_k\}$ with $\phi_k = \{\mu_k, \Sigma_k\}$. Having established how the prior is learned in Section 4, we discuss desirable properties to address the challenges outlined in Section 1 and how GMPs address them.

A key objective is to improve the exploration capabilities of diffusion-based sampling methods to address C1. GMPs allow control over exploration by adjusting the initial variance of each Gaussian component. Additionally, the means of the Gaussian components can be initialized to incorporate prior knowledge of the target density, even if this knowledge is limited to a rough estimate of the target’s support. This aspect will be elaborated on later in this section.

Another important consideration is to adjust the support of the prior such that it matches the target density, which reduces the complexity of the dynamics and, in turn, minimizes the number of diffusion steps required. GMPs demonstrate rapid adaptation capabilities, partially through their small parameter count, making them particularly suitable for addressing C2.

To prevent the model from focusing only on a subset of the target support (C3), which may occur due to the optimization of the mode-seeking reverse KL divergence, we require a more expressive

distribution than a single Gaussian prior. GMPs provide a solution by combining multiple Gaussian components, each of which can focus on different subsets of the target support.

Finally, efficient evaluation of p_0^ϕ is crucial, as it must be performed at each discretization step of the stochastic differential equation (SDE) that governs the diffusion process. This requirement is satisfied by GMPs, particularly when using diagonal covariance matrices.

Iterative Model Refinement. Gradually increasing the model complexity during the optimization process has demonstrated promising results in previous studies (Guo et al., 2016; Miller et al., 2017; Arenz et al., 2018; Cranko & Nock, 2019), and is directly applicable to our approach. We begin with an initial prior distribution $p_0^\phi = p_0^{\phi_1}$ with parameters ϕ_1 . These parameters are optimized using Eq. 48. After a predefined criterion is met, such as a fixed number of iterations, a second distribution $p_0^{\phi_2}$ is added, forming a new prior: $p_0^\phi = \alpha_1 p_0^{\phi_1} + \alpha_2 p_0^{\phi_2}$, with $\alpha \in \mathbb{R}^+$ and $\alpha_1 + \alpha_2 = 1$. This process is repeated, resulting in a mixture model $p_0^\phi(\mathbf{x}) = \sum_{k=1}^K \alpha_k p_0^{\phi_k}(\mathbf{x})$.

We identify the benefits of this iterative scheme as twofold: First, it can simplify optimization by focusing on learning a subset of parameters ϕ_k at a time, rather than jointly optimizing all ϕ_k (Bengio et al., 2009). Second, it enables the initialization of newly added components based on a partially trained model, potentially preventing mixture components to focus on similar parts of the target support. For GMPs, for instance, the mean of a new component μ_{new} can be placed in a promising region, potentially informed by prior knowledge of the task or by running a π -invariant Markov chain to obtain a set of promising samples. More generally, consider a set of candidate samples $\mathcal{C} = \{\mathbf{x}_i\}_{i=1}^C$. We propose initializing the mean of a new component μ_{new} as follows:

$$\mu_{\text{new}} = \arg \max_{\mathbf{x}_0 \in \mathcal{C}} \mathbb{E}_{\mathbf{x}_{1:N} \sim \mathcal{P}^{\theta, \phi, \delta t}} \left[\log \frac{\rho(\mathbf{x}_N)}{p_0^\phi(\mathbf{x}_0)} + \sum_{n=1}^N \log \frac{B_{n-1}^{\gamma, \phi, \delta t}(\mathbf{x}_{n-1} | \mathbf{x}_n)}{F_n^{\theta, \phi, \delta t}(\mathbf{x}_n | \mathbf{x}_{n-1})} \right], \quad (22)$$

where p_0^ϕ is the current model. This heuristic balances exploration and exploitation by favoring samples with high target likelihood and low prior likelihood, while also accounting for the diffusion process.

6 NUMERICAL EVALUATION

In this section, we test the impact of our proposed end-to-end learning scheme for prior distributions. Specifically, we consider three distinct settings: First, we evaluate these methods with a Gaussian prior that is fixed during training. Second and third, we consider learned Gaussian (GP) and Gaussian mixture priors (GMP). We indicate these different settings as X, X-GP, and X-GMP, respectively, where X is the corresponding acronym of the diffusion-based sampling methods. We consider four different methods: Time-Reversed Diffusion Sampler (DIS) (Bernier et al., 2022), Monte Carlo Diffusions (MCD) (Doucet et al., 2022b), Controlled Monte Carlo Diffusions (CMCD) (Vargas et al., 2024) and Diffusion Bridge Sampler (DBS) (Richter et al., 2023). A summary is shown in Table 1. It is worth noting that we do not separately consider the Denoising Diffusion Sampler (DDS) (Vargas et al., 2023a), as it can be viewed as a special case of DIS. For reference, we consider Gaussian (GVI) and Gaussian mixture (GMVI) mean-field approximations (Wainwright & Jordan, 2008), both of which are special cases of the aforementioned methods for $N = 0$ diffusion steps with $K = 1$, $K \geq 1$, respectively (cf. Appendix B). Lastly, we consider three competing state-of-the-art methods, namely, Sequential Monte Carlo (SMC) (Del Moral et al., 2006), Continual Repeated Annealed Flow Transport (CRAFT) (Matthews et al., 2022), and Flow Annealed Importance Sampling Bootstrap (FAB) (Midgley et al., 2022).

For evaluation, we consider the effective sample size (ESS) and the marginal or extended evidence lower bound as performance criteria. Both are denoted as ‘ELBO’ for convenience. Next, if the ground truth normalization constant Z is available, we use an importance-weighted estimate \hat{Z} to

METHOD (X)	f	u^θ	v^γ
MCD	$\nabla \log \pi_t^\phi$	✗	✓
CMCD ³	$\nabla \log \pi_t^\phi$	✓	✓
DIS	$\nabla \log p^\phi$	✓	✗
DBS	ANY	✓	✓

Table 1: Diffusion-based sampling methods considered in this work based on Eq. 2. Crosses indicate that the control is set to zero.

³Vargas et al. (2024) use the same in control in \vec{X} and \vec{X} by leveraging Nelson’s relation (Nelson, 2020).

compute the estimation error $\Delta \log Z = |\log Z - \log \hat{Z}|$. Additionally, if samples from the target π are available, we compute the Sinkhorn distance \mathcal{W}_γ^2 (Cuturi, 2013).

To ensure a fair comparison, all experiments are conducted under identical settings. Our evaluation methodology adheres to the protocol by Blessing et al. (2024). For a comprehensive overview of the experimental setup see Appendix C. Moreover, a comprehensive set of ablation studies and additional experiments, are provided in Appendix D.

Method	Funnel ($d = 10$)			
	ELBO \uparrow	$\Delta \log Z$ \downarrow	ESS \uparrow	\mathcal{W}_γ^2 \downarrow
GVI	-1.841 \pm 0.003	0.691 \pm 0.070	0.092 \pm 0.006	178.007 \pm 0.164
GMVI	<u>-0.212\pm0.001</u>	<u>0.056\pm0.004</u>	<u>0.744\pm0.018</u>	<u>102.826\pm0.109</u>
MCD	-0.721 \pm 0.003	0.201 \pm 0.017	0.207 \pm 0.012	164.882 \pm 0.363
MCD-GP	-0.724 \pm 0.003	0.173 \pm 0.046	0.206 \pm 0.026	164.967 \pm 0.334
MCD-GMP	<u>-0.059\pm0.002</u>	<u>0.014\pm0.001</u>	<u>0.922\pm0.012</u>	<u>100.174\pm0.174</u>
CMCD	-0.210 \pm 0.002	0.020 \pm 0.006	0.588 \pm 0.013	104.652 \pm 0.593
CMCD-GP	-0.211 \pm 0.002	0.023 \pm 0.003	0.567 \pm 0.023	104.644 \pm 0.710
CMCD-GMP	<u>-0.027\pm0.001</u>	<u>0.005\pm0.000</u>	0.950\pm0.004	<u>102.027\pm0.200</u>
DIS	-0.286 \pm 0.002	0.041 \pm 0.008	0.483 \pm 0.025	107.458 \pm 0.670
DIS-GP	-0.296 \pm 0.002	0.047 \pm 0.003	0.498 \pm 0.021	107.458 \pm 0.826
DIS-GMP	<u>-0.058\pm0.002</u>	<u>0.019\pm0.002</u>	<u>0.929\pm0.017</u>	100.093\pm0.028
DBS	-0.180 \pm 0.002	0.019 \pm 0.005	0.600 \pm 0.014	102.964 \pm 0.442
DBS-GP	-0.187 \pm 0.003	0.021 \pm 0.003	0.603 \pm 0.014	102.653 \pm 0.586
DBS-GMP	<u>-0.047\pm0.002</u>	<u>0.012\pm0.002</u>	<u>0.949\pm0.008</u>	<u>100.230\pm0.088</u>
SMC	-0.242 \pm 0.047	0.187 \pm 0.054	-	149.353 \pm 2.973
CRAFT	-0.027 \pm 0.060	0.091 \pm 0.018	-	134.335 \pm 0.663
FAB	<u>-0.014\pm0.003</u>	<u>0.001\pm0.000</u>	-	153.894 \pm 3.916

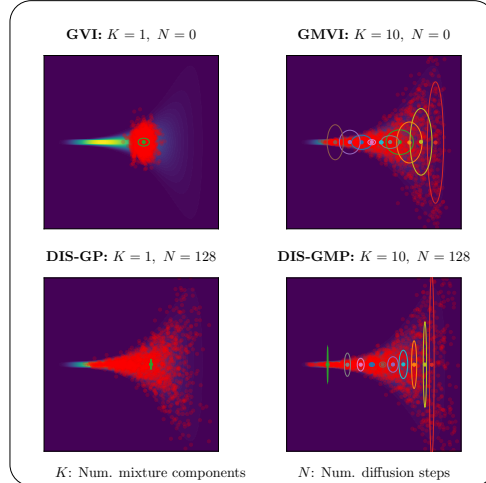


Figure 3: **Left side:** Results for Funnel target, averaged across four seeds. Evaluation criteria include evidence lower bound ELBO, importance-weighted errors for estimating the log-normalizing constant $\Delta \log Z$, effective sample size ESS, Sinkhorn distance \mathcal{W}_γ^2 . The best overall results are highlighted in bold, with category-specific best results underlined. Arrows (\uparrow , \downarrow) indicate whether higher or lower values are preferable, respectively. **Blue** and **green** shading indicate that the method uses learned Gaussian (GP) and Gaussian mixture priors (GMP), respectively. **Red** shading indicate competing state-of-the-art methods. Note that ESS cannot be computed due to the use of resampling schemes. **Right side:** Visualization of the first two dimensions of the Funnel target. Colored ellipses and circles denote standard deviations and means of the Gaussian components, respectively. Red dots illustrate samples of the model.

6.1 BENCHMARK PROBLEMS

We evaluate the different methods on various real-world and synthetic target densities.

Real-World Densities. We consider six real-world target densities: Four Bayesian inference tasks, where inference is performed over the parameters of a logistic regression model, namely *Credit* ($d = 25$), *Cancer* ($d = 31$), *Ionosphere* ($d = 35$), and *Sonar* ($d = 61$). Moreover, *Seeds* ($d = 26$) and *Brownian* ($d = 32$), where the goal is to perform inference over the parameters of a random effect regression model, and the time discretization of a Brownian motion, respectively. For these densities, we do not have access to the ground truth normalizer Z or samples from π preventing us from computing errors for log normalization estimation $\Delta \log Z$ and Sinkhorn distances \mathcal{W}_γ^2 . The resulting ELBO values are presented in Table 2.

Synthetic Densities. The *Funnel* density was introduced by Neal (2003) as has a shape that resembles a funnel, where one part is tight and highly concentrated, while the other is spread out over a wide region, making it challenging for sampling algorithms to explore the distribution effectively. Next, we consider the *Fashion* target which uses NICE (Dinh et al., 2014) to train a normalizing flow on the high-dimensional $d = 28 \times 28 = 784$ MNIST Fashion dataset. A recent study by Blessing et al. (2024) showed that current state-of-the-art methods were not able to generate samples with high quality from multiple modes.

6.2 RESULTS

Impact of Learned Gaussian (GP) and Gaussian Mixture (GMP) Priors. We evaluated the performance of our proposed methods on both real-world tasks and the *Funnel* density, employing $N = 128$ diffusion steps across all methods and $K = 10$ mixture components for X-GMP. To ensure a fair comparison, we initialized the priors of all diffusion-based methods with zero mean

METHOD	CREDIT	SEEDS	CANCER	BROWNIAN	IONOSPHERE	SONAR
GVI	-605.561±0.166	-76.741±0.007	-147.453±0.144	-3.885±0.005	-123.391±0.013	-137.696±0.043
GMVI	<u>-603.424±0.154</u>	<u>-75.221±0.011</u>	<u>-145.456±0.254</u>	<u>-2.250±0.011</u>	<u>-122.019±0.019</u>	<u>-135.959±0.031</u>
MCD	-1399.241±497.114	-75.699±0.015	-148.471±8.565	-15.498±0.158	-114.320±0.007	-112.639±0.025
MCD-GP	-585.350±0.015	-73.542±0.003	-89.676±0.189	0.771±0.008	-111.897±0.004	-109.338±0.004
MCD-GMP	<u>-585.276±0.013</u>	<u>-73.461±0.004</u>	<u>-88.562±0.243</u>	<u>0.993±0.003</u>	<u>-111.827±0.007</u>	<u>-109.197±0.004</u>
CMCD	-586.956±0.018	-74.033±0.010	-80.076±0.118	-1.346±0.013	-112.183±0.006	-109.332±0.006
CMCD-GP	-585.178±0.013	-73.456±0.003	-78.576±0.068	1.043±0.005	-111.687±0.003	-108.669±0.007
CMCD-GMP	<u>-585.162±0.002</u>	<u>-73.429±0.002</u>	<u>-78.402±0.037</u>	<u>1.087±0.001</u>	<u>-111.682±0.000</u>	<u>-108.634±0.000</u>
DIS	-589.636±0.757	-74.400±0.007	-86.592±2.107	-3.503±0.019	-112.525±0.008	-110.153±0.022
DIS-GP	-585.247±0.009	-73.540±0.005	-85.005±1.286	0.588±0.013	-111.847±0.006	-109.280±0.024
DIS-GMP	<u>-585.223±0.006</u>	<u>-73.492±0.003</u>	<u>-84.061±2.117</u>	<u>0.885±0.005</u>	<u>-111.811±0.002</u>	<u>-109.157±0.000</u>
DBS	-587.366±0.683	-73.918±0.008	-82.466±4.090	-0.773±0.010	-112.070±0.005	-109.188±0.005
DBS-GP	-585.524±0.414	-73.437±0.001	-83.395±4.184	1.081±0.004	-111.673±0.002	-108.595±0.006
DBS-GMP	<u>-585.148±0.002</u>	<u>-73.418±0.001</u>	<u>-78.160±0.063</u>	<u>1.118±0.002</u>	<u>-111.657±0.002</u>	<u>-108.548±0.000</u>
SMC	-698.403±4.146	-74.699±0.100	-194.059±0.613	-1.874±0.622	-114.751±0.238	-111.355±1.177
CRAFT	-594.795±0.411	-73.793±0.015	-95.737±1.067	0.886±0.053	-112.386±0.182	-115.618±1.316
FAB	<u>-585.102±0.001</u>	<u>-73.418±0.002</u>	-78.287±0.835	1.031±0.010	-111.678±0.003	-108.593±0.008

Table 2: Evidence lower bound (ELBO) values for various real-world benchmark problems, averaged across four seeds. The best overall results are highlighted in bold, with category-specific best results underlined. Blue and green shading indicate that the method uses learned Gaussian (GP) and Gaussian mixture priors (GMP), respectively. Red shading indicate competing state-of-the-art methods.

and unit variance. Table 2 and Figure 3 present our findings. The analysis demonstrates that GP consistently achieves tighter ELBO values compared to fixed priors, with GMP yielding further improvements over GP. Furthermore, Figure 3 illustrates both qualitatively and quantitatively that GMP effectively combines the strengths of Gaussian mixture and diffusion models, resulting in significant improvements. Specifically, we observed that the Gaussian components adapt well to the target’s support, covering both the neck and opening of the funnel shape. This results in less non-linear dynamics and better target coverage for DIS-GMP compared to using a single Gaussian (DIS-GP). Notably, the combination of DBS and GMP outperforms state-of-the-art methods across the majority of tasks and evaluation metrics.

Method	Fashion ($d = 784$)			
	ELBO \uparrow	$\Delta \log Z \downarrow$	$\mathcal{W}_2^2 \downarrow$	EMC \uparrow
GVI	-73.793±0.032	47.868±0.767	1590.212±0.818	0.000±0.000
GMVI	-72.654±0.176	45.927±0.380	1505.656±1.644	0.034±0.011
GMVI + IMR	<u>-57.021±0.052</u>	<u>30.444±0.571</u>	<u>589.881±1.374</u>	<u>0.761±0.039</u>
DIS	-41.63k±35.37	32.65k±390.6	17.72k±54.21	0.213±0.026
DIS-GP	<u>-24.712±0.253</u>	<u>10.581±0.496</u>	1671.411±2.394	0.007±0.004
DIS-GMP	-38.873±0.175	18.056±0.508	1703.023±3.050	0.012±0.021
DIS-GMP + IMR	-62.482±2.752	27.645±3.118	<u>513.776±13.936</u>	<u>0.780±0.089</u>
SMC	-12.18k±134.6	11.74k±139.2	6696.287±250.4	0.026±0.027
CRAFT	-520.47±5.531	445.10±8.273	1413.303±11.20	0.016±0.027
FAB	-892.97±151.8	350.54±599.0	1186.967±263.4	0.349±0.137

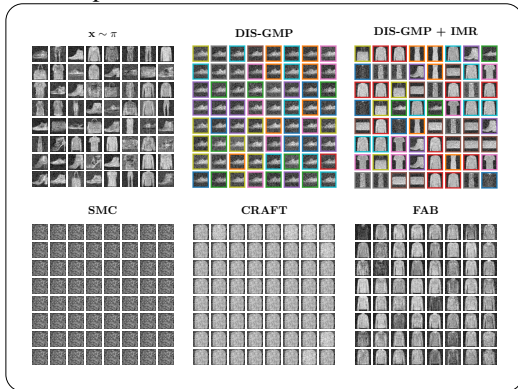


Figure 4: Left side: Results for Fashion target, averaged across four seeds. Evaluation criteria include evidence lower bound ELBO, importance-weighted errors for estimating the log-normalizing constant $\Delta \log Z$, and Sinkhorn distance \mathcal{W}_2^2 . The best overall results are highlighted in bold, with category-specific best results underlined. Arrows (\uparrow , \downarrow) indicate whether higher or lower values are preferable, respectively. Blue and green shading indicate that the method uses learned Gaussian (GP) and Gaussian mixture priors (GMP), respectively. Orange shading indicates that the method uses iterative model refinement (IMR). Red shading indicate competing state-of-the-art methods. Right side: Visualization of the $d = 28 \times 28 = 784$ dimensional Fashion samples. Top left corner visualizes samples from the target distribution. Colored frames indicate samples from different components of the Gaussian mixture.

Ablation Study: Number of Mixture Components K and Diffusion Steps N . We further investigated the effect of varying the number of diffusion steps N and mixture components K on a subset of tasks for DIS. The results, shown in Figure 5, demonstrate consistent improvements in effective sample size (ESS) with increases in both K and N . Additionally, we consistently observed that the combination of a higher number of components and diffusion steps yields the best overall performance. These trends hold across other metrics, as further detailed in Appendix D.

Iterative Model Refinement (IMR). Lastly, we investigated the impact of IMR, as detailed in Section 5, using DIS. For this analysis, we focused on the multi-modal Fashion target, which necessitates exploration in a high-dimensional space ($d = 784$). In addition to the performance criteria outlined in Section 6, we quantify how many of the modes the model discovered via the *entropic mode coverage (EMC)* introduced by Blessing et al. (2024). EMC evaluates the mode coverage of a sampler by leveraging prior knowledge of the target density. It holds that $EMC \in [0, 1]$ where $EMC = 1$ indicates that the model achieves uniform coverage over all modes whereas $EMC = 0$ indicates that the model only produces samples from a single mode. We employed the Metropolis-adjusted Langevin algorithm (MALA) (Cheng et al., 2018) to generate a set of candidate samples, noting that the computational cost of this process is comparable to a single gradient step in most diffusion-based sampling methods. The initial candidate samples as well as the support of DIS without learned prior are initialized such that they roughly cover the target support. Additional details are provided in Appendix C.2. We iteratively increased the number of components to $K = 10$, utilizing $N = 128$ diffusion steps throughout. Figure 4 presents our findings, demonstrating that the absence of IMR leads to mode collapse across all methods, as evidenced by high Sinkhorn distance values. The qualitative results highlight the role of candidate samples in facilitating mode discovery. Notably, the color-coding of DIS-GMP + IMR illustrates that each mixture component concentrates on a distinct mode, validating the effectiveness of the initialization heuristic proposed in Eq. 22 in balancing exploration and exploitation. This finding is also quantitatively reflected by the high EMC and low Sinkhorn distance values. In contrast, the ELBO and $\Delta \log Z$ values are slightly worse when using GMPs and IMR. This is attributed to the fact that these performance criteria are not well-suited for quantifying the model performance for multi-modal targets and tend to favor models that fit a single mode perfectly (Blessing et al., 2024). Moreover, with higher K , the diffusion model has to learn more complex control functions, as it needs to operate over the support of the entire Gaussian Mixture Model (GMM) rather than a single Gaussian. This added complexity can introduce more opportunities for approximation errors, which may negatively impact ELBO and $\Delta \log Z$ values compared to using a single learnable Gaussian. Nevertheless, the resulting samples from DIS-GMP are closer to the target distribution in terms of optimal transport (as indicated by the Sinkhorn distance). Importantly, these errors remain significantly smaller than those observed with non-learnable priors.

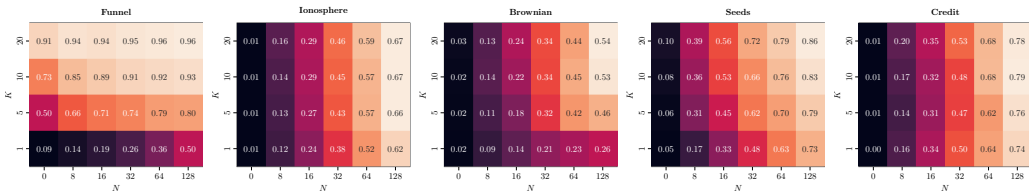


Figure 5: Effective sample size (ESS) of DIS-GMP for various real-world benchmark problems, averaged across four seeds. Here, N denotes the number of discretization steps and K the number of components in den Gaussian mixture.

7 CONCLUSION AND FUTURE WORK

In this paper, we propose a novel approach for improving diffusion-based sampling techniques by introducing end-to-end learnable Gaussian Mixture Priors (GMPs). Our method addresses key challenges in diffusion models—namely, non-linear drifts, mode collapse, and poor exploration—by providing more expressive and adaptable priors compared to the conventional Gaussian priors. We conducted comprehensive experiments on both synthetic and real-world datasets, which consistently demonstrated the superior performance of our proposed method. The results underscore the effectiveness of GMPs in overcoming the limitations of traditional diffusion models while requiring little to no hyperparameter tuning. Furthermore, we developed a novel strategy for iterative model refinement, which involves progressively adding components to the mixture during training, and demonstrated its effectiveness on a challenging high-dimensional problem.

A promising direction for future research is the improvement of the iterative model refinement strategy. While we showed that progressively increasing the number of components in the Gaussian mixture improves performance, optimizing the selection criteria for adding new components, generating better candidate samples, or dynamically adjusting the number of components during training, could lead to further gains in efficiency and accuracy.

REFERENCES

- 540
541
542 Tara Akhound-Sadegh, Jarrid Rector-Brooks, Avishek Joey Bose, Sarthak Mittal, Pablo Lemos,
543 Cheng-Hao Liu, Marcin Sendera, Siamak Ravanbakhsh, Gauthier Gidel, Yoshua Bengio, et al.
544 Iterated denoising energy matching for sampling from boltzmann densities. *arXiv preprint*
545 *arXiv:2402.06121*, 2024.
- 546 Michael Arbel, Alex Matthews, and Arnaud Doucet. Annealed flow transport monte carlo. In
547 *International Conference on Machine Learning*, pp. 318–330. PMLR, 2021.
- 548 Oleg Arenz, Gerhard Neumann, and Mingjun Zhong. Efficient gradient-free variational inference
549 using policy search. In *International conference on machine learning*, pp. 234–243. PMLR, 2018.
- 550 Oleg Arenz, Philipp Dahlinger, Zihan Ye, Michael Volpp, and Gerhard Neumann. A unified per-
551 spective on natural gradient variational inference with gaussian mixture models. *arXiv preprint*
552 *arXiv:2209.11533*, 2022.
- 553 Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In
554 *Proceedings of the 26th annual international conference on machine learning*, pp. 41–48, 2009.
- 555 Julius Berner, Lorenz Richter, and Karen Ullrich. An optimal control perspective on diffusion-based
556 generative modeling. *arXiv preprint arXiv:2211.01364*, 2022.
- 557 David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisti-
558 cians. *Journal of the American statistical Association*, 112(518):859–877, 2017.
- 559 Denis Blessing, Xiaogang Jia, Johannes Esslinger, Francisco Vargas, and Gerhard Neumann. Be-
560 yond elbows: A large-scale evaluation of variational methods for sampling. *arXiv preprint*
561 *arXiv:2406.07423*, 2024.
- 562 James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal
563 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, et al. Jax:
564 Autograd and xla. *Astrophysics Source Code Library*, pp. ascl-2111, 2021.
- 565 Wenlin Chen, Mingtian Zhang, Brooks Paige, José Miguel Hernández-Lobato, and David Barber.
566 Diffusive gibbs sampling. *arXiv preprint arXiv:2402.03008*, 2024.
- 567 Xiang Cheng, Niladri S Chatterji, Peter L Bartlett, and Michael I Jordan. Underdamped langevin
568 mcmc: A non-asymptotic analysis. In *Conference on learning theory*, pp. 300–323. PMLR, 2018.
- 569 Thomas M Cover. *Elements of information theory*. John Wiley & Sons, 1999.
- 570 Zac Cranko and Richard Nock. Boosted density estimation remastered. In *International Conference*
571 *on Machine Learning*, pp. 1416–1425. PMLR, 2019.
- 572 Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. *Advances in neural*
573 *information processing systems*, 26, 2013.
- 574 Marco Cuturi, Laetitia Meng-Papaxanthos, Yingtao Tian, Charlotte Bunne, Geoff Davis, and Olivier
575 Teboul. Optimal transport tools (ott): A jax toolbox for all things wasserstein. *arXiv preprint*
576 *arXiv:2201.12324*, 2022.
- 577 Paolo Dai Pra. A stochastic control approach to reciprocal diffusion processes. *Applied mathematics*
578 *and Optimization*, 23(1):313–329, 1991.
- 579 Pierre Del Moral, Arnaud Doucet, and Ajay Jasra. Sequential Monte Carlo samplers. *Journal of the*
580 *Royal Statistical Society: Series B (Statistical Methodology)*, 68(3):411–436, 2006.
- 581 Adji Bousso Dieng, Dustin Tran, Rajesh Ranganath, John Paisley, and David Blei. Variational
582 inference via χ upper bound minimization. *Advances in Neural Information Processing Systems*,
583 30, 2017.
- 584 Laurent Dinh, David Krueger, and Yoshua Bengio. Nice: Non-linear independent components esti-
585 mation. *arXiv preprint arXiv:1410.8516*, 2014.
- 586
587
588
589
590
591
592
593

- 594 Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp. *arXiv*
595 *preprint arXiv:1605.08803*, 2016.
- 596
- 597 Carles Domingo-Enrich. A taxonomy of loss functions for stochastic optimal control. *arXiv preprint*
598 *arXiv:2410.00345*, 2024.
- 599
- 600 Arnaud Doucet, Will Grathwohl, Alexander G Matthews, and Heiko Strathmann. Score-based diffu-
601 sion meets annealed importance sampling. *Advances in Neural Information Processing Systems*,
35:21482–21494, 2022a.
- 602
- 603 Arnaud Doucet, Will Grathwohl, Alexander G de G Matthews, and Heiko Strathmann. Score-based
604 diffusion meets annealed importance sampling. In *Advances in Neural Information Processing*
605 *Systems*, 2022b.
- 606
- 607 Hans Föllmer. An entropy approach to the time reversal of diffusion processes. In *Stochastic*
608 *Differential Systems Filtering and Control: Proceedings of the IFIP-WG 7/1 Working Conference*
609 *Marseille-Luminy, France, March 12–17, 1984*, pp. 156–163. Springer, 2005.
- 610
- 611 Daan Frenkel and Berend Smit. *Understanding molecular simulation: from algorithms to applica-*
612 *tions*. Elsevier, 2023.
- 613
- 614 Tomas Geffner and Justin Domke. MCMC variational inference via uncorrected Hamiltonian an-
615 nealing. In *Advances in Neural Information Processing Systems*, 2021.
- 616
- 617 Tomas Geffner and Justin Domke. Langevin diffusion variational inference. *arXiv preprint*
618 *arXiv:2208.07743*, 2022.
- 619
- 620 Fangjian Guo, Xiangyu Wang, Kai Fan, Tamara Broderick, and David B Dunson. Boosting varia-
621 tional inference. *arXiv preprint arXiv:1611.05559*, 2016.
- 622
- 623 John Hammersley. *Monte carlo methods*. Springer Science & Business Media, 2013.
- 624
- 625 Carsten Hartmann, Omar Kebiri, Lara Neureither, and Lorenz Richter. Variational approach to rare
626 event simulation using least-squares regression. *Chaos: An Interdisciplinary Journal of Nonlinear*
627 *Science*, 29(6), 2019.
- 628
- 629 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
630 *neural information processing systems*, 33:6840–6851, 2020.
- 631
- 632 Robert E Kass, Bradley P Carlin, Andrew Gelman, and Radford M Neal. Markov chain monte carlo
633 in practice: a roundtable discussion. *The American Statistician*, 52(2):93–100, 1998.
- 634
- 635 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
636 *arXiv:1412.6980*, 2014.
- 637
- 638 Xuechen Li, Ting-Kam Leonard Wong, Ricky TQ Chen, and David Duvenaud. Scalable gradients
639 for stochastic differential equations. In *International Conference on Artificial Intelligence and*
640 *Statistics*, pp. 3870–3882. PMLR, 2020.
- 641
- 642 Yingzhen Li and Richard E Turner. Rényi divergence variational inference. *Advances in neural*
643 *information processing systems*, 29, 2016.
- 644
- 645 Jun S Liu and Jun S Liu. *Monte Carlo strategies in scientific computing*, volume 75. Springer, 2001.
- 646
- 647 Jun S Liu, Rong Chen, and Tanya Logvinenko. A theoretical framework for sequential importance
sampling with resampling. In *Sequential Monte Carlo methods in practice*, pp. 225–246. Springer,
2001.
- 648
- 649 Xingchao Liu, Lemeng Wu, Mao Ye, and Qiang Liu. Let us build bridges: Understanding and
650 extending diffusion generative models. *arXiv preprint arXiv:2208.14699*, 2022.
- 651
- 652 Alex Matthews, Michael Arbel, Danilo Jimenez Rezende, and Arnaud Doucet. Continual repeated
annealed flow transport monte carlo. In *International Conference on Machine Learning*, pp.
15196–15219. PMLR, 2022.

- 648 Laurence Illing Midgley, Vincent Stimper, Gregor NC Simm, Bernhard Schölkopf, and
649 José Miguel Hernández-Lobato. Flow annealed importance sampling bootstrap. *arXiv preprint*
650 *arXiv:2208.01893*, 2022.
- 651 Andrew C Miller, Nicholas J Foti, and Ryan P Adams. Variational boosting: Iteratively refining
652 posterior approximations. In *International Conference on Machine Learning*, pp. 2420–2429.
653 PMLR, 2017.
- 654 Christian Naeseth, Scott Linderman, Rajesh Ranganath, and David Blei. Variational sequential
655 monte carlo. In *International conference on artificial intelligence and statistics*, pp. 968–977.
656 PMLR, 2018.
- 657 Christian Naeseth, Fredrik Lindsten, and David Blei. Markovian score climbing: Variational infer-
658 ence with kl (p— q). *Advances in Neural Information Processing Systems*, 33:15499–15510,
659 2020.
- 660 Radford M Neal. Annealed importance sampling. *Statistics and Computing*, 11(2):125–139, 2001.
- 661 Radford M Neal. Slice sampling. *The annals of statistics*, 31(3):705–767, 2003.
- 662 Edward Nelson. *Dynamical theories of Brownian motion*, volume 101. Princeton university press,
663 2020.
- 664 Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models.
665 In *International conference on machine learning*, pp. 8162–8171. PMLR, 2021.
- 666 Kim A Nicoli, Christopher J Anders, Lena Funcke, Tobias Hartung, Karl Jansen, Pan Kessel,
667 Shinichi Nakajima, and Paolo Stornati. Estimation of thermodynamic observables in lattice field
668 theories with deep generative models. *Physical review letters*, 126(3):032001, 2021.
- 669 George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji Lak-
670 shminarayanan. Normalizing flows for probabilistic modeling and inference. *The Journal of*
671 *Machine Learning Research*, 22(1):2617–2680, 2021.
- 672 Angus Phillips, Hai-Dang Dau, Michael John Hutchinson, Valentin De Bortoli, George Deligianni-
673 dis, and Arnaud Doucet. Particle denoising diffusion sampler. *arXiv preprint arXiv:2402.06320*,
674 2024.
- 675 Maziar Raissi, Paris Perdikaris, and George E Karniadakis. Physics-informed neural networks: A
676 deep learning framework for solving forward and inverse problems involving nonlinear partial
677 differential equations. *Journal of Computational physics*, 378:686–707, 2019.
- 678 Lorenz Richter, Ayman Boustati, Nikolas Nüsken, Francisco Ruiz, and Omer Deniz Akyildiz. Var-
679 grad: a low-variance gradient estimator for variational inference. *Advances in Neural Information*
680 *Processing Systems*, 33:13481–13492, 2020.
- 681 Lorenz Richter, Julius Berner, and Guan-Hong Liu. Improved sampling via learned diffusions.
682 *arXiv preprint arXiv:2307.01198*, 2023.
- 683 Gareth O Roberts and Richard L Tweedie. Exponential convergence of langevin distributions and
684 their discrete approximations. 1996.
- 685 Simo Särkkä and Arno Solin. *Applied stochastic differential equations*, volume 10. Cambridge
686 University Press, 2019.
- 687 Yuyang Shi, Valentin De Bortoli, Andrew Campbell, and Arnaud Doucet. Diffusion schrödinger
688 bridge matching. *Advances in Neural Information Processing Systems*, 36, 2024.
- 689 Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution.
690 *Advances in neural information processing systems*, 32, 2019.
- 691 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
692 Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint*
693 *arXiv:2011.13456*, 2020.

- 702 Gabriel Stoltz, Mathias Rousset, et al. *Free energy computations: A mathematical perspective*.
703 World Scientific, 2010.
704
- 705 Jingtong Sun, Julius Berner, Lorenz Richter, Marius Zeinhofer, Johannes Müller, Kamyar Aziz-
706 zadenesheli, and Anima Anandkumar. Dynamical measure transport and neural pde solvers for
707 sampling. *arXiv preprint arXiv:2407.07873*, 2024.
- 708 Achille Thin, Nikita Kotelevskii, Alain Durmus, Eric Moulines, Maxim Panov, and Arnaud Doucet.
709 Monte Carlo variational auto-encoders. In *International Conference on Machine Learning*, 2021.
710
- 711 Francisco Vargas, Will Grathwohl, and Arnaud Doucet. Denoising diffusion samplers. *arXiv*
712 *preprint arXiv:2302.13834*, 2023a.
- 713 Francisco Vargas, Andrius Ovsianas, David Fernandes, Mark Girolami, Neil D Lawrence, and Niko-
714 las Nüsken. Bayesian learning via neural schrödinger–föllmer flows. *Statistics and Computing*,
715 33(1):3, 2023b.
- 716 Francisco Vargas, Shreyas Padhy, Denis Blessing, and Nikolas Nüsken. Transport meets variational
717 inference: Controlled monte carlo diffusions. In *The Twelfth International Conference on Learn-*
718 *ing Representations*, 2024.
719
- 720 Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural compu-*
721 *tation*, 23(7):1661–1674, 2011.
722
- 723 Martin J. Wainwright and Michael I. Jordan. Graphical Models, Exponential Families, and Varia-
724 tional Inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, November 2008.
- 725 Neng Wan, Dapeng Li, and Naira Hovakimyan. F-divergence variational inference. *Advances in*
726 *neural information processing systems*, 33:17370–17379, 2020.
727
- 728 Hao Wu, Jonas Köhler, and Frank Noé. Stochastic normalizing flows. *Advances in Neural Informa-*
729 *tion Processing Systems*, 33:5933–5944, 2020.
- 730 Dinghuai Zhang, Ricky Tian Qi Chen, Cheng-Hao Liu, Aaron Courville, and Yoshua Bengio. Diffu-
731 sion generative flow samplers: Improving learning signals through partial trajectory optimization.
732 *arXiv preprint arXiv:2310.02679*, 2023.
- 733 Guodong Zhang, Kyle Hsu, Jianing Li, Chelsea Finn, and Roger Grosse. Differentiable annealed im-
734 portance sampling and the perils of gradient noise. In *Advances in Neural Information Processing*
735 *Systems*, 2021.
736
- 737 Qinsheng Zhang and Yongxin Chen. Path integral sampler: a stochastic control approach for sam-
738 pling. *arXiv preprint arXiv:2111.15141*, 2021.
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

A PROOFS

A.1 PROOF OF PROPOSITION 1

We will use the Fokker-Planck equation (FPE) and show that the given stationary distribution satisfies it when the time derivative is set to zero:

First, recall the FPE for the probability density $p(\mathbf{x}, t)$ of a process described by the stochastic differential equation (SDE)

$$d\mathbf{x}_t = -f(\mathbf{x}_t)dt + \sigma d\mathbf{w}_t \quad (23)$$

is given as

$$\frac{\partial p(\mathbf{x}, t)}{\partial t} = \nabla \cdot [f(\mathbf{x})p(\mathbf{x}, t)] + \frac{\sigma^2}{2} \nabla^2 p(\mathbf{x}, t). \quad (24)$$

For the stationary distribution $p_s(\mathbf{x})$, we set $\frac{\partial p(\mathbf{x}, t)}{\partial t} = 0$:

$$0 = \nabla \cdot [f(\mathbf{x})p_s(\mathbf{x})] + \frac{\sigma^2}{2} \nabla^2 p_s(\mathbf{x}) \quad (25)$$

Next, recall the proposed stationary distribution:

$$p_s(\mathbf{x}) \propto \exp\left(-\frac{2}{\sigma^2} \int f(\mathbf{x})d\mathbf{x}\right) \quad (26)$$

Next, we verify that this satisfies stationary FPE (Eq. 25). First, let's compute the gradient and Laplacian of $p_s(\mathbf{x})$:

$$\nabla p_s(\mathbf{x}) = p_s(\mathbf{x}) \cdot \left(-\frac{2}{\sigma^2}\right) f(\mathbf{x}) \quad (27)$$

$$\nabla^2 p_s(\mathbf{x}) = \nabla \cdot [p_s(\mathbf{x}) \cdot \left(-\frac{2}{\sigma^2}\right) f(\mathbf{x})] \quad (28)$$

$$= p_s(\mathbf{x}) \cdot \left(-\frac{2}{\sigma^2}\right)^2 [f(\mathbf{x})]^2 + p_s(\mathbf{x}) \cdot \left(-\frac{2}{\sigma^2}\right) \nabla \cdot f(\mathbf{x}) \quad (29)$$

Finally, we substitute these into the left side of Eq. 25, that is,

$$\nabla \cdot [f(\mathbf{x})p_s(\mathbf{x})] + \frac{\sigma^2}{2} \nabla^2 p_s(\mathbf{x}) \quad (30)$$

$$= \nabla \cdot [f(\mathbf{x})p_s(\mathbf{x})] + \frac{\sigma^2}{2} \left[p_s(\mathbf{x}) \cdot \left(-\frac{2}{\sigma^2}\right)^2 [f(\mathbf{x})]^2 + p_s(\mathbf{x}) \cdot \left(-\frac{2}{\sigma^2}\right) \nabla \cdot f(\mathbf{x}) \right] \quad (31)$$

$$= p_s(\mathbf{x}) \nabla \cdot f(\mathbf{x}) + f(\mathbf{x}) \nabla p_s(\mathbf{x}) + p_s(\mathbf{x}) \left[-\frac{2}{\sigma^2} \right] [f(\mathbf{x})]^2 - p_s(\mathbf{x}) \nabla \cdot f(\mathbf{x}) \quad (32)$$

$$= p_s(\mathbf{x}) \nabla \cdot f(\mathbf{x}) + f(\mathbf{x}) p_s(\mathbf{x}) \left(-\frac{2}{\sigma^2}\right) f(\mathbf{x}) + p_s(\mathbf{x}) \left[-\frac{2}{\sigma^2} \right] [f(\mathbf{x})]^2 - p_s(\mathbf{x}) \nabla \cdot f(\mathbf{x}) \quad (33)$$

$$= 0, \quad (34)$$

which yields the desired result. \square

B ADDITIONAL DETAILS FOR DIFFUSION-BASED SAMPLER

Pseudocode: We additionally provide pseudocode in Algorithm 1 for a generic diffusion sampler with learnable prior p_0^ϕ . For clarity, we present an update step for a single sample. In practice, however, one would use mini-batches for these updates.

Special Cases of X-GMP: Consider the generic (extended) ELBO for X-GMP, that is,

$$\mathcal{L}_{\text{GMP}}(\theta, \gamma, \phi, \delta t) = \mathbb{E}_{\mathbf{x}_{0:N} \sim \mathcal{P}^{\theta, \phi, \delta t}} \left[\log \frac{\rho(\mathbf{x}_N)}{\sum_{k=1}^K \alpha_k p_0^{\phi_k}(\mathbf{x}_0)} + \sum_{n=1}^N \log \frac{B_{n-1}^{\gamma, \phi, \delta t}(\mathbf{x}_{n-1} | \mathbf{x}_n)}{F_n^{\theta, \phi, \delta t}(\mathbf{x}_n | \mathbf{x}_{n-1})} \right], \quad (35)$$

with $p_0^{\phi_k}(\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_0 | \mu_k, \Sigma_k)$. We obtain the following special cases:

Algorithm 1 Training of diffusion sampler with learnable prior**Require:**

- control functions u_θ, v_γ with initial parameters θ_0, γ_0
- prior distribution p_0^ϕ with initial parameters ϕ_0
- initial step size δt_0
- number of gradient steps G , number of diffusion steps N , step size η

 $\Theta_0 = \{\theta_0, \gamma_0, \phi_0, \delta t_0\}$
for $i \leftarrow 0, \dots, G - 1$ **do**
 $\mathbf{x}_0 \leftarrow g(\xi, \phi_i), \quad \xi \sim p(\cdot) \quad \triangleright$ sample p_i^ϕ via reparameterization (batched in practice)

 $\mathcal{L} \leftarrow \log p^{\phi_i}(\mathbf{x}_0)$
for $n \leftarrow 0, \dots, N - 1$ **do**
 $\mathbf{x}_{n+1} = \mathbf{x}_n + [f(\mathbf{x}_n, n) + \sigma u^{\theta_i}(\mathbf{x}_n, n)] \delta t_i + \sigma \sqrt{2\delta t_i} \epsilon_n$
 $\mathcal{L} \leftarrow \mathcal{L} + \log F_{n+1}^{\theta_i, \phi_i, \delta t_i}(\mathbf{x}_{n+1} | \mathbf{x}_n) - \log B_n^{\gamma_i, \phi_i, \delta t_i}(\mathbf{x}_n | \mathbf{x}_{n+1})$
 $\mathcal{L} \leftarrow \mathcal{L} - \log \rho(\mathbf{x}_N)$
 $\Theta_{i+1} \leftarrow \Theta_i + \eta \nabla_{\Theta} \mathcal{L}$
 \triangleright maximize (extended) ELBO

return optimized parameters Θ_G

- **GVI** ($K = 1, N = 0$): For a single Gaussian mixture component and zero diffusion steps, Equation (35) reduces to the marginal ELBO objective in Equation (9) for a Gaussian distribution, that is,

$$\mathcal{L}_{\text{GVI}}(\phi) = \mathbb{E}_{\mathbf{x}_0 \sim p_0^\phi} \left[\log \frac{\rho(\mathbf{x}_0)}{p_0^\phi(\mathbf{x}_0)} \right]. \quad (36)$$

- **GVI** ($K > 1, N = 0$): Similarly, if we have zero diffusion steps, but multiple Gaussian mixture components we obtain the marginal ELBO for Gaussian mixture models, i.e.,

$$\mathcal{L}_{\text{GMVI}}(\phi) = \mathbb{E}_{\mathbf{x}_0 \sim p_0^\phi} \left[\log \frac{\rho(\mathbf{x}_0)}{\sum_{k=1}^K \alpha_k p_0^{\phi_k}(\mathbf{x}_0)} \right]. \quad (37)$$

Please note that there are more sophisticated methods to train Gaussian mixture models for VI, see [Arenz et al. \(2018; 2022\)](#).

- **X-GP** ($K = 1, N > 0$): For a single mixture component and multiple diffusion steps, we obtain the objective for X-GP, i.e., for a diffusion-model with learned Gaussian prior, given by

$$\mathcal{L}_{\text{GP}}(\theta, \gamma, \phi, \delta t) = \mathbb{E}_{\mathbf{x}_0: N \sim p^{\theta, \phi, \delta t}} \left[\log \frac{\rho(\mathbf{x}_N)}{p_0^\phi(\mathbf{x}_0)} + \sum_{n=1}^N \log \frac{B_{n-1}^{\gamma, \phi, \delta t}(\mathbf{x}_{n-1} | \mathbf{x}_n)}{F_n^{\theta, \phi, \delta t}(\mathbf{x}_n | \mathbf{x}_{n-1})} \right]. \quad (38)$$

- **X-GMP** ($K > 1, N > 0$): Having multiple multiple mixture components K and diffusion steps N results in the full X-GMP objective, as in Equation (35).

Forward and Backward Transitions. We provide further information about the diffusion-based sampling methods considered in this work in Table 3. Specifically, we provide expressions for the forward and backward transitions.

Time complexity. Diffusion-based samplers that use a Gaussian prior have a time complexity of $\mathcal{O}(N)$, whereas Gaussian Mixture Priors (GMPs) incur a time complexity of $\mathcal{O}(NK)$. The additional factor K arises from the need to compute the likelihood of the GMP at each diffusion step. However, in practice, the evaluation of the likelihood of the GMP can be parallelized across its components, which substantially reduces the computational overhead. This parallelization allows for efficient implementation despite the increased theoretical complexity.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

Method	$F_{n+1}^{\theta, \phi, \Delta}(\mathbf{x}_{n+1} \mathbf{x}_n)$	$B_{n-1}^{\gamma, \phi, \Delta}(\mathbf{x}_{n-1} \mathbf{x}_n)$
DIS	$\mathcal{N}(\mathbf{x}_{n+1} \mathbf{x}_n + [-\sigma^2 \nabla \log p_0^\phi(\mathbf{x}_n) + \sigma u^\theta(\mathbf{x}_n, n)] \delta t, 2\sigma^2 \delta t I)$	$\mathcal{N}(\mathbf{x}_{n-1} \mathbf{x}_n + \sigma^2 \nabla \log p_0^\phi(\mathbf{x}_n) \delta t, 2\sigma^2 \delta t I)$
MCD	$\mathcal{N}(\mathbf{x}_{n+1} \mathbf{x}_n + \sigma^2 \nabla_{\mathbf{x}} \log \pi_n^\phi(\mathbf{x}_n) \delta t, 2\sigma^2 \delta t I)$	$\mathcal{N}(\mathbf{x}_{n-1} \mathbf{x}_n - [\sigma^2 \nabla_{\mathbf{x}} \log \pi_n^\phi(\mathbf{x}_n) - \sigma v^\gamma(\mathbf{x}_n, n)] \delta t, 2\sigma^2 \delta t I)$
CMCD	$\mathcal{N}(\mathbf{x}_{n+1} \mathbf{x}_n + [\sigma^2 \nabla_{\mathbf{x}} \log \pi_n^\phi(\mathbf{x}_n) + \sigma u^\theta(\mathbf{x}_n, n)] \delta t, 2\sigma^2 \delta t I)$	$\mathcal{N}(\mathbf{x}_{n-1} \mathbf{x}_n - [\sigma^2 \nabla_{\mathbf{x}} \log \pi_n^\phi(\mathbf{x}_n) - \sigma u^\theta(\mathbf{x}_n, n)] \delta t, 2\sigma^2 \delta t I)$
DBS	$\mathcal{N}(\mathbf{x}_{n+1} \mathbf{x}_n + [f(\mathbf{x}_n, n) + \sigma u^\theta(\mathbf{x}_n, n)] \delta t, 2\sigma^2 \delta t I)$	$\mathcal{N}(\mathbf{x}_{n-1} \mathbf{x}_n - [f(\mathbf{x}_n, n) - \sigma v^\gamma(\mathbf{x}_n, n)] \delta t, 2\sigma^2 \delta t I)$

Table 3: Comparison of different forward and backward transitions $F^{\theta, \phi, \delta t}$, and $B^{\gamma, \phi, \delta t}$, respectively, for diffusion-based sampling methods based on f , π_n^ϕ , p_0^ϕ , u^θ and v^γ as defined in the text.

Memory consumption. When using the standard "discrete-then-optimize" approach to minimize the KL divergence in Eq. 11, which requires differentiation through the SDE, memory consumption scales linearly with both K (number of components) and N (number of diffusion steps). In contrast, methods like the stochastic adjoint approach for KL optimization (Li et al., 2020) achieve constant memory consumption, making them more suitable for scenarios with a large number of components or diffusion steps.

In our experiments, we opted for the former approach due to its simplicity. However, for tasks involving extensive components or steps, the stochastic adjoint method or similar approaches may be more practical.

Additionally, constant memory consumption can also be achieved by using alternative loss functions such as the log-variance loss (Richter et al., 2020; 2023) or moment-loss (Hartmann et al., 2019).

C EXPERIMENTAL DETAILS

C.1 BENCHMARKING TARGETS

This section introduces the target densities considered in our experiments. Please note that the majority of tasks are taken from the recent benchmark study from Blessing et al. (2024). For convenience, we provide a brief explanation of the target densities.

Bayesian Logistic Regression: We evaluate a Bayesian logistic regression model on four standardized binary classification datasets:

- **Ionosphere** ($d = 35, 351 (x_i, y_i)$ pairs)
- **Sonar** ($d = 61, 208 (x_i, y_i)$ pairs)
- **German Credit** ($d = 25, 1000 (x_i, y_i)$ pairs)
- **Breast Cancer** ($d = 31, 569 (x_i, y_i)$ pairs)

The model assumes:

$$\omega \sim \mathcal{N}(0, \sigma_\omega^2 I),$$

$$y_i \sim \text{Bernoulli}(\text{sigmoid}(\omega^\top x_i)),$$

where features are standardized for linear logistic regression. Here, we perform inference over the parameters ω of the (linear) logistic regression model. In Blessing et al. (2024), the authors used an uninformative prior for the parameters of the Bayesian logistic regression models for the *Credit* and *Cancer* tasks, which frequently caused numerical instabilities. To maintain the challenge of the tasks while ensuring stability, we opted for a Gaussian prior with zero mean and variance of $\sigma_\omega^2 = 100$.

918 **Random Effect Regression:** We apply random effect regression to the **Seeds** dataset ($d = 26$):

919
920 $\tau \sim \text{Gamma}(0.01, 0.01),$
921 $a_0, a_1, a_2, a_{12} \sim \mathcal{N}(0, 10),$
922 $b_i \sim \mathcal{N}(0, \frac{1}{\sqrt{\tau}}), \quad i = 1, \dots, 21,$
923 $\text{logits}_i = a_0 + a_1 x_i + a_2 y_i + a_{12} x_i y_i + b_1,$
924 $r_i \sim \text{Binomial}(\text{logits}_i, N_i),$

925
926
927 with inference conducted over model parameters given observed data.

928
929 **Time Series Models:** For time series analysis, we use the **Brownian** model ($d = 32$):

930
931 $\alpha_{\text{inn}} \sim \text{LogNormal}(0, 2),$
932 $\alpha_{\text{obs}} \sim \text{LogNormal}(0, 2),$
933 $x_1 \sim \mathcal{N}(0, \alpha_{\text{inn}}),$
934 $x_i \sim \mathcal{N}(x_{i-1}, \alpha_{\text{inn}}), \quad i = 2, \dots, 20,$
935 $y_i \sim \mathcal{N}(x_i, \alpha_{\text{obs}}), \quad i = 1, \dots, 30,$
936

937 with inference focusing on parameters $\alpha_{\text{inn}}, \alpha_{\text{obs}}$, and latent states $\{x_i\}_{i=1}^{30}$.

938
939 **Funnel:** ($d = 10$), a funnel-shaped distribution defined by:

940 $\pi(x) = \mathcal{N}(x_1; 0, \sigma_f^2) \mathcal{N}(x_{2:10}; 0, \exp(x_1)I),$

941
942 with $\sigma_f^2 = 9$.

943
944 **Fashion and Digits.** MNIST variants (**DIGITS**) and Fashion MNIST (**Fashion**) datasets using
945 NICE (Dinh et al., 2014) to train normalizing flows, with resolutions 14×14 and DIGITS and
946 28×28 for Fashion.

947
948 C.2 DIFFUSION-BASED METHODS: DETAILS AND TUNING

949
950 **General setting:** All experiments are conducted using the Jax library (Bradbury et al., 2021). Our
951 default experimental setup, unless specified otherwise, is as follows: We use a batch size of 2000
952 (halved if memory-constrained) and train for 140k gradient steps to ensure approximate conver-
953 gence. We use the Adam optimizer (Kingma & Ba, 2014), gradient clipping with a value of 1, and a
954 learning rate scheduler that starts at 8×10^{-3} and uses a cosine decay starting at 60k gradient steps.
955 We utilized 128 discretization steps and the Euler-Maruyama method for integration. The control
956 functions u^θ and v^γ were parameterized as two-layer neural networks with 128 neurons. **For DBS,**
957 **we set the drift to $f = \sigma^2 \nabla \log \pi$.**

958 Unlike Zhang & Chen (2021), we did not include the gradient of the target density in the network
959 architecture. Inspired by Nichol & Dhariwal (2021), we applied a cosine-square scheduler for the
960 discretization step size: $\delta t = a \cos^2(\frac{\pi}{2} \frac{n}{N})$, where $a : [0, \infty) \rightarrow (0, \infty)$ is learned for all methods.
961 We enforced non-negativity of a via an element-wise softplus transformation. The diffusion coeffi-
962 cient σ was set to 1 for all experiments. Furthermore, we set the initial a to 0.1 for all experiments
963 except Brownian, where we set 0.01. We did not perform any hyperparameter tuning since most
964 parameters are learned end-to-end.

965
966 **Gaussian Priors (GP) and Gaussian Mixture Priors (GMP):** We learn diagonal Gaussian priors
967 and ensure positive definiteness with an element-wise softplus transformation. We use a separate
968 learning rate of 10^{-2} for all experiments to allow for quick adaptation of the Gaussian components.
969 Furthermore, the mean was initialized at 0 and the initial covariance matrix was set to the identity
970 **except for Fashion where we set the initial variance to 5 which roughly covers the support of the**
971 **target.** The individual components in the Gaussian mixture follow the setup of Gaussian priors. The
mixture weights are uniformly initialized and fixed during training. If not otherwise specified, we
use $K = 10$ mixture components for X-GMP.

Methods / Parameters	Grid	Funnel	Fashion	Credit	Cancer	Brownian	Sonar	Seeds	Ionosphere	ϕ^4
SMC										
Initial Scale	{0.1, 1, 10}	1	5 [†]	0.1	1	1	1	1	1	5 [†]
HMC stepsize ($\beta \leq 0.5$)	{0.005, 0.001, 0.01, 0.05, 0.1, 0.2}	0.001	0.2	0.01	0.01	0.001	0.05	0.2	0.2	0.1
HMC stepsize ($\beta > 0.5$)	{0.005, 0.001, 0.01, 0.05, 0.1, 0.2}	0.1	0.2	0.005	0.005	0.05	0.001	0.05	0.2	0.05
CRAFT										
Initial Scale	{0.1, 1, 10}	1	5 [†]	1	1	1	1	0.1	0.1	
Learning Rate	$\{10^{-3}, 10^{-4}, 5 \times 10^{-4}, 10^{-5}\}$	10^{-3}	10^{-4}	5×10^{-4}	5×10^{-4}	10^{-3}	10^{-3}	10^{-3}	10^{-3}	
FAB										
Initial Scale	{0.1, 1, 10}	1	5 [†]	1	1	1	0.1	0.1	1	
Learning Rate	$\{10^{-3}, 10^{-4}, 5 \times 10^{-4}, 10^{-5}\}$	10^{-4}	10^{-3}	10^{-4}	10^{-3}	10^{-3}	10^{-3}	10^{-4}	10^{-3}	

Table 4: Hyperparameter selection for all different sampling algorithms. The ‘Grid’ column indicates the values over which we performed a grid search. The values in the column which are marked with experiment names indicate which values were chosen for the reported results. The values for parameters indicated with † are set by using prior knowledge about the task.

Iterative Model Refinement (IMR): For IMR, we add a new component after 500 training iterations starting with a single component. The initial means were selected with the heuristic presented in Equation (22). **The variance of the newly added components was set to be 1.** The candidate sample set was generated using the Metropolis Adjusted Langevin Algorithm (MALA) (Cheng et al., 2018). For that, we used 2000 random samples from a Gaussian with zero mean and variance 5, which roughly covers the support of the *Fashion* target. Please note that competing methods also use this prior knowledge for initialization of the prior, see Table 4. We use 128 steps steps, that is,

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \tilde{\sigma}^2 \nabla \log \pi(\mathbf{x}_i) \delta t + \tilde{\sigma} \sqrt{2\delta t} \epsilon, \quad \epsilon \sim \mathcal{N}(\cdot|0, I) \quad (39)$$

with $\tilde{\sigma} = 5$ and an additional Metropolis adjustment step. Here, $\tilde{\sigma}$ was chosen such that the final set of samples yields high target log-likelihoods $\log \rho(\mathbf{x})$. The final samples are used as candidate set. We note that this procedure brings the new components close to different modes in the target distribution and therefore facilitates exploration. Moreover, the computation of such a candidate set is very cheap, i.e., the equivalent of a single gradient step for e.g. MCD or CMCD.

C.3 COMPETING METHODS: DETAILS AND TUNING

The results for competing methods presented in this work are primarily drawn from Blessing et al. (2024), where hyperparameters were carefully optimized. For convenience, we repeat the details. Since our experimental setup differs for the *Credit* and *Cancer* tasks (detailed in Section C.1), we adhered to the tuning recommendations provided by Blessing et al. (2024). Details about hyperparameters can be found in Table 4.

Sequential Monte Carlo (SMC) and Continual Repeated Annealed Flow Transport (CRAFT): The Sequential Monte Carlo (SMC) approach was implemented with 2000 particles and 128 annealing steps, matching the number of sequential steps used in diffusion-based sampling methods. Resampling was performed with a threshold of 0.3, and one Hamiltonian Monte Carlo (HMC) step was applied per temperature, using 5 leapfrog steps. The HMC step size was tuned according to Table 4, with different step sizes based on the annealing parameter β_t . Additionally, the scale of the initial proposal distribution was tuned. As CRAFT builds on the SMC framework, it used the same SMC specifications, incorporating diagonal affine flows (Papamakarios et al., 2021) as transition models.

Flow Annealed Importance Sampling Bootstrap (FAB): Automatic step size tuning for the SMC sampler was applied on top of the normalizing flow (Papamakarios et al., 2021). The flow architecture utilized RealNVP (Dinh et al., 2016), with an 8-layer MLP serving as the conditioner. FAB’s replay buffer was employed to accelerate computations. The learning rate and base distribution scale were adjusted for target specificity as outlined in Table 4. A batch size of 2000 was used, and FAB was trained until reaching approximate convergence, which was sufficient to achieve approximate convergence.

C.4 EVALUATION

Evaluation protocol and model selection We follow the evaluation protocol of prior work (Blessing et al., 2024) and evaluate all performance criteria 100 times during training, using 2000 samples for each evaluation. To smooth out short-term fluctuations and obtain more robust results within a

single run, we apply a running average with a window of 5 evaluations. We conduct each experiment using four different random seeds and average the best results of each run.

Performance Criteria: In order to define the performance criteria, we first define the unnormalized (extended) importance weights \tilde{w} , that is,

$$\tilde{w} := \frac{\rho(\mathbf{x}_N) \prod_{n=1}^N B_{n-1}^{\gamma, \phi, \delta t}(\mathbf{x}_{n-1} | \mathbf{x}_n)}{p_0^\phi(\mathbf{x}_0) \prod_{n=1}^N F_n^{\theta, \phi, \delta t}(\mathbf{x}_n | \mathbf{x}_{n-1})}. \quad (40)$$

We consider the following following performance criteria:

- **Evidence lower bound (ELBO):** We compute the (extended) ELBO as

$$\text{ELBO} := \mathbb{E}_{\mathbf{x}_{0:N} \sim \mathcal{P}^{\theta, \phi, \delta t}} [\log \tilde{w}] \approx \frac{1}{m} \sum_{i=1}^m \log \tilde{w}^{(i)}. \quad (41)$$

- **Evidence upper bound (EUBO):** We compute the (extended) EUBO as

$$\text{EUBO} := \mathbb{E}_{\mathbf{x}_{0:N} \sim \mathcal{Q}^{\gamma, \phi, \delta t}} [\log \tilde{w}] \approx \frac{1}{m} \sum_{i=1}^m \log \tilde{w}^{(i)}. \quad (42)$$

Please note that we need samples from the target, i.e., $\mathbf{x}_T \sim \pi$ to compute the expectation in Eq. 42 by simulating the backward process $\tilde{\mathbf{X}}$. Moreover, it is straightforward to see that the EUBO serves as an upper bound on the log normalization constant since

$$D_{\text{KL}}(\mathcal{Q}^{\gamma, \phi, \delta t} \| \mathcal{P}^{\theta, \phi, \delta t}) = \mathbb{E}_{\mathbf{x}_{0:N} \sim \mathcal{Q}^{\gamma, \phi, \delta t}} [\log \tilde{w}] - \log Z = \text{EUBO} - \log Z \quad (43)$$

and thus $\text{EUBO} \geq \log Z$ due to $D_{\text{KL}}(\cdot \| \cdot) \geq 0$. Since the evidence upper bound is based on the mode-seeking forward KL, it is well suited for quantifying mode-collapse. For further details, see Blessing et al. (2024).

- **Estimation error $\Delta \log Z$:** When having access to the ground truth normalization constant $\log Z$, we can compute the estimation error $\Delta \log Z = |\log Z - \log \hat{Z}|$ using an importance weighted estimate, that is,

$$\log \hat{Z} := \log \mathbb{E}_{\mathbf{x}_{0:N} \sim \mathcal{P}^{\theta, \phi, \delta t}} [\tilde{w}] \approx \log \frac{1}{m} \sum_{i=1}^m \tilde{w}^{(i)}. \quad (44)$$

- **Effective sample size (ESS):** Moreover, we compute the (normalized) ESS as

$$\text{ESS} := \frac{(\sum_{i=1}^m \tilde{w}^{(i)})^2}{m \sum_{i=1}^m (\tilde{w}^{(i)})^2}. \quad (45)$$

- **Sinkhorn distance:** We estimate the Sinkhorn distance \mathcal{W}_γ^2 (Cuturi, 2013), i.e., an entropy regularized optimal transport distance between a set of samples from the model and target using the Jax ott library (Cuturi et al., 2022). Note that computing \mathcal{W}_γ^2 requires samples from the target density which are typically not available for real-world target densities.

- **Entropic mode coverage (EMC):** EMC evaluates the mode coverage of a sampler by leveraging prior knowledge of the target density. It holds that $\text{EMC} \in [0, 1]$ where $\text{EMC} = 1$ indicates that the model achieves uniform coverage over all modes whereas $\text{EMC} = 0$ indicates that the model only produces samples from a single mode. Please note that EMC does not provide any information about the sample quality. For further details, we refer the interested reader to Blessing et al. (2024).

D FURTHER NUMERICAL RESULTS

Here, we provide further numerical results.

1080
1081
1082
1083
1084
1085

K	DIS-GMP		CMCD-GMP	
	ABS. [s]	REL. [%]	ABS. [s]	REL. [%]
1	0.103	-	1.123	-
5	0.128	24.27	1.166	3.82
10	0.155	50.48	1.203	7.12

Table 5: Wallclock time of DIS-GMP and CMCD-GMP for the Fashion target for $N = 128$. Here, K the number of components in den Gaussian mixture, ‘abs.’ denotes the absolute time per gradient step in seconds, and ‘rel.’ denotes the relative increase in percent compared to $K = 1$.

1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113

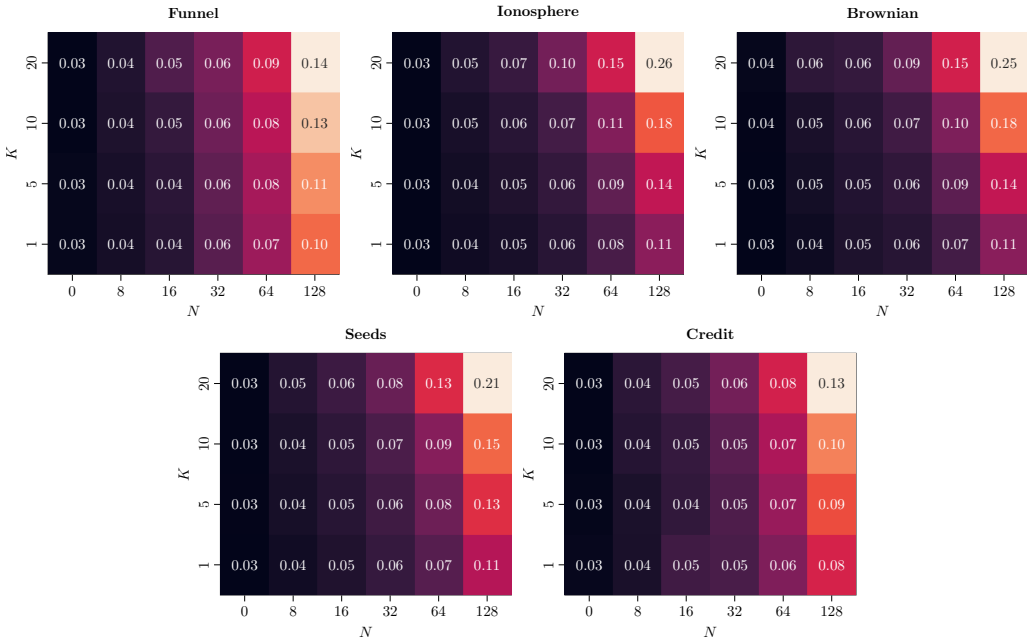


Figure 6: Wallclock time per gradient step of DIS-GMP for various benchmark problems. Here, N denotes the number of discretization steps and K the number of components in den Gaussian mixture.

1114 **Wallclock time** We further report the wallclock time per gradient step for DIS for a different
 1115 number of diffusion steps N and mixture components K . The results are shown in Figure 6. For
 1116 $N \leq 64$, the Gaussian mixture prior barely influences the wallclock time where using $K = 10$
 1117 components roughly adds a 20 percent increase. Considering the performance improvements this is
 1118 a good trade-off. For $N = 128$, Using $K = 10$ roughly results in a 50 percent increase as the like-
 1119 lihood of the prior has to be evaluated in every diffusion step. However, since most diffusion-based
 1120 methods apart from DIS additionally require evaluating the target density at every step, the relative
 1121 costs of using GMPs reduce if the target is more expensive to evaluate. We empirically validated
 1122 this by additionally including a comparison between the wallclock time for DIS and CMCD on the
 1123 Fashion target in Table 5. In this setting, the relative cost added by the GMP is minor.

1124 **Additional results for DIS-GMP** We present further details regarding the ablation study from
 1125 Section 6. Specifically, we report ELBO values for the real-world benchmark problems in Figure 7
 1126 and various metrics for the *Funnel* target in Figure 8. The results are consistent with the results
 1127 in Figure 5, where the performance improves with a higher number of mixture components K and
 1128 diffusion steps N .

1130 **DIS-GMP on *Digits* target** We additionally investigate the performance of DIS-GMP on the syn-
 1131 thetic digits target. The results are reported in Figure 9. Here, we observe that the ELBO get looser
 1132 when using more mixture components K , while $\Delta \log Z$ stays roughly constant. However, the sam-
 1133 ple diversity improves significantly as shown quantitatively from the Sinkhorn distance \mathcal{W}_γ^2 and
 qualitatively in the Figure on the right-hand side.

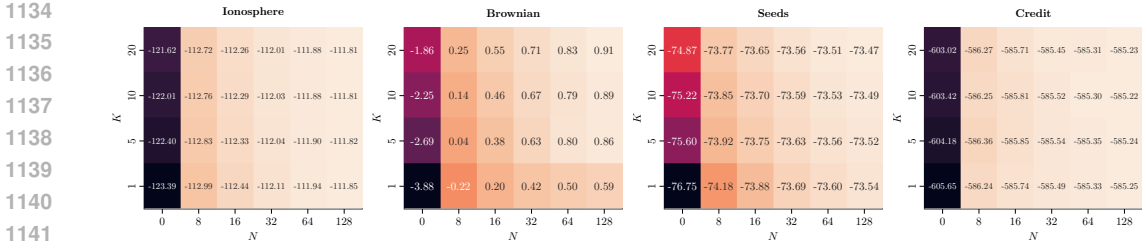


Figure 7: Evidence Lower Bound (ELBO) of DIS-GMP for various real-world benchmark problems, averaged across four seeds. Here, N denotes the number of discretization steps and K the number of components in den Gaussian mixture.

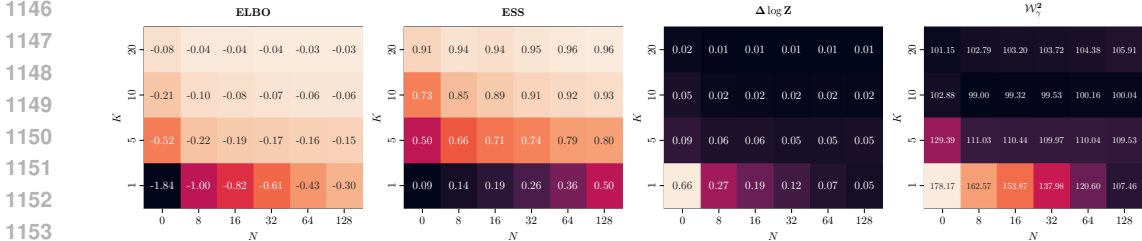


Figure 8: Various performance criteria of DIS-GMP for the Funnel target, averaged across four seeds. Here, N denotes the number of discretization steps and K the number of components in den Gaussian mixture.

Digits ($d = 196$)			
K	ELBO \uparrow	$\Delta \log Z$ \downarrow	\mathcal{W}_2^2 \downarrow
1	-12.090±0.050	5.269±0.416	197.566±0.340
5	-12.303±0.350	4.419±0.316	183.241±8.776
10	-13.820±0.831	4.658±0.260	164.827±2.626
20	-15.413±0.317	5.663±0.085	151.006±0.640

$K = 1$

$K = 5$

$K = 10$

$K = 20$

K : Num. mixture components

Figure 9: **Left side**: Results for Digits target, averaged across four seeds using DIS-GMP+IMR. Evaluation criteria include evidence lower bound ELBO, importance-weighted errors for estimating the log-normalizing constant $\Delta \log Z$, and Sinkhorn distance \mathcal{W}_2^2 . The best results are highlighted in bold. Arrows (\uparrow , \downarrow) indicate whether higher or lower values are preferable, respectively. **Right side**: Visualization of the $d = 14 \times 14 = 196$ dimensional Digits samples for a different number of mixture components K .

Ablation on Gaussian Mixture Target We additionally experiment with using a two-dimensional Gaussian mixture model (GMM) as the target density. The GMM has ten components where the means are uniformly sampled in $[-12, 12]$ and the covariance matrices are sampled from a Wishart distribution. In addition to the performance criteria outlined in Section 6, we quantify the variation of the dynamics over time using the spectral norm of the Jacobian of the learned control, i.e.,

$$S = \mathbb{E}_{\mathbf{x}_0, T \sim \mathcal{P}^\theta} \left[\int_0^T \left\| \frac{\partial \sigma u^\theta(\mathbf{x}, t)}{\partial \mathbf{x}} \right\|_2 dt \right]. \tag{46}$$

For DIS we initialized the prior with a standard deviation of 12 such that the prior covers the support of the target. For DIS-GMP, we use $K = 10$ components that are initialized with a standard deviation of 1. We report qualitative and qualitative results in Figure 10 and Figure 11. We find that DIS without learned prior and sufficiently large prior support is able to cover all modes as indicated by

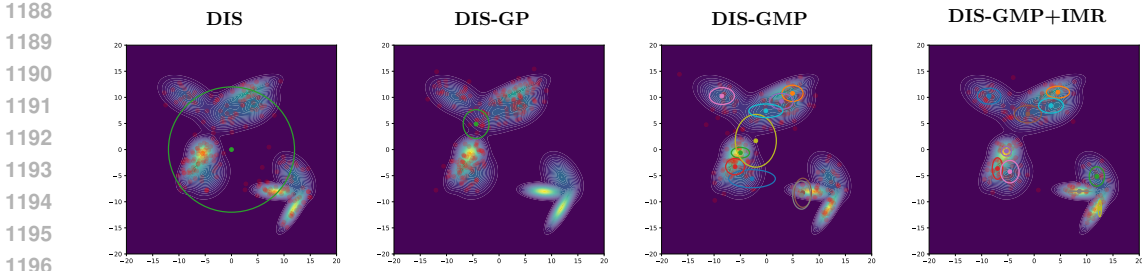


Figure 10: Visualization of a two-dimensional Gaussian mixture target density for different variants of DIS with $N = 128$ diffusion steps and $K = 10$ components for GMP versions. Colored ellipses and circles denote standard deviations and means of the Gaussian components, respectively. Red dots illustrate samples of the learned model.

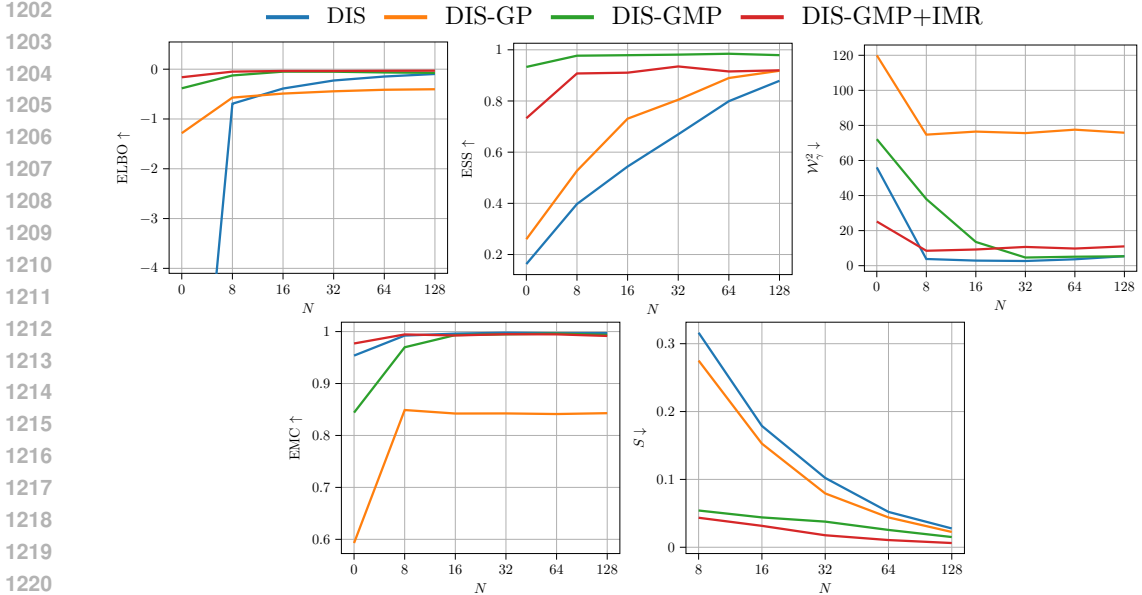


Figure 11: Results for the two-dimensional Gaussian mixture target, averaged across four seeds and reported across different numbers of diffusion steps N for different variants of DIS. Evaluation criteria include evidence lower bound ELBO, importance-weighted errors for estimating the log-normalizing constant $\Delta \log Z$, and Sinkhorn distance \mathcal{W}_2^2 , entropic mode coverage EMC, and the time-integrated spectral norm of the control S (see Equation (46)).

EMC ≈ 1 . While the sample quality is similar between DIS and GMP counterparts (see \mathcal{W}_γ^2), plain DIS needs significantly more diffusion steps in order to achieve similar ELBO/ESS values compared to the GMP counterparts which achieve good performance with as few as 8 diffusion steps. The requirement of plain DIS for more discretization steps is additionally reflected in the variation of the dynamics over time S . Lastly, the GP version is not able to cover all modes due to the mode-seeking nature of the reverse KL as indicated by the EMC and \mathcal{W}_γ^2 values.

Ablation: Influence of the control architecture We further evaluate the performance using the architecture by Zhang & Chen (2021) which additionally incorporates the score of the target, i.e. $\nabla \log \pi$, into the architecture via

$$u^\theta(\mathbf{x}, t) = f_1^\theta(\mathbf{x}, t) + f_2^\theta(t) \nabla \log \pi(\text{stop_gradient}(\mathbf{x})). \quad (47)$$

where f_1 and f_2 are parameterized function approximator with parameters θ . Zhang & Chen (2021); Vargas et al. (2023a) found that detaching, that is, using a stop-gradient operator on \mathbf{x} yields superior

METHOD	$\nabla \log \pi$	CREDIT	SEEDS	CANCER	BROWNIAN	IONOSPHERE	SONAR
DIS-GP	✗	-585.247 ± 0.009	-73.540 ± 0.005	-85.005 ± 1.286	0.588 ± 0.013	-111.847 ± 0.006	-109.280 ± 0.024
DIS-GP	✓	-592.262 ± 0.794	-73.497 ± 0.001	-96.180 ± 10.044	N/A	-111.957 ± 0.090	-109.473 ± 0.143
DIS-GMP	✗	-585.223 ± 0.006	-73.492 ± 0.003	-84.061 ± 2.117	0.885 ± 0.005	-111.811 ± 0.002	-109.157 ± 0.000
DIS-GMP	✓	-586.817 ± 0.906	-73.475 ± 0.002	-84.732 ± 0.466	N/A	-112.108 ± 0.002	-109.248 ± 0.001

Table 6: Evidence lower bound (ELBO) values for various real-world benchmark problems, averaged across four seeds. Here, $\nabla \log \pi$ indicates if the model architecture uses target score as described in Equation (47). The best overall results are highlighted in bold, with category-specific best results underlined. Blue and green shading indicate that the method uses learned Gaussian (GP) and Gaussian mixture priors (GMP), respectively.

METHOD	Div.	CREDIT	SEEDS	CANCER	BROWNIAN	IONOSPHERE	SONAR
DIS	KL	-589.636 ± 0.757	-74.400 ± 0.007	-86.592 ± 2.107	-3.503 ± 0.019	-112.525 ± 0.008	-110.153 ± 0.022
DIS	LV	-5170.845 ± 5.627	-74.654 ± 0.022	-88.379 ± 1.491	-5.682 ± 0.303	-112.609 ± 0.053	-110.622 ± 0.071
DIS-GP	KL	-585.247 ± 0.009	-73.540 ± 0.005	-85.005 ± 1.286	0.588 ± 0.013	-111.847 ± 0.006	-109.280 ± 0.024
DIS-GP	LV	-5163.451 ± 3.296	-73.703 ± 0.177	-549.071 ± 466.902	<u>0.729 ± 0.004</u>	-111.839 ± 0.006	-109.498 ± 0.005
DIS-GMP	KL	-585.223 ± 0.006	-73.492 ± 0.003	-84.061 ± 2.117	0.885 ± 0.005	-111.811 ± 0.002	-109.157 ± 0.000
DIS-GMP	LV	-5152.728 ± 18.004	-73.777 ± 0.007	-86.456 ± 0.557	0.722 ± 0.005	-111.844 ± 0.000	-109.443 ± 0.000

Table 7: Evidence lower bound (ELBO) values for various real-world benchmark problems, averaged across four seeds. Here, ‘Div.’ indicates if the model is trained using the Kullback-Leibler (KL) or log-variance (LV) divergence. The best overall results are highlighted in bold, with category-specific best results underlined. Blue and green shading indicate that the method uses learned Gaussian (GP) and Gaussian mixture priors (GMP), respectively.

results due to the simplification of the computational graph. We adopt this change and report the results in Table 6. We find that using the score of the target leads, in the majority of experiments, slightly worse results, with one exception where it yields superior results.

Ablation: Kullback-Leibler vs. Log-Variance divergence We further compare the KL divergence to the log-variance divergence introduced in Richter et al. (2020) and later extended to diffusion models in Richter et al. (2023). The log-variance divergence is defined as

$$\mathcal{L}(\theta, \phi, \delta t) = \mathbb{V}_{\mathbf{x}_{0:N} \sim \mathcal{R}} \left[\log \frac{Q^{\gamma, \phi, \delta t}(\mathbf{x}_{0:N})}{\mathcal{P}^{\theta, \phi, \delta t}(\mathbf{x}_{0:N})} \right], \quad (48)$$

where \mathcal{R} describes a reference process, e.g. Equation (2a) where u^θ is replaced with an arbitrary control. In practice, one typically uses the generative process $\mathcal{P}^{\theta, \phi, \delta t}$ with an additional stop gradient operator on the parameters (Richter et al., 2023). Not computing the expectations with respect to samples from the generative process significantly reduces memory consumption and does not require the prior distribution to be amendable to the reparameterization trick. The results are reported in Table 7 and follow the same experimental setting as outlined in the main part of the paper. We find that the KL divergence typically performs better than the log-variance divergence. Most significantly, the log-variance divergence seems to be numerically unstable for the *Credit* target.

Comparison to long-run Sequential Monte Carlo We additionally compare diffusion samplers with a learned Gaussian mixture prior to a Sequential Monte Carlo with a high number of discretization steps N . The results are shown in Table 8 and Table 9. While long-run SMC significantly increases ELBO values, GMP-based diffusion sampler yield superior results in most experiments. Moreover, the results on the *Fashion* target indicate that more discretization steps yield better ELBO, $\Delta \log Z$ and Sinkorn distances, but are not able to prevent mode collapse as indicated by the low EMC values.

METHOD	N	CREDIT	SEEDS	CANCER	BROWNIAN	IONOSPHERE	SONAR
MCD-GMP	128	-585.276±0.013	-73.461±0.004	-88.562±0.243	0.993±0.003	-111.827±0.007	-109.197±0.004
CMCD-GMP	128	-585.162±0.002	-73.429±0.002	-78.402±0.037	1.087±0.001	-111.682±0.000	-108.634±0.000
DIS-GMP	128	-585.223±0.006	-73.492±0.003	-84.061±2.117	0.885±0.005	-111.811±0.002	-109.157±0.000
DBS-GMP	128	-585.148±0.002	-73.418±0.001	-78.160±0.063	1.118±0.002	-111.657±0.002	-108.548±0.000
SMC	128	-698.403±4.146	-74.699±0.100	-194.059±0.613	-1.874±0.622	-114.751±0.238	-111.355±1.177
	256	-708.185±14.225	-73.972±0.034	-140.757±7.041	-0.360±0.136	-113.110±0.046	-109.822±0.630
	512	-686.335±18.333	-73.667±0.015	-137.028±2.336	0.414±0.048	-112.353±0.036	-109.197±0.420
	1024	-690.011±12.879	-73.532±0.038	-128.809±6.046	0.786±0.047	-111.962±0.018	-108.291±0.325
	2048	-672.602±15.229	-73.496±0.017	-128.376±3.504	0.992±0.036	-111.785±0.022	-108.261±0.565
	4096	-665.973±19.849	-73.438±0.004	-121.950±3.315	1.088±0.029	-111.692±0.013	-108.736±0.227

Table 8: Evidence lower bound (ELBO) values for various real-world benchmark problems and different numbers of discretization steps N , averaged across four seeds. The best overall results are highlighted in bold, with category-specific best results underlined. green shading indicate that the method Gaussian mixture priors (GMP).

		FASHION ($d = 784$)			
METHOD	N	ELBO \uparrow	$\Delta \log Z$ \downarrow	\mathcal{W}_2^2 \downarrow	EMC \uparrow
DIS-GMP+IMR	128	-62.482±2.752	27.645±3.118	513.776±13.936	0.780±0.089
SMC	128	-12181.932±134.611	11747.518±139.292	6696.287±250.4	0.026±0.027
	256	-10095.076±1076.723	9901.113±1078.916	6018.423±197.144	<u>0.191±0.112</u>
	512	-9340.232±803.694	9254.499±804.027	5821.422±396.492	0.141±0.140
	1024	-9229.557±742.223	9190.558±742.656	5610.511±283.885	0.075±0.129
	2048	-8472.660±281.288	8454.062±281.475	5718.030±328.909	0.257±0.038
	4096	<u>-8399.465±153.637</u>	<u>8390.302±153.707</u>	<u>5583.099±179.370</u>	0.102±0.108

Table 9: Results for Fashion target, averaged across four seeds and reported across different numbers of discretization steps N . Evaluation criteria include evidence lower bound ELBO, importance-weighted errors for estimating the log-normalizing constant $\Delta \log Z$, and Sinkhorn distance \mathcal{W}_2^2 and entropic mode coverage EMC. The best overall results are highlighted in bold, with category-specific best results underlined. Arrows (\uparrow, \downarrow) indicate whether higher or lower values are preferable, respectively. Orange shading indicates that the method uses iterative model refinement (IMR).

Ablation: Iterations for IMR We additionally conducted an ablation study which considers different numbers of iterations for the iterative model refinement scheme at which new components are added. The results are reported in Table 10 and indicate that the performance remains stable for different choices of the hyperparameter. Nevertheless, the concept of iteratively adding components is important, such that when a new component is added, the initialization is informed by the already existing mixture model (note that the heuristic in Eq. 22 depends on the likelihood of the GMP) to prevent initializing components at the same location.

Ablation: Variation of dynamics We additionally compare the variability in the dynamics of the learned model between DIS and DIS-GMP via time-integrated spectral norm of the control S (see Equation (46)). The results are shown in Table 11 and show that DIS-GMP indeed has less variation in the dynamics. These findings are also in line with those in Figure 3 and Table 2 where DIS-GMP has significantly higher ELBO values compared to DIS without learned prior.

E LATTICE ϕ^4 THEORY

We apply our method to simulate a statistical lattice field theory near and beyond the phase transition. This phase transition marks the progression of the lattice from disordered to semi-ordered and ultimately to a fully ordered state, where neighboring sites exhibit strong correlations in sign and magnitude.

FASHION ($d = 784$)				
IMR ITER.	ELBO \uparrow	$\Delta \log Z \downarrow$	$\mathcal{W}_\gamma^2 \downarrow$	EMC \uparrow
100	-60.129 ± 3.045	26.473 ± 3.765	520.483 ± 14.923	0.764 ± 0.091
500	-62.482 ± 2.752	27.645 ± 3.118	513.776 ± 13.936	0.780 ± 0.089
1000	-65.784 ± 3.567	32.134 ± 4.054	538.145 ± 15.678	0.751 ± 0.097
2000	-61.478 ± 3.321	25.129 ± 3.832	505.781 ± 14.342	0.785 ± 0.094

Table 10: Results for Fashion target, averaged across four seeds and reported for different numbers of iterations at which components are added (IMR. iter). Evaluation criteria include evidence lower bound ELBO, importance-weighted errors for estimating the log-normalizing constant $\Delta \log Z$, and Sinkhorn distance \mathcal{W}_γ^2 and entropic mode coverage EMC. Arrows (\uparrow , \downarrow) indicate whether higher or lower values are preferable, respectively.

METHOD	FUNNEL	SEEDS	BROWNIAN	IONOSPHERE	SONAR
DIS	2.993 ± 0.042	4.688 ± 0.055	6.266 ± 0.329	4.394 ± 0.066	4.840 ± 0.031
DIS-GMP	1.898 ± 0.002	2.367 ± 0.008	2.445 ± 0.004	2.736 ± 0.004	3.861 ± 0.036

Table 11: Variability in the dynamics of the learned model via time-integrated spectral norm of the control $S \times 10^2$ (see Equation (46)) for various benchmark problems. Both DIS and DIS-GMP use $N = 128$ diffusion steps. Here, DIS-GMP uses $K = 10$ components. Lower values indicate lower variability in the dynamics of the learned model.

We study the lattice ϕ^4 theory in $D = 2$ spacetime dimensions (distinct from the problem’s dimensionality as described below). The random variables in this setting are field configurations $\phi \in \mathbb{R}^{L \times L}$, where L represents the lattice extent in space and time. The density of these configurations is defined as

$$\pi(\phi) = \frac{e^{-U(\phi)}}{Z},$$

where the potential $U(\phi)$ is given by:

$$U(\phi) = -2\kappa \sum_x \sum_\mu \phi_x \phi_{x+\mu} + (1 - 2\lambda) \sum_x \phi_x^2 + \lambda \sum_x \phi_x^4. \quad (49)$$

Here, the summation over x runs over all lattice sites, and the summation over μ considers the neighbors of each site. The parameters λ and κ are referred to as the bare coupling constant and the hopping parameter, respectively. Following Nicoli et al. (2021), we set $\lambda = 0.022$, identifying the critical threshold of the theory (the transition from ordered to disordered states) at $\kappa \geq 0.3$. Near this threshold, sampling becomes increasingly challenging due to the multimodality of the density, with modes becoming more separated for larger values of κ .

We conduct experiments for $\kappa \in \{0.2, 0.3, 0.5\}$ across various problem dimensions $d = L \times L$. The methods compared include DIS, DIS-GP, and DIS-GMP, each with $N = 128$ diffusion steps, as well as a long-run SMC sampler with $N = 4096$. For all methods, the initial support is set to 5, approximately covering the target’s support for all tested values of κ . The tuned parameters of the HMC kernel for the SMC sampler are detailed in Table 4, while additional parameter settings are provided in Appendix C.2. Note that DIS (and its extensions) do not undergo hyperparameter tuning due to their end-to-end learning framework.

To compare the different methods, we utilize the negative variational free energy of the system, defined as:

$$-\mathcal{F} = \frac{1}{L^2} \log Z \geq \frac{1}{L^2} \text{ELBO}. \quad (50)$$

This bound follows from the inequality $\log Z \geq \text{ELBO}$ as discussed in Section 3.2 and provides a means of comparison between sampling methods. However, as the ELBO (and thus \mathcal{F}) is not

METHOD	κ	$d = 16$	$d = 64$	$d = 100$	$d = 144$	$d = 196$	$d = 256$
DIS	0.2	0.6263±0.0000	0.6186±0.0009	0.6166±0.0008	0.6145±0.0014	0.5997±0.0011	0.5749±0.0013
DIS-GP		0.6274±0.0000	0.6219±0.0001	0.6200±0.0001	0.6175±0.0000	0.6158±0.0001	0.6113±0.0015
DIS-GMP		0.6276±0.0000	0.6232±0.0004	0.6231±0.0002	0.6166±0.0027	0.6139±0.0006	0.6156±0.0005
SMC		0.6167±0.0022	0.6175±0.0005	0.6186±0.0003	0.6164±0.0072	0.6087±0.0029	0.6142±0.0099
DIS	0.3	1.0653±0.0195	1.0277±0.0004	1.0217±0.0004	1.0040±0.0024	0.9534±0.0017	0.8905±0.0014
DIS-GP		1.0831±0.0021	1.0459±0.0001	1.0411±0.0001	1.0372±0.0011	1.0343±0.0010	1.0319±0.0004
DIS-GMP		1.0940±0.0000	1.0496±0.0033	1.0461±0.0001	1.0406±0.0001	1.0380±0.0007	1.0347±0.0006
SMC		1.0848±0.0013	1.0514±0.0007	1.0488±0.0002	1.0437±0.0052	1.0339±0.0043	1.0389±0.0098
DIS	0.5	12.2545±0.0004	12.2097±0.0023	9.2610±2.9035	11.7179±0.0799	9.4392±1.9044	10.8546±0.0716
DIS-GP		12.2735±0.0000	12.2715±0.0000	12.2692±0.0012	12.2670±0.0010	10.1714±2.0923	12.2629±0.0001
DIS-GMP		12.3167±0.0001	12.2806±0.0003	12.2761±0.0003	12.2730±0.0002	12.2722±0.0000	12.2588±0.0001
SMC		12.3049±0.0025	12.2707±0.0013	12.2679±0.0027	12.2499±0.0101	12.2416±0.0040	12.2407±0.0024

Table 12: Lower bound values for negative variational free energy $-\mathcal{F}$ as defined in Eq. 50 for the lattice ϕ^4 theory problem with different values for the space-time extend $\sqrt{d} = L$ averaged across two seeds. The best (i.e. the highest) overall results are highlighted in bold for each configuration of the hopping parameter κ and space-time extend.

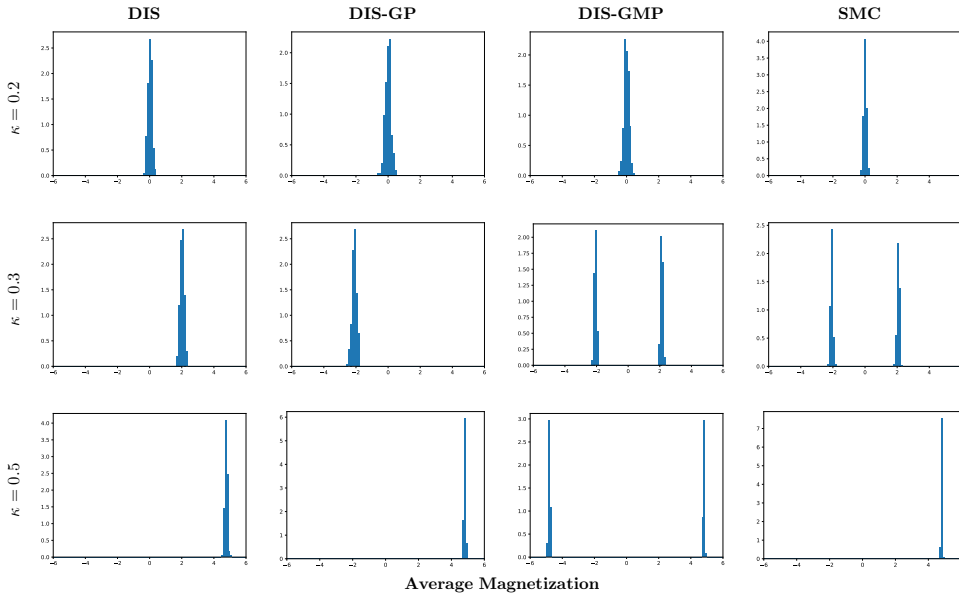


Figure 12: Normalized histogram of the average magnetization $M(\phi) = \sum_x \phi_x$ for 2000 samples $\phi \in \mathbb{R}^{L \times L}$ and space-time extend $L = 14$ for DIS, DIS-GP, DIS-GMP and long-run SMC for different values of the hopping parameter κ . The plots are generated using the same random seed 0.

sensitive to mode collapse (Blessing et al., 2024), and since samples from the target distribution are unavailable, we also qualitatively assess the methods by visualizing the (normalized) histogram of the average magnetization $M(\phi) = \sum_x \phi_x$ across lattice configurations ϕ .

Quantitative results are presented in Table 12, with qualitative findings illustrated in Figure 12. The results indicate that learning the prior (i.e. DIS-GP/DIS-GMP) significantly improves free energy estimates compared to DIS without a learned prior. Moreover, Figure 12 demonstrates that DIS-GMP avoids mode collapse while achieving comparable or better free energy estimates than both DIS-GP and the long-run SMC sampler in the majority of settings. While SMC captures multimodality at the phase transition ($\kappa = 0.3$), it struggles with the multimodality in the fully ordered phase ($\kappa = 0.5$). By contrast, DIS and DIS-GP are prone to mode collapse. Lastly, the performance of DIS degrades significantly with increasing problem dimension d , which is mitigated when using a learned Gaussian or Gaussian mixture prior.