



HAL
open science

Diffusion models for Gaussian distributions: Exact solutions and Wasserstein errors

Emile Pierret, Bruno Galerne

► **To cite this version:**

Emile Pierret, Bruno Galerne. Diffusion models for Gaussian distributions: Exact solutions and Wasserstein errors. 2024. hal-04584209v1

HAL Id: hal-04584209

<https://hal.science/hal-04584209v1>

Preprint submitted on 23 May 2024 (v1), last revised 7 Jun 2024 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Diffusion models for Gaussian distributions: Exact solutions and Wasserstein errors

Emile Pierret

Institut Denis Poisson
Université d'Orléans
Orléans, France

`emile.pierret@univ-orleans.fr`

Bruno Galerne

Institut Denis Poisson
Université d'Orléans
Orléans, France

Institut universitaire de France (IUF)
`bruno.galerie@univ-orleans.fr`

Abstract

Diffusion or score-based models recently showed high performance in image generation. They rely on a forward and a backward stochastic differential equations (SDE). The sampling of a data distribution is achieved by solving numerically the backward SDE or its associated flow ODE. Studying the convergence of these models necessitates to control four different types of error: the initialization error, the truncation error, the discretization and the score approximation. In this paper, we study theoretically the behavior of diffusion models and their numerical implementation when the data distribution is Gaussian. In this restricted framework where the score function is a linear operator, we can derive the analytical solutions of the forward and backward SDEs as well as the associated flow ODE. This provides exact expressions for various Wasserstein errors which enable us to compare the influence of each error type for any sampling scheme, thus allowing to monitor convergence directly in the data space instead of relying on Inception features. Our experiments show that the recommended numerical schemes from the diffusion models literature are also the best sampling schemes for Gaussian distributions.

1 Introduction

Over the last five years, diffusion models have proven to be a highly efficient and reliable framework for generative modeling [25, 14, 24, 26, 7, 15]. First introduced as a discrete process, Denoising Diffusion Probabilistic Models (DDPM) [14] can be studied as a reversal of a continuous Stochastic Differential Equation (SDE) [26]. A forward SDE progressively transforms the initial data distribution by adding more and more noise as time grows. Then, the reversal of this process, called backward SDE, allows us to approximately sample the data distribution starting from Gaussian white noise. Moreover, the SDE is associated with an Ordinary Differential Equations (ODE) called probability flow [26]. This flow preserves the same marginal distributions as the backward SDE and provides another way to sample the score-based generative model.

An important issue about diffusion models is the theoretical guarantees of convergence of the model: How close to the data distribution the generate distribution is? There are four main source of errors to study for deriving theoretical guarantees for diffusion models. First, the *initialization error* is induced when approximating the marginal distribution at the end of the forward process by a standard Gaussian distribution. Then, the *discretization error* comes from the resolution of the SDE or the ODE by a numerical method. The *truncation error* occurs because the backward time integration is stopped at a small time $\varepsilon > 0$ to avoid numerical instabilities due to ill-defined score function near the origin. Finally the *score approximation error* accounts for the mismatch between the ideal score function and the one given by the network trained using denoising score-matching.

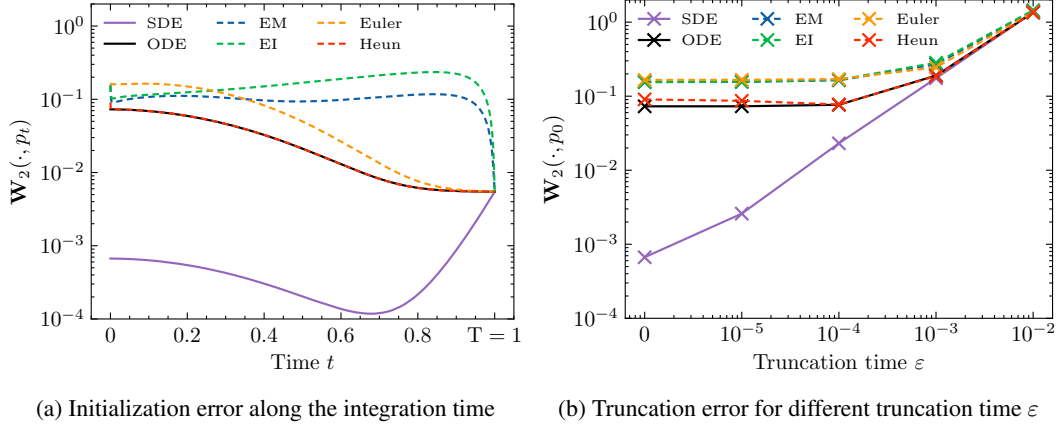


Figure 1: **Wasserstein errors for the diffusion models associated with the CIFAR-10 Gaussian.** Left: Evolution of the Wasserstein distance between p_t and the distributions associated with the continuous SDE, the continuous flow ODE and four discrete sampling schemes with standard \mathcal{N}_0 initialization, either stochastic (Euler-Maruyama (EM) and Exponential Integrator (EI)) or deterministic (Euler and Heun). While the continuous SDE is less sensible than the continuous ODE (as proved by Proposition 4), the initialization error impacts all discrete schemes with a comparable order of magnitude. Heun’s method has the lowest error and is very close to the theoretical ODE, except for the last step that is usually discarded when using time truncation. Right: Wasserstein errors due to time truncation for various truncation times ε . Using time truncation increases the error for all the methods except Heun’s scheme due to instability near the origin. Interestingly, for the standard practice truncation time $\varepsilon = 10^{-3}$, all numerical schemes have a comparable error close to their continuous counterparts.

Despite these numerous sources of errors, a lot of numerical and theoretical research has been led to assess the generative capacity of diffusion models. Several articles [3, 10, 15] provide strong experimental studies for the choices of sampling parameters. On the theoretical side, several works derive upper bounds on the 1-Wasserstein or TV distance between the data and the model distributions by making assumptions on the L^2 -error between the ideal and learned score functions and on the compacity of the support of the data [19, 6, 4, 18, 1], eventually under an additional manifold assumption [5, 28, 2]. Yet, on one hand, to the best of our knowledge, the derived theoretical bounds mostly rely on worst case scenario and are not tight enough to explain the practical efficiency of diffusion models. On the other hand, numerical considerations mostly rely on Inception feature distributions through the FID metric [13].

Ideally, given a data distribution of interest, one would like to have an adapted estimation of the discrepancy between the data and the diffusion model samples, thus enabling adaptive hyperparameter selection for the sampling procedure. As a first step towards reaching this goal, in the present work we study diffusion models applied to Gaussian data distributions. While this setting has a priori no practical interest, since simulating Gaussian variates does not require a diffusion model, it provides a large parametric family of distributions for which the errors involved in diffusion model sampling can be completely understood.

When restricting the data distribution to be Gaussian, the resulting score function is a simple linear operator. Exploiting this specificity allows us to derive the following contributions **under the assumption that the data is Gaussian**:

- We give the exact solutions for both the backward SDE and probability flow ODE.
- We fully describe the Gaussian processes that occurs when using classical sampling discretization schemes.
- We derive exact 2-Wasserstein errors for the corresponding sample distributions and are able to assert for the influence of each error type on these errors, as illustrated by Figure 1.

Our theoretical study allows for a thorough experimental analysis of any numerical sampler, either stochastic or deterministic. In particular, it confirms the strength of best practice scheme such as Heun’s method for the ODE flow [15].

While our theoretical analysis relies on an exact known score function, we conduct additional experiments to assess for the influence of the score approximation error. Surprisingly, in the context of texture synthesis, we show that with a score neural network trained for modeling a specific Gaussian micro-texture a stochastic Euler-Maruyama sampler is more faithful to the data distribution than Heun’s method, thus highlighting the importance of the score approximation error in practical situations.

Plan of the paper: First, we recall in Section 2 the continuous framework for SDE-based diffusion models. Section 3 presents our main theoretical results detailing the exact backward SDE and probability flow ODE solutions when supposing the data distribution to be Gaussian. Section 4 gives explicit Wasserstein error formulas when sampling the corresponding processes, yielding to an ablation study for comparing the influence of each error type on several sampling schemes. Finally, in Section 5 we study numerically a special case of Gaussian distribution for texture synthesis in order to evaluate the influence of the score approximation error occurring with a standard network architecture.

2 Preliminaries: Score-based models through diffusion SDEs

This preliminary section follows the seminal work of Song *et al.* [26] and introduces specific notation to differentiate the exact backward process and the generative backward process obtained when starting from a white noise. Given a target distribution p_{data} over \mathbb{R}^d , the forward diffusion process is the following variance preserving SDE

$$d\mathbf{x}_t = -\beta_t \mathbf{x}_t dt + \sqrt{2\beta_t} d\mathbf{w}_t, \quad 0 \leq t \leq T, \quad \mathbf{x}_0 \sim p_{\text{data}} \quad (1)$$

where $(\mathbf{w}_t)_{t \geq 0}$ is a d -dimension Brownian motion and β is a positive weight function. The distribution p_{data} is noised progressively and the function β is the variance of the added noise by time unit. We denote by p_t the density of (\mathbf{x}_t) for $t > 0$ since p_{data} can be supported on a lower-dimensional manifold [5]. The SDE is designed so that p_T is close to the Gaussian standard distribution that we denote \mathcal{N}_0 in whole paper. Under some assumptions on the distribution p_{data} [21], the backward process $(\mathbf{x}_{T-t})_{0 \leq t \leq T}$ verifies the backward SDE

$$d\mathbf{y}_t = \beta_{T-t}(\mathbf{y}_t + 2\nabla \log p_{T-t}(\mathbf{y}_t))dt + \sqrt{2\beta_{T-t}}d\mathbf{w}_t, \quad 0 \leq t < T, \quad \mathbf{y}_0 \sim p_T. \quad (2)$$

The objective is now to solve this reverse equation to sample $\mathbf{y}_T \sim p_{\text{data}}$. However, the distribution p_T is in general not known, and image¹ generation is achieved by sampling

$$d\tilde{\mathbf{y}}_t = \beta_{T-t}(\tilde{\mathbf{y}}_t + 2\nabla \log p_{T-t}(\tilde{\mathbf{y}}_t))dt + \sqrt{2\beta_{T-t}}d\mathbf{w}_t, \quad 0 \leq t < T, \quad \tilde{\mathbf{y}}_0 \sim \mathcal{N}_0. \quad (3)$$

Note that approximating p_T by \mathcal{N}_0 for the initialization \mathbf{y}_0 makes that the solution of the SDE (3) is not exactly the target distribution p_{data} . An alternative way to approximately sample p_{data} is to use that every diffusion process is associated with a deterministic process whose trajectories share the same marginal probability densities $(p_t)_{0 < t \leq T}$ as the SDE [26]. The deterministic process associated with Equation (2) is

$$d\mathbf{x}_t = [-\beta_t \mathbf{x}_t - \beta_t \nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t)] dt, \quad 0 < t \leq T, \quad \mathbf{x}_0 \sim p_{\text{data}}. \quad (4)$$

This ODE can be solved in reverse-time to sample \mathbf{x}_0 from $\mathbf{x}_T \sim p_T$. Given $(\mathbf{x}_t)_{0 \leq t \leq T}$ solution of Equation (4), $(\mathbf{x}_{T-t})_{0 \leq t \leq T}$ is solution of

$$d\mathbf{y}_t = [\beta_{T-t} \mathbf{y}_t + \beta_{T-t} \nabla_{\mathbf{y}} \log p_{T-t}(\mathbf{y}_t)] dt, \quad 0 \leq t < T. \quad (5)$$

Again, in practice, the ODE which is considered to achieve image generation is

$$d\hat{\mathbf{y}}_t = [\beta_{T-t} \hat{\mathbf{y}}_t + \beta_{T-t} \nabla_{\hat{\mathbf{y}}} \log p_{T-t}(\hat{\mathbf{y}}_t)] dt, \quad 0 \leq t < T, \quad \hat{\mathbf{y}}_T \sim \mathcal{N}_0, \quad (6)$$

where p_T is replaced by \mathcal{N}_0 . As a consequence of this approximation, the property of conservation of the marginals $(p_t)_{0 \leq t \leq T}$ does not occur. We denote by $(\tilde{q}_t)_{0 \leq t \leq T}$, respectively $(\hat{q}_t)_{0 \leq t \leq T}$, the marginals of $(\tilde{\mathbf{y}}_t)_{0 \leq t \leq T}$ and $(\hat{\mathbf{y}}_t)_{0 \leq t \leq T}$ and $\tilde{p}_t = \tilde{q}_{T-t}$, $\hat{p}_t = \hat{q}_{T-t}$ the marginals of $(\tilde{\mathbf{y}}_{T-t})_{0 \leq t \leq T}$ and $(\hat{\mathbf{y}}_{T-t})_{0 \leq t \leq T}$ such that \tilde{p}_t and \hat{p}_t are approximations of p_t .

¹Although we may refer to data as images, our analysis is fully general and applies to any vector-valued diffusion model.

3 Exact SDE and ODE solutions

Our approach relies on deriving explicit solutions to the various SDE and ODE. We begin with the forward SDE in full generality obtained in applying the variation of constants (see the proof in Appendix B.1). This resolution also provides an ODE verified by the covariance matrix of \mathbf{x}_t , that we denote $\Sigma_t = \text{Cov}(\mathbf{x}_t)$.

Proposition 1 (Solution of the forward SDE). *The strong solution of Equation (1) can be written as:*

$$\mathbf{x}_t = e^{-B_t} \mathbf{x}_0 + \boldsymbol{\eta}_t, \quad 0 \leq t \leq T, \quad (7)$$

where $B_t = \int_0^t \beta_s ds$ and $\boldsymbol{\eta}_t = e^{-B_t} \int_0^t e^{B_s} \sqrt{2\beta_s} d\mathbf{w}_s$ is a Gaussian process independent of \mathbf{x}_0 whose covariance matrix is $(1 - e^{-2B_t})\mathbf{I}$. Consequently, the covariance matrix Σ_t of \mathbf{x}_t is

$$\Sigma_t = e^{-2B_t} \Sigma + (1 - e^{-2B_t})\mathbf{I}. \quad (8)$$

where Σ is the covariance matrix of $\mathbf{x}_0 \sim p_{\text{data}}$. Furthermore, Σ_t is invertible for $t > 0$ and verifies the matrix-valued ODE

$$d\Sigma_t = -2(\beta_t \Sigma_t + \beta_t \mathbf{I}) dt, \quad 0 < t \leq T. \quad (9)$$

For a general data distribution p_{data} , solving the backward SDE is infeasible, the main reason being that the expression of the score function to integrate is unknown. To circumvent this obstacle, we now suppose that the data distribution is Gaussian.

Assumption 1 (Gaussian assumption). p_{data} is a centered Gaussian distribution $\mathcal{N}(\mathbf{0}, \Sigma)$.

Note that Σ may be non-invertible and thus p_{data} supported on a strict subspace of \mathbb{R}^d , a special case of manifold hypothesis. Consequently, the matrix Σ_t is in general only invertible for $t > 0$. Under Gaussian assumption, (\mathbf{x}_t) is a Gaussian process with marginal distribution $p_t = \mathcal{N}(\mathbf{0}, \Sigma_t)$ and consequently the score is a linear function

$$\nabla \log p_t(\mathbf{x}) = -\Sigma_t^{-1} \mathbf{x}, \quad 0 < t \leq T. \quad (10)$$

Note that the linearity of the diffusion score characterizes Gaussian distributions as detailed by Proposition 5 in Appendix A.

The cornerstone of our work is that under Gaussian assumption we can derive an exact solution of the backward SDE, without supposing that the initial condition is Gaussian.

Proposition 2 (Solution of the backward SDE under Gaussian assumption). *Under Gaussian assumption, the strong solution to Equation (2) can be written as:*

$$\mathbf{y}_t = e^{-(B_T - B_{T-t})} \Sigma_{T-t}^{-1} \Sigma_T^{-1} \mathbf{y}_0 + \boldsymbol{\xi}_t, \quad 0 \leq t \leq T \quad (11)$$

where $\boldsymbol{\xi}_t = e^{-(B_T - B_{T-t})} \Sigma_{T-t}^{-1} \int_0^t \Sigma_{T-s}^{-1} e^{-(B_T - B_{T-s})} \sqrt{2\beta_{T-s}} d\mathbf{w}_s$ is a Gaussian process with covariance matrix $\text{Cov}(\boldsymbol{\xi}_t) = \Sigma_{T-t} - e^{-2(B_T - B_{T-t})} \Sigma_{T-t}^2 \Sigma_T^{-1}$. Finally:

$$\text{Cov}(\mathbf{y}_t) = \Sigma_{T-t} + e^{-2(B_T - B_{T-t})} \Sigma_{T-t}^2 \Sigma_T^{-1} (\Sigma_{T-t}^{-1} \text{Cov}(\mathbf{y}_0) \Sigma_T^{-1} \Sigma_{T-t} - \mathbf{I}), \quad (12)$$

and in particular, if $\text{Cov}(\mathbf{y}_0)$ and Σ commute,

$$\text{Cov}(\mathbf{y}_t) = \Sigma_{T-t} + e^{-2(B_T - B_{T-t})} \Sigma_{T-t}^2 \Sigma_T^{-1} [\Sigma_T^{-1} \text{Cov}(\mathbf{y}_0) - \mathbf{I}] \quad (13)$$

While not as straightforward as the forward case, the proof also relies on applying the variation of constants and is given in Appendix B.2. As shown by the following proposition (proved in Appendix B.3), the flow ODE also has an explicit solution under Gaussian assumption which related to optimal transport.

Proposition 3 (Solution of the ODE probability flow under Gaussian assumption). *The solution to the probability flow ODE (4) under Gaussian assumption corresponds to the optimal transport map between p_T and p_{data} . More precisely, for any \mathbf{y}_0 ,*

$$\mathbf{y}_t = \Sigma_T^{-1/2} \Sigma_{T-t}^{1/2} \mathbf{y}_0, \quad 0 \leq t \leq T,$$

is the solution of the reverse-time ODE (5). Consequently, the covariance matrix $\text{Cov}(\mathbf{y}_t)$ verifies

$$\text{Cov}(\mathbf{y}_t) = \Sigma_T^{-1/2} \Sigma_{T-t}^{1/2} \text{Cov}(\mathbf{y}_0) \Sigma_{T-t}^{1/2} \Sigma_T^{-1/2}, \quad 0 \leq t \leq T, \quad (14)$$

and in particular, if $\text{Cov}(\mathbf{y}_0)$ and Σ commute,

$$\text{Cov}(\mathbf{y}_t) = \Sigma_T^{-1} \Sigma_{T-t} \text{Cov}(\mathbf{y}_0), \quad 0 \leq t \leq T. \quad (15)$$

Here we must highlight a subtle issue: Whatever the initial distribution is, the ODE solution consists in applying the optimal transport map between p_T and p_{data} . Since in practice one cannot truly sample p_T , the resulting flow is not an optimal transport flow and the distribution of \mathbf{y}_T differs from p_{data} .

Links with related work. Some parts of the previous propositions have been stated in previous work. The expression of the SDE solution of Proposition 1 is given without proof in [12] and the ODE verified by the variance is given in [26] (citing [22]) but it is generalized here to the full covariance matrix (Equation (9)). To the best of our knowledge, Proposition 2 is new and one of the most important contribution of the paper. Similar to our approach, Gaussian mixtures have been studied in the context of diffusion model [29, 30, 23] since they also provide an explicit analytical score. However, solving the backward SDE is not feasible for Gaussian mixtures as far as we know. The ODE [17, 16] can be interpreted in the infinite time as an optimal transport (OT) between the prior distribution and the Gaussian standard distribution. The relation between optimal transport and probability flow ODE (also called Fokker-Planck ODE) has been discussed in [17, 16] in the asymptotic case where $T \mapsto +\infty$. Both these papers discussed the Gaussian case, but our Proposition 3 highlights that the generated process is not an optimal transport flow due to the initialization error.

4 Exact Wasserstein errors

The specificity of the Gaussian case allows us to study precisely the different types of error with the expression of the explicit solution of the backward SDE. In what follows, we designate by Wasserstein distance the 2-Wasserstein distance which is known in closed forms when applied to Gaussian distributions[8]. For two centered Gaussians $\mathcal{N}(\mathbf{0}, \Sigma_1)$ and $\mathcal{N}(\mathbf{0}, \Sigma_2)$ such that Σ_1, Σ_2 are simultaneously diagonalizable with respective eigenvalues $(\lambda_{i,1})_{1 \leq i \leq d}, (\lambda_{i,2})_{1 \leq i \leq d}$,

$$\mathbf{W}_2(\mathcal{N}(\mathbf{0}, \Sigma_1), \mathcal{N}(\mathbf{0}, \Sigma_2))^2 = \sum_{1 \leq i \leq d} (\sqrt{\lambda_{i,1}} - \sqrt{\lambda_{i,2}})^2. \quad (16)$$

In the literature, the quality of the diffusion models is measured with FID [13] which is the \mathbf{W}_2 -error between Gaussians fitted to the Inception features [27] of two discrete datasets. Here we use the \mathbf{W}_2 -errors directly in data space, which is more informative and allows us to provide theoretical \mathbf{W}_2 -errors. To illustrate our theoretical results, we consider the CIFAR-10 Gaussian distribution, that is, the Gaussian distribution such that Σ is the empirical covariance of the CIFAR-10 dataset. As shown in Appendix C, images produced by this model are not interesting due to a lack of structure, but the corresponding covariance has the advantage of reflecting the complexity of real data.

The initialization error. As discussed in Sections 2 and 3, the marginals of both generative processes $\tilde{\mathbf{y}}$ and $\hat{\mathbf{y}}$ following respectively Equation (6) and Equation (3) slightly differs from p_t due to their common white noise initial condition. This implies an error that we call the initialization error. The distance between $(\tilde{p}_t)_{0 \leq t \leq T}, (\hat{p}_t)_{0 \leq t \leq T}$ and $(p_t)_{0 \leq t \leq T}$ can be explicitly studied in the Gaussian case with the following proposition (proved in Appendix B.4).

Proposition 4 (Marginals of the generative processes under Gaussian assumption). *Under Gaussian assumption, $(\tilde{\mathbf{y}}_t)_{0 \leq t \leq T}$ and $(\hat{\mathbf{y}}_t)_{0 \leq t \leq T}$ are Gaussian processes. At each time t , \tilde{p}_t is the Gaussian distribution $\mathcal{N}(\mathbf{0}, \tilde{\Sigma}_t)$ with $\tilde{\Sigma}_t = \Sigma_t + e^{-2(B_T - B_t)} \Sigma_t^2 \Sigma_T^{-1} (\Sigma_T^{-1} - \mathbf{I})$ and \hat{p}_t is the Gaussian distribution $\mathcal{N}(\mathbf{0}, \hat{\Sigma}_t)$ with $\hat{\Sigma}_t = \Sigma_T^{-1} \Sigma_t$. For all $0 \leq t \leq T$, the three covariance matrices $\Sigma_t, \tilde{\Sigma}_t$ and $\hat{\Sigma}_t$ share the same range. Furthermore, for all $0 \leq t \leq T$,*

$$\mathbf{W}_2(\tilde{p}_t, p_t) \leq \mathbf{W}_2(\hat{p}_t, p_t) \quad (17)$$

which shows for $t = 0$ that the SDE sampler is a better sampler than the ODE sampler when the exact score is known.

In practice the initialization error for the SDE and ODE samplers may vary by several orders of magnitude, as shown for the CIFAR-10 example in Figure 1.(a) (solid lines).

The discretization error. The implementation of the SDE and the ODE implies to choose a discrete numerical scheme. We propose to study four different schemes presented in Table 1. The classical

SDE schemes	Euler-Maruyama (EM)	$\begin{cases} \hat{\mathbf{y}}_0^{\Delta, \text{EM}} \sim \mathcal{N}_0 \\ \hat{\mathbf{y}}_{k+1}^{\Delta, \text{EM}} = \hat{\mathbf{y}}_k^{\Delta, \text{EM}} + \Delta_t \beta_{T-t_k} \left(\hat{\mathbf{y}}_k^{\Delta, \text{EM}} - 2\boldsymbol{\Sigma}_{T-t_k}^{-1} \hat{\mathbf{y}}_k^{\Delta, \text{EM}} \right) + \sqrt{2\Delta_t \beta_{T-t_k}} \mathbf{z}_k, \mathbf{z}_k \sim \mathcal{N}_0 \end{cases} \quad (19)$
	Exponential integrator (EI)	$\begin{cases} \hat{\mathbf{y}}_0^{\Delta, \text{EI}} \sim \mathcal{N}_0 \\ \hat{\mathbf{y}}_{k+1}^{\Delta, \text{EI}} = \hat{\mathbf{y}}_k^{\Delta, \text{EI}} + \gamma_{1,k} \left(\hat{\mathbf{y}}_k^{\Delta, \text{EI}} - 2\boldsymbol{\Sigma}_{T-t_k}^{-1} \hat{\mathbf{y}}_k^{\Delta, \text{EI}} \right) + \sqrt{2\gamma_{2,k}} \mathbf{z}_k, \mathbf{z}_k \sim \mathcal{N}_0 \end{cases} \quad (20)$ <p>where $\gamma_{1,k} = \exp(B_{T-t_k} - B_{T-t_{k+1}}) - 1$ and $\gamma_{2,k} = \frac{1}{2}(\exp(2B_{T-t_k} - 2B_{T-t_{k+1}}) - 1)$</p>
ODE schemes	Explicit Euler	$\begin{cases} \hat{\mathbf{y}}_0^{\Delta, \text{Euler}} \sim \mathcal{N}_0 \\ \hat{\mathbf{y}}_{k+1}^{\Delta, \text{Euler}} = \hat{\mathbf{y}}_k^{\Delta, \text{Euler}} + \Delta_t f(t_k, \hat{\mathbf{y}}_k^{\Delta, \text{Euler}}) \quad \text{with } f(t, \mathbf{y}) = \beta_{T-t} \mathbf{y} - \beta_{T-t} \boldsymbol{\Sigma}_{T-t}^{-1} \mathbf{y} \end{cases} \quad (21)$
	Heun's method	$\begin{cases} \hat{\mathbf{y}}_0^{\Delta, \text{Heun}} \sim \mathcal{N}_0 \\ \hat{\mathbf{y}}_{k+1/2}^{\Delta, \text{Heun}} = \hat{\mathbf{y}}_k^{\Delta, \text{Heun}} + \Delta_t f(t_k, \hat{\mathbf{y}}_k^{\Delta, \text{Heun}}) \quad \text{with } f(t, \mathbf{y}) = \beta_{T-t} \mathbf{y} - \beta_{T-t} \boldsymbol{\Sigma}_{T-t}^{-1} \mathbf{y} \\ \hat{\mathbf{y}}_{k+1}^{\Delta, \text{Heun}} = \hat{\mathbf{y}}_k^{\Delta, \text{Heun}} + \frac{\Delta_t}{2} \left(f(t_k, \hat{\mathbf{y}}_k^{\Delta, \text{Heun}}) + f(t_{k+1}, \hat{\mathbf{y}}_{k+1/2}^{\Delta, \text{Heun}}) \right) \end{cases} \quad (22)$

Table 1: **Stochastic and deterministic discretization schemes.** EM and EI discretize the backward SDE (3), Euler and Heun schemes discretize of the probability flow ODE (6).

Euler-Maruyama (EM) is used in [26] and the exponential integrator (EI) in [5] to sample from the SDE (3). The Euler method is the simplest ODE solver and Heun's scheme is recommended in [15] to model the ODE (6). Under Gaussian assumption, the eigenvalues of the covariance matrix can be computed numerically recursively for each scheme to evaluate the Wasserstein distance. More precisely, all the covariance matrices are diagonalizable in the diagonalization basis of $\boldsymbol{\Sigma}$, and thus Equation (16) is valid. For example, denoting $(\lambda_i^t)_{1 \leq i \leq d}$ the eigenvalues of $\boldsymbol{\Sigma}_t$ and $(\lambda_i^{\text{EM},k})_{1 \leq i \leq d}$ the eigenvalues of the covariance matrix of the Euler-Maruyama discretization of the SDE at the k th step, $1 \leq k \leq N-1$, the relation verified by these eigenvalues is

$$\lambda_i^{\text{EM},k+1} = \left(1 + \Delta_t \beta_{T-t_k} \left(1 - \frac{2}{\lambda_i^{T-t_k}} \right) \right)^2 \lambda_i^{\text{EM},k} + 2\Delta_t \beta_{T-t_k}, 1 \leq i \leq d, 0 \leq k \leq N-2 \quad (18)$$

with initialization $\lambda_i^{0,\text{EM}} = 1, 1 \leq i \leq d$. For each scheme, we recursively compute the eigenvalues at each time discretization and present the observed Wasserstein distance in Figure 1.(a). We can observe that Heun's method provide the lower Wasserstein distance, followed by EM, EI and the Euler scheme. Note that the discrete schemes does not preserve the range of the covariance matrix, contrary to the continuous formulas. This explains the fact that the Wasserstein distance increases at the final step.

The truncation error. As discussed in [26], it is preferable to study the backward process on $[\varepsilon, T]$ instead of $[0, T]$ because the score is a priori not defined for $t = 0$, which occurs in our case if $\boldsymbol{\Sigma}$ is not invertible. This approximation is called the truncation error. As a consequence, even without error initialization, the backward process leads to p_ε and not p_0 . Under Gaussian assumption, it is possible to explicit this error with the expression given in Proposition 3 and 2 as done in Figure 5.(b) for both continuous and numerical solutions. For the standard practice truncation time $\varepsilon = 10^{-3}$ [26, 15], all numerical schemes have an error close to the corresponding continuous solution. Using a lower ε value is only relevant for the continuous SDE solution.

Ablation study. We propose in Table 2 an ablation study to monitor the magnitude of each error and their accumulation for various sampling schemes for the CIFAR-10 example. In accordance with Proposition 4, the initialization error influences the ODE schemes, while SDE schemes are not affected. Schemes having a sufficient number of steps are not sensitive to the truncation error for $\varepsilon < 10^{-2}$. The discretization error is the more important approximation but it becomes very low for a sufficient number of steps. The lower Wasserstein error is provided by Heun's method with 1000 steps, $\varepsilon = 10^{-5}$. As [15], our conclusions lead to the choice of Heun's scheme as the go-to method.

Influence of eigenvalues. The above observations and conclusions are observed on the CIFAR-10 Gaussian. However, in general, they depend on the eigenvalues of the covariance matrix $\boldsymbol{\Sigma}$. Indeed, as seen in Equation (16), the Wasserstein distance is separable and each eigenvalue contributes to increase it. In Figure 2, we evaluate the contribution of each eigenvalue by plotting $\lambda \mapsto |\sqrt{\lambda} - \sqrt{\lambda^{\text{scheme}}}|$ for each scheme. Figure 2.(a) demonstrates that for the continuous equations, the error increases with the eigenvalues except for a strong decrease for $\lambda = 1$. Besides, as proved in the proof of Proposition 4 (see Appendix B.4), the error for the SDE is always lower than the error for the ODE.

		Continuous		N = 50		N = 250		N = 500		N = 1000	
		p_T	\mathcal{N}_0	p_T	\mathcal{N}_0	p_T	\mathcal{N}_0	p_T	\mathcal{N}_0	p_T	\mathcal{N}_0
EM	$\varepsilon = 0$	0	6.7E-4	4.77	4.77	0.65	0.65	0.31	0.31	0.15	0.16
	$\varepsilon = 10^{-5}$	2.5E-3	2.6E-3	4.77	4.77	0.65	0.65	0.31	0.31	0.16	0.16
	$\varepsilon = 10^{-3}$	0.17	0.17	4.67	4.67	0.69	0.69	0.39	0.39	0.27	0.27
	$\varepsilon = 10^{-2}$	1.35	1.35	4.56	4.56	1.69	1.69	1.50	1.50	1.42	1.42
EI	$\varepsilon = 0$	0	6.7E-4	2.81	2.81	0.57	0.57	0.30	0.30	0.16	0.16
	$\varepsilon = 10^{-5}$	2.5E-3	2.6E-3	2.81	2.81	0.57	0.57	0.30	0.30	0.16	0.16
	$\varepsilon = 10^{-3}$	0.17	0.17	2.91	2.91	0.66	0.66	0.41	0.41	0.28	0.28
	$\varepsilon = 10^{-2}$	1.35	1.35	3.93	3.93	1.76	1.76	1.55	1.55	1.45	1.45
Euler	$\varepsilon = 0$	0	0.07	1.72	1.78	0.38	0.44	0.19	0.26	0.10	0.17
	$\varepsilon = 10^{-5}$	2.5E-3	0.07	1.72	1.78	0.38	0.44	0.20	0.26	0.10	0.17
	$\varepsilon = 10^{-3}$	0.17	0.19	1.72	1.78	0.42	0.48	0.27	0.32	0.21	0.25
	$\varepsilon = 10^{-2}$	1.35	1.36	2.21	2.25	1.41	1.43	1.37	1.38	1.36	1.37
Heun	$\varepsilon = 0$	0	0.07	7.09	7.09	0.72	0.73	0.21	0.22	0.05	0.09
	$\varepsilon = 10^{-5}$	2.5E-3	0.07	6.48	6.48	0.64	0.65	0.18	0.20	0.05	0.09
	$\varepsilon = 10^{-3}$	0.17	0.19	0.56	0.57	0.13	0.15	0.16	0.18	0.17	0.19
	$\varepsilon = 10^{-2}$	1.35	1.36	1.37	1.38	1.35	1.36	1.35	1.36	1.35	1.36

Table 2: **Ablation study of Wasserstein errors for the CIFAR-10 Gaussian.** For a given discretization scheme, the table presents the Wasserstein distance associated with the truncation error for different values of ε . The columns p_T and \mathcal{N}_0 show the influence of the initialization error. The continuous column corresponds to the continuous SDE or ODE linked with the scheme (identical values for EM, EI and Euler, Heun).

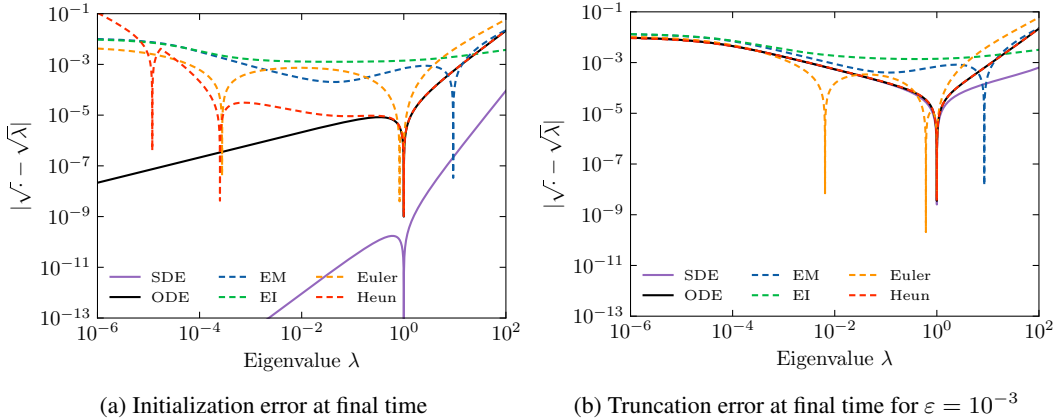


Figure 2: **Eigenvalue contribution to the Wasserstein error.** The magnitude of the Wasserstein error is influenced by the eigenvalues of the covariance of the Gaussian distribution. Left: Contribution to the Wasserstein error for the continuous equations and the discretization schemes with standard initialization \mathcal{N}_0 . Right: Same plot when using a truncation time $\varepsilon = 10^{-3}$. All schemes use $N = 1000$ steps. While we prove that the continuous SDE is always better than the continuous ODE (Proposition 4), it is not the same for the discrete schemes. With a truncation time $\varepsilon = 10^{-3}$ (b), Heun’s method is nearly as good as the continuous ODE solution for all eigenvalues, which shows it is well-adapted to any Gaussian distribution.

Unfortunately, once discretized the stochastic schemes are not as good as the continuous solutions. The EI scheme is the more stable along the range of eigenvalues but in the end it is in general more costly than the others in terms of Wasserstein error. Without truncation time, Heun’s method fails for low eigenvalues because Σ is not stably invertible. However, as seen in Figure 2.(b), with a truncation time $\varepsilon = 10^{-3}$, Heun’s method is very close to the continuous ODE solution. This shows that for any Gaussian distribution Heun’s method introduces nearly no additional discretization error, making this scheme the one to favor in practice.

5 Numerical study of the score approximation

So far our theoretical and numerical study has been conducted under the hypothesis that the score function is known, thus discarding the evaluation of the score approximation. In practice, for general data distribution, the score function is parameterized by a neural network trained using denoising score-matching. This learned score function is not perfect and while theoretical studies assume the network to be close to the theoretical one (with uniform or adaptative bounds, see the discussion in [5]), such an hypothesis is hard to check in practice, especially in our non compact setting. Thus, we propose in this section to train a diffusion models on a Gaussian distribution and evaluate numerically the impact of the score approximation.

The Gaussian ADSN distribution for microtextures. So far our running example was the CIFAR-10 Gaussian but we will now turn to another example that produces visually interesting images, namely Gaussian micro-textures. We consider the asymptotic discrete spot noise (ADSN) distribution [11] associated with an RGB texture $\mathbf{u} \in \mathbb{R}^{3 \times M \times N}$ which is defined as the stationary Gaussian distribution that has covariance equal the autocorrelation of \mathbf{u} . More precisely, this distribution is sampled using convolution with a white Gaussian noise [11]: Denoting $m \in \mathbb{R}^3$ the chanelwise mean of \mathbf{u} and $\mathbf{t}_c = \frac{1}{\sqrt{MN}}(\mathbf{u}_c - m_c)$, $1 \leq c \leq 3$, its associated *texton*, for $\mathbf{w} \sim \mathcal{N}_0$ of size $M \times N$ the chanelwise convolution $\mathbf{x} = m + \mathbf{t} \star \mathbf{w} \in \mathbb{R}^{3 \times M \times N}$ follows $\text{ADSN}(\mathbf{u})$. This distribution is the Gaussian $\mathcal{N}(m, \Sigma)$. To deal with zero mean Gaussian, adding the mean m is considered as a post-processing to visualize samples and we study $\mathcal{N}(\mathbf{0}, \Sigma)$. The matrix Σ is a well-known convolution matrix [9], its eigenvectors and associated eigenvalues can be computed in the Fourier domain, as done in Appendix E.2. Σ admits the eigenvalues $\lambda_1^{\xi, \text{ADSN}} = |\hat{\mathbf{t}}_1|^2(\xi) + |\hat{\mathbf{t}}_2|^2(\xi) + |\hat{\mathbf{t}}_3|^2(\xi)$, $\xi \in \mathbb{R}^{M \times N}$ and 0 with multiplicity $2MN$ and we can conduct the same analysis as before (see Appendix D). To evaluate if a set of N_{samples} sampled images is close to the ADSN distribution p_{data} , we evaluate a problem-specific empirical Wasserstein distance: Supposing that the N_{samples} are drawn from a Gaussian distribution $p^{\text{emp.}} = \mathcal{N}(\mathbf{0}, \Gamma)$ such that Γ admits the same eigenvectors as Σ , we compute

$$\mathbf{W}_2^{\text{emp.}}(p^{\text{emp.}}, p_{\text{data}}) = \sqrt{\sum_{\xi \in \mathbb{R}^{3M \times N}} \left(\sqrt{\lambda_1^{\xi, \text{emp.}}} - \sqrt{\lambda_1^{\xi, \text{ADSN}}} \right)^2 + \lambda_2^{\xi, \text{emp.}} + \lambda_3^{\xi, \text{emp.}}} \quad (23)$$

where $(\lambda_i^{\xi, \text{emp.}})_{\xi \in \mathbb{R}^{M \times N}, 1 \leq i \leq 3}$ are estimators of the eigenvalues of Γ given in Appendix E.3.

Learning the score function. We train the network using the code² associated with the paper [26]. We choose the architecture of DDPM, which is a U-Net described in [14], with the parameters proposed for the dataset CelebaHQ256 to deal with the 256×256 ADSN model associated with the top-left image of Figure 3. We use the training procedure corresponding to DDPM cont. in [26]. β is linear from 0.05 to 10 with $T = 1$. We train over 1.3M iterations, and we generate at each iteration a new batch of ADSN samples. We implement the stochastic EM and deterministic Heun schemes replacing the exact score by its learned version with $N = 1000$ steps and a truncation time $\varepsilon = 10^{-3}$. We name p_θ^{EM} and p_θ^{Heun} , the corresponding distributions and present samples in Figure 3. Both distributions accumulate the four error types.

Evaluation of the score approximation. It is not possible to compute theoretically the Wasserstein distance between $p_{\text{data}} = \text{ADSN}(\mathbf{u})$ and $p_\theta^{\text{EM}}, p_\theta^{\text{Heun}}$ due to the non-linearity of the learned score. To compute an empirical Wasserstein error between it, we use Equation (23). Let us precise that this approximation underestimates the real Wasserstein distance since it wrongly assumes that the distributions $p_\theta^{\text{EM}}, p_\theta^{\text{Heun}}$ are Gaussian with a covariance matrix diagonalizable in the same basis than the covariance matrix Σ of $\text{ADSN}(\mathbf{u})$. We complete this dedicated empirical measure with the standard FID. This metrics are reported in Table 3 where for theoretical distributions that are fast to sample we add the standard deviations computed on 25 different $50k$ -samplings. For this Gaussian distribution, the score approximation is by far the most impactful source of error. We observe that the stochastic EM sampling is more resilient to score approximation than the deterministic Heun’s scheme, resulting in out-of-distribution samples (Figure 3). We may explain this behavior by recalling the results of Proposition 4 that shows that SDE solutions are less sensitive to initialization errors than ODE. Indeed, adding noise at each iteration tends to mitigate the accumulated errors, and score approximation may be consider as some initialization error ocuring at each step.

²Code available at https://github.com/yang-song/score_sde_pytorch

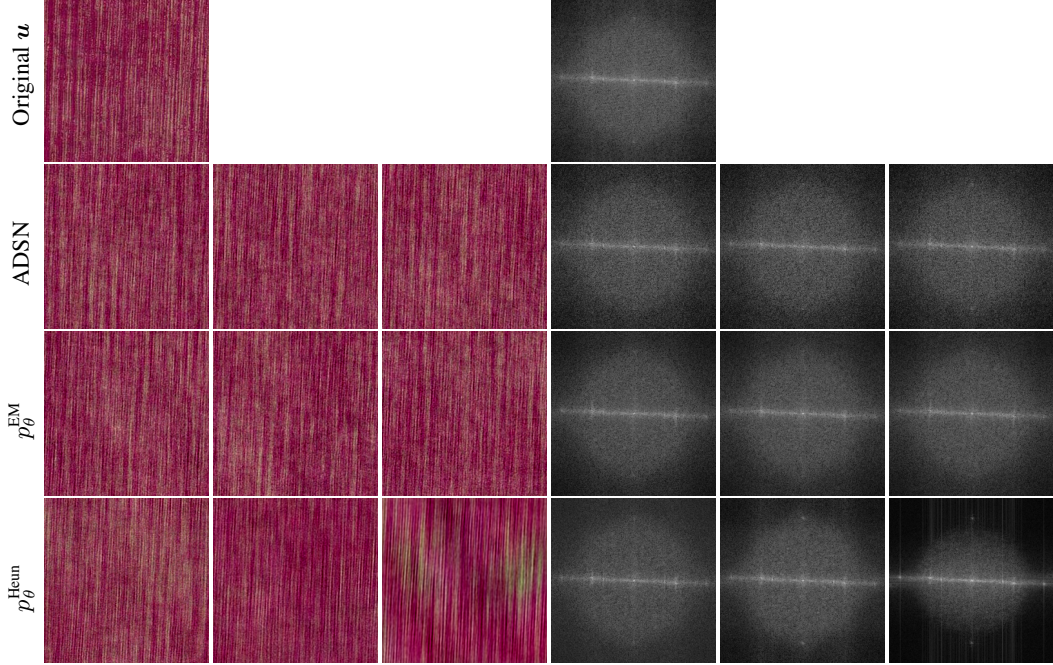


Figure 3: **Texture samples generated with the learned score.** First row: original image u and its DFT modulus (for all DFT modulus we display the sum of the DFT modulus of the three color chanel and apply a logarithmic contrast change). Second row: three samples of $\text{ADSN}(u)$ with their associated DFT moduli. Third and fourth row: Samples generated with the learned score with EM and Heun’s discretization schemes and their associated DFT moduli. While both schemes use the same learned score function, the generation with Heun’s scheme can produce out-of-distribution samples, as seen with the third sample.

p	Exact score distribution			Learned score distribution	
	$W_2(p, p_{\text{data}}) \downarrow$	$W_2^{\text{emp.}}(p^{\text{emp.}}, p_{\text{data}}) \downarrow$	$\text{FID}(p^{\text{emp.}}, p_{\text{data}}) \downarrow$	$W_2^{\text{emp.}}(p_{\theta}^{\text{emp.}}, p_{\text{data}}) \downarrow$	$\text{FID}(p_{\theta}^{\text{emp.}}, p_{\text{data}}) \downarrow$
EM	5.16	$5.1630 \pm 7\text{E-}5$	$0.0891 \pm 8\text{E-}4$	15.6	1.02
Heun	3.73	$3.7323 \pm 2\text{E-}4$	$0.0447 \pm 6\text{E-}4$	56.7	19.4

Table 3: **Numerical evaluation of the score approximation for a Gaussian microtexture model.** For two schemes, the EM discretization of the backward SDE and Heun’s method associated with the flow ODE, the table shows the Wasserstein distance and FID for theoretical and learned distributions. The theoretical W_2 value is computed with explicit formulas, as done in Table 4. The FID and empirical W_2 w.r.t the theoretical distribution are computed on 25 samplings of 50k images while only one sampling of 50k images is drawn for the parametric distributions (to limit computation time).

6 Conclusion

By restricting the analysis of diffusion models to the specific case of Gaussian distributions, we were able to derive exact solutions for both the backward SDE and its associated probability flow ODE. Additionally, we characterized the discrete Gaussian processes arising when discretizing these equations. This allowed us to provide exact Wasserstein errors for the initialization error, the discretization error, and the truncation error as well as any of their combinations. This theoretical analysis led to conclude that Heun’s scheme is the best numerical solution, in accordance with empirical previous work [15].

To conclude our work we conducted an empirical analysis with a learned score function using standard architecture which showed that the score approximation error may be the most important one in practice. This suggests that assessing the quality of learned score functions is an important research direction for future work.

Acknowledgements: The authors acknowledge the support of the project MISTIC (ANR-19-CE40-005).

References

- [1] Joe Benton, Valentin De Bortoli, Arnaud Doucet, and George Deligiannidis. Nearly \mathcal{L}_1 -linear convergence bounds for diffusion models via stochastic localization. In *The Twelfth International Conference on Learning Representations*, 2024.
- [2] Minshuo Chen, Kaixuan Huang, Tuo Zhao, and Mengdi Wang. Score approximation, estimation and distribution recovery of diffusion models on low-dimensional data. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 4672–4712. PMLR, 23–29 Jul 2023.
- [3] Sitan Chen, Sinho Chewi, Holden Lee, Yuanzhi Li, Jianfeng Lu, and Adil Salim. The probability flow ODE is provably fast. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [4] Sitan Chen, Sinho Chewi, Jerry Li, Yuanzhi Li, Adil Salim, and Anru Zhang. Sampling is as easy as learning the score: theory for diffusion models with minimal data assumptions. In *The Eleventh International Conference on Learning Representations*, 2023.
- [5] Valentin De Bortoli. Convergence of denoising diffusion models under the manifold hypothesis. *Transactions on Machine Learning Research*, 2022.
- [6] Valentin De Bortoli, James Thornton, Jeremy Heng, and Arnaud Doucet. Diffusion schrödinger bridge with applications to score-based generative modeling. In *Advances in Neural Information Processing Systems*, volume 34, pages 17695–17709. Curran Associates, Inc., 2021.
- [7] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat GANs on image synthesis. In *Advances in Neural Information Processing Systems*, volume 34, pages 8780–8794. Curran Associates, Inc., 2021.
- [8] D.C Dowson and B.V Landau. The fréchet distance between multivariate normal distributions. *Journal of Multivariate Analysis*, 12(3):450–455, 1982.
- [9] Sira Ferradans, Gui-Song Xia, Gabriel Peyré, and Jean-François Aujol. Static and dynamic texture mixing using optimal transport. In *Scale Space and Variational Methods in Computer Vision*, 2013.
- [10] Giulio Franzese, Simone Rossi, Lixuan Yang, Alessandro Finamore, Dario Rossi, Maurizio Filippone, and Pietro Michiardi. How much is enough? a study on diffusion times in score-based generative models. *Entropy*, 25(4), 2023.
- [11] Bruno Galerne, Yann Gousseau, and Jean-Michel Morel. Random Phase Textures: Theory and Synthesis. *IEEE Transactions on Image Processing*, 20(1):257–267, 2011.
- [12] Xuefeng Gao and Lingjiong Zhu. Convergence analysis for general probability flow odes of diffusion models in wasserstein distances. *ArXiv*, abs/2401.17958, 2024.
- [13] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local Nash equilibrium. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [14] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, volume 33, pages 6840–6851. Curran Associates, Inc., 2020.
- [15] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. In *Proc. NeurIPS*, 2022.
- [16] Valentin Khrulkov, Gleb Ryzhakov, Andrei Chertkov, and Ivan Oseledets. Understanding DDPM latent codes through optimal transport. In *The Eleventh International Conference on Learning Representations*, 2023.

- [17] Hugo Lavenant and Filippo Santambrogio. The flow map of the fokker–planck equation does not provide optimal transport. *Applied Mathematics Letters*, 133:108225, 2022.
- [18] Holden Lee, Jianfeng Lu, and Yixin Tan. Convergence of score-based generative modeling for general data distributions. In *NeurIPS 2022 Workshop on Score-Based Methods*, 2022.
- [19] Holden Lee, Jianfeng Lu, and Yixin Tan. Convergence for score-based generative modeling with polynomial complexity. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS ’22, Red Hook, NY, USA, 2024. Curran Associates Inc.
- [20] B. Øksendal. *Stochastic Differential Equations: An Introduction with Applications*. Universitext. Springer Berlin Heidelberg, 2010.
- [21] E. Pardoux. Grossissement d’une filtration et retournement du temps d’une diffusion. In Jacques Azéma and Marc Yor, editors, *Séminaire de Probabilités XX 1984/85*, pages 48–55, Berlin, Heidelberg, 1986. Springer Berlin Heidelberg.
- [22] S. Särkkä and A. Solin. *Applied Stochastic Differential Equations*. Institute of Mathematical Statistics Textbooks. Cambridge University Press, 2019.
- [23] Kulin Shah, Sitan Chen, and Adam Klivans. Learning mixtures of gaussians using the DDPM objective. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [24] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2021.
- [25] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’ Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [26] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021.
- [27] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [28] Li Kevin Wenliang and Ben Moran. Score-based generative model learn manifold-like structures with constrained mixing. In *NeurIPS 2022 Workshop on Score-Based Methods*, 2022.
- [29] Martin Zach, Erich Kobler, Antonin Chambolle, and Thomas Pock. Product of gaussian mixture diffusion models. *Journal of Mathematical Imaging and Vision*, Mar 2024.
- [30] Martin Zach, Thomas Pock, Erich Kobler, and Antonin Chambolle. Explicit diffusion of gaussian mixture model based image priors. In Luca Calatroni, Marco Donatelli, Serena Morigi, Marco Prato, and Matteo Santacesaria, editors, *Scale Space and Variational Methods in Computer Vision*, pages 3–15, Cham, 2023. Springer International Publishing.

A Characterization of Gaussian distributions through diffusion models

The following proposition shows that our Gaussian assumption occurs if and only if the score function is linear.

Proposition 5. *The three following propositions are equivalent:*

- (i) $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \Sigma)$ for some covariance Σ .
- (ii) $\forall t > 0, \nabla_{\mathbf{x}} \log p_t(\mathbf{x})$ is linear w.r.t \mathbf{x} .
- (iii) $\exists t > 0, \nabla_{\mathbf{x}} \log p_t(\mathbf{x})$ is linear w.r.t \mathbf{x} .

In this case, for $t > 0, \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) = -\Sigma_t^{-1} \mathbf{x}$, with Σ_t defined in Proposition 1.

Proof. (ii) \Rightarrow (iii) is clear.

If (i), for $t > 0$, $p_t(\mathbf{x}) = C_t \exp(-\frac{1}{2}\mathbf{x}^T \Sigma_t^{-1} \mathbf{x})$. Consequently, $\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) = -\Sigma_t^{-1} \mathbf{x}$ and (i) \Rightarrow (ii)

If (iii), there exists A such that $\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) = A\mathbf{x}$. Consequently, $p_t(\mathbf{x}) = C_t \exp(-\frac{1}{2}\mathbf{x}^T A\mathbf{x})$ and \mathbf{x}_t is Gaussian. This provides that $\mathbf{x}_0 = e^{B_t} \mathbf{x}_t - \boldsymbol{\eta}_t$ is Gaussian and (iii) \Rightarrow (i). \square

B Proofs of Section 3

B.1 Proposition 1: Solution of the forward SDE

We aim at solving:

$$d\mathbf{x}_t = -\beta_t \mathbf{x}_t dt + \sqrt{2\beta_t} d\mathbf{w}_t, \quad \mathbf{x}_0 \sim p_{\text{data}}. \quad (24)$$

By considering $\mathbf{z}_t = e^{B_t} \mathbf{x}_t$ where $B_t = \int_0^t \beta_s ds$,

$$d\mathbf{z}_t = \beta_t e^{B_t} \mathbf{x}_t + e^{B_t} d\mathbf{x}_t = \beta_t e^{B_t} \mathbf{x}_t + e^{B_t} (-\beta_t \mathbf{x}_t dt + \sqrt{2\beta_t} d\mathbf{w}_t) = \sqrt{2\beta_t} e^{B_t} d\mathbf{w}_t. \quad (25)$$

Consequently, for $0 \leq t \leq T$,

$$\mathbf{z}_t = \mathbf{z}_0 + \int_0^t \sqrt{2\beta_s} e^{B_s} d\mathbf{w}_s, \quad \mathbf{z}_0 = e^{B_0} \mathbf{x}_0 = \mathbf{x}_0 \quad (26)$$

and for $0 \leq t \leq T$,

$$\mathbf{x}_t = e^{-B_t} \mathbf{z}_t = e^{-B_t} \mathbf{x}_0 + e^{-B_t} \int_0^t e^{B_s} \sqrt{2\beta_s} d\mathbf{w}_s = e^{-B_t} \mathbf{x}_0 + \boldsymbol{\eta}_t. \quad (27)$$

By Itô's isometry (see e.g [20]),

$$\text{Var} \left(\int_0^t e^{B_s} \sqrt{2\beta_s} d\mathbf{w}_s \right) = \int_0^t 2\beta_s e^{2B_s} ds = [e^{2B_s}]_0^t = e^{2B_t} - e^{2B_0} = e^{2B_t} - 1 \quad (28)$$

which provides the covariance matrix of $\boldsymbol{\eta}_t$:

$$\text{Cov}(\boldsymbol{\eta}_t) = e^{-2B_t} (e^{2B_t} - 1) \mathbf{I} = (1 - e^{-2B_t}) \mathbf{I}. \quad (29)$$

Because \mathbf{x}_0 and $\boldsymbol{\eta}_t$ are independent, $\Sigma_t = e^{-2B_t} \Sigma + \Sigma_t = e^{-2B_t} \Sigma + (1 - e^{-2B_t}) \mathbf{I}$.

And,

$$d\Sigma_t = -2\beta_t e^{-2B_t} (\Sigma - \mathbf{I}) dt = [-2\beta_t \Sigma_t + 2\beta_t \mathbf{I}] dt = -(2\beta_t \Sigma_t - 2\beta_t \Sigma_T^{-1}) \Sigma_t dt. \quad (30)$$

B.2 Proposition 2: Solution of the ODE probability flow under Gaussian assumption

We aim at solving

$$d\mathbf{y}_t = \beta_{T-t} (\mathbf{y}_t + 2\nabla \log p_{T-t}(\mathbf{y}_t)) dt + \sqrt{2\beta_{T-t}} d\mathbf{w}_t, \quad 0 \leq t \leq T \quad (31)$$

Denoting $C_t = \int_0^t \beta_{T-s} ds$, by considering $\mathbf{z}_t = \Sigma_{T-t}^{-1} e^{C_t} \mathbf{y}_t$,

$$dz_t = e^{Ct} \Sigma_{T-t}^{-1} dy_t - e^{Ct} d[\Sigma^{-1}]_{T-t} y_t dt + \beta_{T-t} z_t dt \quad (32)$$

$$= \left[\Sigma_{T-t}^{-1} e^{Ct} \beta_{T-t} (\mathbf{y}_t - 2\Sigma_{T-t}^{-1} \mathbf{y}_t) - \beta_{T-t} z_t + 2\beta_{T-t} \Sigma_{T-t}^{-1} z_t \right] dt + \sqrt{2\beta_{T-t}} e^{Ct} \Sigma_{T-t}^{-1} d\mathbf{w}_t \quad (33)$$

$$\text{(using Equation (9))} \quad (34)$$

$$= \beta_{T-t} (1 - 2\Sigma_{T-t}^{-1}) z_t dt - \beta_{T-t} z_t dt + 2\beta_{T-t} \Sigma_{T-t}^{-1} z_t dt + e^{Ct} \sqrt{2\beta_{T-t}} \Sigma_{T-t}^{-1} d\mathbf{w}_t \quad (35)$$

$$= \sqrt{2\beta_{T-t}} e^{Ct} \Sigma_{T-t}^{-1} d\mathbf{w}_t. \quad (36)$$

$$(37)$$

Consequently,

$$\mathbf{z}_t = \mathbf{z}_0 + \int_0^t \sqrt{2\beta_{T-s}} e^{Cs} \Sigma_{T-s}^{-1} d\mathbf{w}_s = \Sigma_T^{-1} \mathbf{y}_0 + \int_0^t \sqrt{2\beta_{T-s}} e^{Cs} \Sigma_{T-s}^{-1} d\mathbf{w}_s. \quad (38)$$

And,

$$\mathbf{y}_t = e^{-Ct} \Sigma_{T-t} \mathbf{z}_t = e^{-Ct} \Sigma_{T-t} \Sigma_T^{-1} \mathbf{y}_0 + e^{-Ct} \Sigma_{T-t} \int_0^t \Sigma_{T-s}^{-1} e^{Cs} \sqrt{2\beta_{T-s}} d\mathbf{w}_s. \quad (39)$$

Finally,

$$\mathbf{y}_t = e^{-Ct} \Sigma_{T-t} \Sigma_T^{-1} \mathbf{y}_0 + \boldsymbol{\xi}_t \quad \text{with} \quad \boldsymbol{\xi}_t = e^{-Ct} \Sigma_{T-t} \int_0^t \Sigma_{T-s}^{-1} e^{Cs} \sqrt{2\beta_{T-s}} d\mathbf{w}_s. \quad (40)$$

By the multidimensional Itô's isometry,

$$\text{Cov} \left(\int_0^t \Sigma_{T-s}^{-1} e^{Cs} \sqrt{2\beta_{T-s}} d\mathbf{w}_s \right) = 2 \int_0^t e^{2Cs} \beta_{T-s} \Sigma_{T-s}^{-2} ds. \quad (41)$$

Now, remark that for $A_s = e^{2Cs} \Sigma_{T-s}^{-1}$,

$$dA_s = 2\beta_{T-s} A_s ds - e^{2Cs} d[\Sigma^{-1}]_{T-s} \quad (42)$$

$$= 2\beta_{T-s} A_s ds + e^{2Cs} [-2\beta_{T-s} \Sigma_{T-s}^{-1} + 2\beta_{T-s} \Sigma_{T-s}^{-2}] ds \quad (43)$$

$$= 2e^{2Cs} \beta_{T-s} \Sigma_{T-s}^{-2} ds. \quad (44)$$

$$\text{Cov} \left(\int_0^t \Sigma_{T-s}^{-1} e^{Cs} \sqrt{\beta_{T-s}} d\mathbf{w}_s \right) = \int_0^t dA_s = [A_s]_0^t = e^{2Ct} \Sigma_{T-t}^{-1} - \Sigma_T^{-1}. \quad (45)$$

Finally, $\text{Cov}(\boldsymbol{\xi}_t) = \Sigma_{T-t}^2 (\Sigma_{T-t}^{-1} - e^{-2Ct} \Sigma_T^{-1}) = \Sigma_{T-t} - e^{-2Ct} \Sigma_{T-t}^2 \Sigma_T^{-1}$

We have the final formula considering:

$$C_t = \int_0^t \beta_{T-s} ds = \int_{T-t}^T \beta_x dx = \int_0^T \beta_x dx - \int_0^{T-t} \beta_x dx = B_T - B_{T-t} \quad (46)$$

that provides

$$\text{Cov}(\mathbf{y}_t) = \Sigma_{T-t} + e^{-2(B_T - B_{T-t})} \Sigma_{T-t}^2 \Sigma_T^{-1} (\Sigma_{T-t}^{-1} \text{Cov}(\mathbf{y}_0) \Sigma_T^{-1} \Sigma_{T-t} - \mathbf{I}). \quad (47)$$

In particular, if $\text{Cov}(\mathbf{y}_0)$ and Σ commute,

$$\text{Cov}(\mathbf{y}_t) = \Sigma_{T-t} + e^{-2(B_T - B_{T-t})} \Sigma_{T-t}^2 \Sigma_T^{-1} (\Sigma_T^{-1} \text{Cov}(\mathbf{y}_0) - \mathbf{I}). \quad (48)$$

B.3 Proposition 3: Solution of the ODE probability flow under Gaussian assumption

As done in [16], the matrix $\Sigma_t^{1/2}$ admits a derivative which is $d\left[\Sigma^{1/2}\right]_t = \frac{1}{2}d\Sigma_t\Sigma_t^{-1/2}$ because it is diagonalisable. Let us check that $\mathbf{y}_t = \Sigma_T^{-1/2}\Sigma_{T-t}^{1/2}\mathbf{y}_0$ is solution of the ODE (4),

$$d\mathbf{y}_t = -\Sigma_T^{-1/2}\frac{1}{2}d\Sigma_{T-t}\Sigma_t^{-1/2}\mathbf{y}_0dt \quad (49)$$

$$= \Sigma_T^{-1/2}[\beta_{T-t}\Sigma_{T-t} - 2\beta_{T-t}\mathbf{I}]\Sigma_{T-t}^{-1/2}\mathbf{y}_0dt \quad (\text{using Equation(9)}) \quad (50)$$

$$= [\beta_{T-t} - 2\beta_{T-t}\Sigma_{T-t}^{-1}]\mathbf{y}_tdt \quad (51)$$

$$= [\beta_{T-t} + 2\beta_{T-t}\nabla_{\mathbf{y}}\log p_{T-t}(\mathbf{y}_t)]\mathbf{y}_tdt. \quad (52)$$

Finally,

$$\text{Cov}(\mathbf{y}_t) = \Sigma_T^{-1/2}\Sigma_{T-t}^{1/2}\text{Cov}(\mathbf{y}_0)\Sigma_{T-t}^{-1/2}\Sigma_T^{1/2}. \quad (53)$$

In particular, if $\text{Cov}(\mathbf{y}_0)$ and Σ commute,

$$\text{Cov}(\mathbf{y}_t) = \Sigma_T^{-1}\Sigma_{T-t}\text{Cov}(\mathbf{y}_0). \quad (54)$$

B.4 Proof of Proposition 4

For $0 \leq t \leq T$, denoting $(\lambda_i^t)_{1 \leq i \leq d}$ the eigenvalues of Σ_t , the eigenvalues of $\tilde{\Sigma}_t = \text{Cov}(\tilde{\mathbf{y}}_{T-t})$ are

$$\tilde{\lambda}_i^t = \lambda_i^t + e^{-2(B_T - B_t)}(\lambda_i^t)^2 \frac{1}{\lambda_i^T} \left(\frac{1}{\lambda_i^T} - 1 \right), \quad i = 1, \dots, d. \quad (55)$$

and the eigenvalues of $\hat{\Sigma}_t = \text{Cov}(\hat{\mathbf{y}}_{T-t})$ are

$$\hat{\lambda}_i^t = \frac{\lambda_i^t}{\lambda_i^T}, \quad i = 1, \dots, d. \quad (56)$$

Consequently, $\mathbf{W}_2(p_t, \tilde{p}_t)$ is the sum of the squares of all:

$$\sqrt{\lambda_i^t} - \sqrt{\tilde{\lambda}_i^t} = \sqrt{\lambda_i^t} \left(1 - \sqrt{1 + e^{-2(B_T - B_t)}\lambda_i^t \frac{1}{\lambda_i^T} \left(\frac{1}{\lambda_i^T} - 1 \right)} \right). \quad (57)$$

Similarly, $\mathbf{W}_2(p_t, \hat{p}_t)$ is the sum of the squares of all:

$$\sqrt{\lambda_i^t} - \sqrt{\hat{\lambda}_i^t} = \sqrt{\lambda_i^t} \left(1 - \sqrt{\frac{1}{\lambda_i^T}} \right) \quad (58)$$

$$= \sqrt{\lambda_i^t} \left(1 - \sqrt{1 + \left(\frac{1}{\lambda_i^T} - 1 \right)} \right). \quad (59)$$

Let us now compare individually these differences.

$$\frac{e^{-2(B_T-B_t)}\lambda_i^t\frac{1}{\lambda_i^T}\left(\frac{1}{\lambda_i^T}-1\right)}{\frac{1}{\lambda_i^T}-1} = e^{-2(B_T-B_t)}\frac{\lambda_i^t}{\lambda_i^T} \quad (60)$$

$$= e^{-2(B_T-B_t)}\frac{e^{-2Bt}(\lambda_i-1)+1}{e^{-2B_T}(\lambda_i-1)+1} \quad (61)$$

$$= \frac{(\lambda_i-1)+e^{2Bt}}{(\lambda_i-1)+e^{2B_T}} \quad (62)$$

$$< 1. \quad (63)$$

Case 1: $0 < \lambda_i < 1$ and $t > 0$

In this case, $\lambda_i^T < 1$ and:

$$0 < e^{-2(B_T-B_t)}\lambda_i^t\frac{1}{\lambda_i^T}\left(\frac{1}{\lambda_i^T}-1\right) < \frac{1}{\lambda_i^T}-1. \quad (64)$$

Thus,

$$\left|\sqrt{\lambda_i^t}-\sqrt{\widehat{\lambda}_i^t}\right| = \sqrt{\widehat{\lambda}_i^t}-\sqrt{\lambda_i^t} \quad (65)$$

$$= \sqrt{\lambda_i^t}\left(\sqrt{1+e^{-2(B_T-B_t)}\lambda_i^t\frac{1}{\lambda_i^T}\left(\frac{1}{\lambda_i^T}-1\right)}-1\right) \quad (66)$$

$$< \sqrt{\lambda_i^t}\left(\sqrt{1+\left(\frac{1}{\lambda_i^T}-1\right)}-1\right) \quad (67)$$

$$= \sqrt{\widehat{\lambda}_i^t}-\sqrt{\lambda_i^t} \quad (68)$$

$$= \left|\sqrt{\lambda_i^t}-\sqrt{\widehat{\lambda}_i^t}\right|. \quad (69)$$

Case 2: $\lambda_i = 0$ and $t = 0$.

In this case, for $1 \leq i \leq d$, $\widehat{\lambda}_i^T = \widetilde{\lambda}_i^T = 0$.

Case 3: $\lambda_i = 1$.

In this case, for $1 \leq i \leq d$, $\widehat{\lambda}_i^t = \widetilde{\lambda}_i^t = 1$.

Case 4: $1 < \lambda_i$.

In this case, $\lambda_i^T \geq 1$, and $\frac{e^{-2(B_T-B_t)}\lambda_i^t\frac{1}{\lambda_i^T}\left(\frac{1}{\lambda_i^T}-1\right)}{\frac{1}{\lambda_i^T}-1} = e^{-2(B_T-B_t)}\frac{\lambda_i^t}{\lambda_i^T} < 1$ provides

$$e^{-2(B_T-B_t)}\lambda_i^t\frac{1}{\lambda_i^T}\left(\frac{1}{\lambda_i^T}-1\right) > \frac{1}{\lambda_i^T}-1. \quad (70)$$

Finally,

$$\left| \sqrt{\lambda_i^t} - \sqrt{\tilde{\lambda}_i^t} \right| = \sqrt{\lambda_i^t} - \sqrt{\tilde{\lambda}_i^t} \quad (71)$$

$$= \sqrt{\lambda_i^t} \left(1 - \sqrt{1 + e^{-2(B_T - B_t)} \lambda_i^T \frac{1}{\lambda_i^T} \left(\frac{1}{\lambda_i^T} - 1 \right)} \right) \quad (72)$$

$$< \sqrt{\lambda_i^t} \left(1 - \sqrt{1 + \left(\frac{1}{\lambda_i^T} - 1 \right)} \right) \quad (73)$$

$$= \sqrt{\lambda_i^t} - \sqrt{\hat{\lambda}_i^t} \quad (74)$$

$$= \left| \sqrt{\lambda_i^t} - \sqrt{\hat{\lambda}_i^t} \right|. \quad (75)$$

This case study provides:

$$\mathbf{W}_2(\tilde{p}_t, p_t) \leq \mathbf{W}_2(\hat{p}_t, p_t). \quad (76)$$

C Gaussian CIFAR-10 samples

The Gaussian CIFAR-10 produces unstructured images. A grid of samples is presented in Figure 4.



Figure 4: **CIFAR-10 Gaussian samples.** Samples are generated from the Gaussian distribution fitting the CIFAR-10 dataset.

D Theoretical Wasserstein distance for the ADSN model

As done for the Gaussian CIFAR-10, the Wasserstein errors can be computed for the ADSN model as shown in Figure 5 and Table 4.

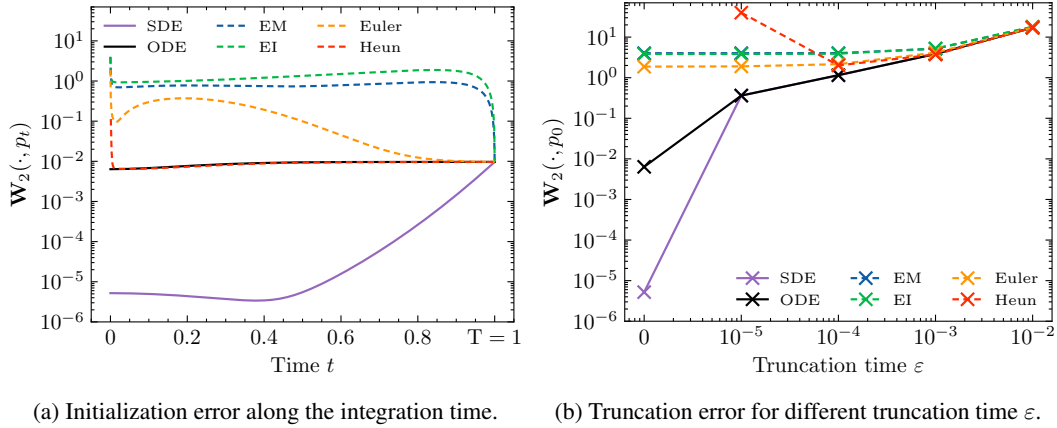


Figure 5: **Wasserstein errors for the diffusion models associated with the Gaussian microtextures.** Left: Evolution of the Wasserstein distance between p_t and the distributions associated with the continuous SDE, the continuous flow ODE and four discrete sampling schemes with standard \mathcal{N}_0 initialization, either stochastic (Euler-Maruyama (EM) and Exponential Integrator (EI)) or deterministic (Euler and Heun). While the continuous SDE is less sensible than the continuous ODE (as proved by Proposition 4), the initialization error impacts all discrete schemes. Heun’s method has the lowest error and is very close to the theoretical ODE, except for the last step that is usually discarded when using time truncation. Right: Wasserstein errors due to time truncation for various truncation times ε . Heun’s scheme is not defined without truncation time due to the zero eigenvalue. Interestingly, for the standard practice truncation time $\varepsilon = 10^{-3}$, all numerical schemes have a comparable error close to their continuous counterparts.

		Continuous		$N = 50$		$N = 250$		$N = 500$		$N = 1000$	
		p_T	\mathcal{N}_0	p_T	\mathcal{N}_0	p_T	\mathcal{N}_0	p_T	\mathcal{N}_0	p_T	\mathcal{N}_0
EM	$\varepsilon = 0$	0	5.2E-6	53.37	53.37	10.58	10.58	6.27	6.27	4.02	4.02
	$\varepsilon = 10^{-5}$	0.36	0.36	53.35	53.35	10.57	10.57	6.26	6.26	4.02	4.02
	$\varepsilon = 10^{-3}$	3.84	3.84	51.92	51.92	10.55	10.55	6.80	6.80	5.16	5.16
	$\varepsilon = 10^{-2}$	17.09	17.09	48.24	48.24	20.39	20.39	18.57	18.57	17.79	17.79
EI	$\varepsilon = 0$	0	5.2E-6	30.91	30.91	8.85	8.85	5.71	5.71	3.84	3.84
	$\varepsilon = 10^{-5}$	0.36	0.36	30.92	30.92	8.85	8.85	5.72	5.72	3.84	3.84
	$\varepsilon = 10^{-3}$	3.84	3.84	31.94	31.94	9.74	9.74	6.76	6.76	5.24	5.24
	$\varepsilon = 10^{-2}$	17.09	17.09	41.49	41.49	21.02	21.02	18.95	18.95	17.99	17.99
Euler	$\varepsilon = 0$	0	6.4E-3	5.69	5.70	3.27	3.27	2.50	2.51	1.87	1.87
	$\varepsilon = 10^{-5}$	0.36	0.36	5.70	5.71	3.28	3.28	2.53	2.53	1.90	1.90
	$\varepsilon = 10^{-3}$	3.84	3.84	6.79	6.79	4.85	4.85	4.41	4.41	4.14	4.14
	$\varepsilon = 10^{-2}$	17.09	17.09	18.52	18.52	17.35	17.35	17.22	17.22	17.15	17.15
Heun	$\varepsilon = 0$	0	6.4E-3	-	-	-	-	-	-	-	-
	$\varepsilon = 10^{-5}$	0.36	0.36	2.4E+3	2.4E+3	3.0E+2	3.0E+2	1.1E+2	1.1E+2	40.00	40.00
	$\varepsilon = 10^{-3}$	3.84	3.84	15.42	15.42	2.25	2.25	3.40	3.40	3.73	3.73
	$\varepsilon = 10^{-2}$	17.09	17.09	16.59	16.59	17.07	17.07	17.09	17.09	17.09	17.09

Table 4: **Ablation study of Wasserstein errors for the Gaussian microtextures.** For a given discretization scheme, the table presents the Wasserstein distance associated with the truncation error for different values of ε . The columns p_T and \mathcal{N}_0 show the influence of the initialization error. The continuous column corresponds to the continuous SDE or ODE linked with the scheme (identical values for EM, EI and Euler, Heun). Note that the Heun scheme is not defined without truncation time due to the zero eigenvalue.

E Study of the covariance matrix of the ADSN distribution

E.1 Reminders on the Discrete Fourier Transform (DFT)

For a given image $\mathbf{v} \in \mathbb{R}^{3 \times M \times N}$, we define the DFT of \mathbf{v} , $\widehat{\mathbf{v}} \in \mathbb{R}^{3 \times M \times N}$ such that for $1 \leq c \leq 3, \xi \in \mathbb{R}^{M \times N}$

$$\widehat{\mathbf{v}}_{c,\xi} = \sum_{x \in M \times N} \mathbf{v}_{c,x} \exp\left(-\frac{2i\pi x_1 \xi_1}{M}\right) \exp\left(-\frac{2i\pi x_2 \xi_2}{N}\right), \quad i^2 = -1 \quad (77)$$

where $\widehat{\mathbf{v}}_{c,\xi}$ is the value of $\widehat{\mathbf{v}}$ at coordinate ξ of the c -th channel of $\widehat{\mathbf{v}}$. For $\mathbf{u} \in \mathbb{R}^{3 \times M \times N}$, by defining $\mathbf{u} \star \mathbf{v}$ the periodic convolution such that for $1 \leq c \leq 3, x \in \mathbb{R}^{M \times N}$:

$$(\mathbf{u} \star \mathbf{v})_{c,x} = \sum_{y \in M \times N} \mathbf{u}_{c,x-y} \mathbf{v}_{c,y} \quad (78)$$

we have:

$$\widehat{\mathbf{u} \star \mathbf{v}} = \widehat{\mathbf{u}} \odot \widehat{\mathbf{v}}, \quad (79)$$

where \odot is the componentwise product.

E.2 Eigenvectors of the covariance matrix of the ADSN distribution

Let $\mathbf{u} \in \mathbb{R}^{3 \times M \times N}$ and its associated texton $\mathbf{t} \in \mathbb{R}^{3 \times M \times N}$. The distribution $\text{ADSN}(\mathbf{u})$ is the Gaussian distribution of $\mathbf{X} = \mathbf{t} \star \mathbf{w}$ such that:

$$\mathbf{X}_i = \mathbf{t}_i \star \mathbf{w} \in \mathbb{R}^{M \times N}, 1 \leq i \leq 3, \mathbf{w} \sim \mathcal{N}_0 \quad (80)$$

Consequently, denoting Σ the covariance of $\text{ADSN}(\mathbf{u})$, for $\mathbf{v} \in \mathbb{R}^{3 \times M \times N}$,

$$\widehat{\Sigma} \mathbf{v}_i = \widehat{\mathbf{t}}_i \widehat{\mathbf{t}}_1 \widehat{\mathbf{v}}_1 + \widehat{\mathbf{t}}_i \widehat{\mathbf{t}}_2 \widehat{\mathbf{v}}_2 + \widehat{\mathbf{t}}_i \widehat{\mathbf{t}}_3 \widehat{\mathbf{v}}_3 = \widehat{\mathbf{t}}_i \left(\widehat{\mathbf{t}}_1 \widehat{\mathbf{v}}_1 + \widehat{\mathbf{t}}_2 \widehat{\mathbf{v}}_2 + \widehat{\mathbf{t}}_3 \widehat{\mathbf{v}}_3 \right) \quad (81)$$

This equation proves that the kernel of Σ contains the kernel of $\mathbf{v} \in \mathbb{R}^{3 \times M \times N} \mapsto \widehat{\mathbf{t}}_1 \widehat{\mathbf{v}}_1 + \widehat{\mathbf{t}}_2 \widehat{\mathbf{v}}_2 + \widehat{\mathbf{t}}_3 \widehat{\mathbf{v}}_3 \in \mathbb{R}^{M \times N}$ which has a dimension greater than $2MN$. Consequently, 0 is eigenvalue of Σ with multiplicity greater than $2MN$. Furthermore, for $\xi \in \mathbb{R}^{M \times N}$, denoting $\mathbf{u}^{1,\xi}$ such that:

$$\widehat{\mathbf{u}}_i^{1,\xi}(\omega) = \mathbf{1}_{\omega=\xi} \widehat{\mathbf{t}}_i(\omega), 1 \leq i \leq 3, \omega \in \mathbb{R}^{M \times N} \quad (82)$$

we have,

$$\Sigma \mathbf{u}^{1,\xi} = (|\widehat{\mathbf{t}}_1(\xi)|^2 + |\widehat{\mathbf{t}}_2(\xi)|^2 + |\widehat{\mathbf{t}}_3(\xi)|^2) \mathbf{u}^{1,\xi}. \quad (83)$$

Furthermore, the family $(\mathbf{u}^{1,\xi})_{\xi \in M \times N}$ is orthogonal. Thus, the eigenvalues of Σ are $(|\widehat{\mathbf{t}}_1(\xi)|^2 + |\widehat{\mathbf{t}}_2(\xi)|^2 + |\widehat{\mathbf{t}}_3(\xi)|^2)_{\xi \in M \times N}$ and 0 with multiplicity $2MN$.

For $\xi \in \mathbb{R}^{M \times N}$, we denote $\mathbf{u}^{2,\xi}, \mathbf{u}^{3,\xi}$ such that for $\omega \in \mathbb{R}^{M \times N}$:

$$\begin{cases} \widehat{\mathbf{u}}_1^{2,\xi}(\omega) &= -\mathbf{1}_{\omega=\xi} \widehat{\mathbf{t}}_3(\omega) \\ \widehat{\mathbf{u}}_2^{2,\xi}(\omega) &= 0 \\ \widehat{\mathbf{u}}_3^{2,\xi}(\omega) &= \mathbf{1}_{\omega=\xi} \widehat{\mathbf{t}}_1(\omega) \end{cases} \quad (84)$$

$$\begin{cases} \widehat{\mathbf{u}}_1^{3,\xi}(\omega) &= 0 \\ \widehat{\mathbf{u}}_2^{3,\xi}(\omega) &= -\mathbf{1}_{\omega=\xi} \widehat{\mathbf{t}}_3(\omega) \\ \widehat{\mathbf{u}}_3^{3,\xi}(\omega) &= \mathbf{1}_{\omega=\xi} \widehat{\mathbf{t}}_2(\omega) \end{cases} \quad (85)$$

We have

$$\Sigma \mathbf{u}^{2,\xi} = 0. \mathbf{u}^{2,\xi} \quad (86)$$

$$\Sigma \mathbf{u}^{3,\xi} = 0. \mathbf{u}^{3,\xi}. \quad (87)$$

Then, applying the orthonormalization of Gram-Schmidt on each tuple $(\mathbf{u}^{1,\xi}, \mathbf{u}^{2,\xi}, \mathbf{u}^{3,\xi})_{\xi \in \mathbb{R}^{M \times N}}$, we obtain an orthonormal basis in the Fourier domain $(\mathbf{v}^{1,\xi}, \mathbf{v}^{2,\xi}, \mathbf{v}^{3,\xi})_{\xi \in \mathbb{R}^{M \times N}}$ of eigenvectors of Σ . More precisely, for $\xi_1, \xi_2 \in \mathbb{R}^{M \times N}$, $1 \leq j_1, j_2 \leq 3$,

$$\left(\widehat{\mathbf{v}}^{j_1, \xi_1}\right)^T \widehat{\mathbf{v}}^{j_2, \xi_2} = \sum_{\substack{x_1 \in M \times N \\ x_2 \in M \times N}} \widehat{\mathbf{v}}_{x_1}^{j_1, \xi_1} \widehat{\mathbf{v}}_{x_2}^{j_2, \xi_2} \quad (88)$$

$$= \mathbf{1}_{\substack{j_1=j_2 \\ \xi_1=\xi_2}} \quad (89)$$

which is applying the square root of Σ to the white Gaussian noise \mathbf{w} . Furthermore, we can ensure that for $\xi \neq \omega \in \mathbb{R}^{M \times N}$, $1 \leq j \leq 3$, $\widehat{\mathbf{v}}^{j,\xi}(\omega) = 0$ such that only the frequency ξ is active in the Fourier transform of $\mathbf{v}^{j,\xi}$. Consequently, for $\mathbf{w} \in \mathbb{R}^{3M \times N}$,

$$\widehat{\mathbf{w}}^T \mathbf{v}^{j,\xi} = \sum_{1 \leq i \leq 3} \widehat{\mathbf{w}}_i(\xi) \widehat{\mathbf{v}}_i^{j,\xi}(\xi). \quad (90)$$

In particular,

$$\left(\widehat{\mathbf{v}}^{j,\xi}\right)^T = \|\widehat{\mathbf{v}}^{j,\xi}\|^2 = \sum_{1 \leq i \leq 3} \left|\mathbf{v}_i^{j,\xi}(\xi)\right|^2 = 1. \quad (91)$$

E.3 Computation of the empirical Wasserstein error in the ADSN covariance diagonalization basis

Let consider a Gaussian distribution $\mathcal{N}(\mathbf{0}, \Gamma)$ such that there exists $(\lambda_1^\xi, \lambda_2^\xi, \lambda_3^\xi)_{\xi \in \mathbb{R}^{M \times N}}$ such that for all $\xi \in \mathbb{R}^{M \times N}$,

$$\Gamma \mathbf{v}^{j,\xi} = \lambda_j^\xi \mathbf{v}^{j,\xi}, \quad 1 \leq j \leq 3. \quad (92)$$

Let $\mathbf{w} \sim \mathcal{N}_0 \in \mathbb{R}^{3M \times N}$, $(\mathbf{v}^{1,\xi}, \mathbf{v}^{2,\xi}, \mathbf{v}^{3,\xi})_{\xi \in \mathbb{R}^{M \times N}}$ is an orthonormal basis in the Fourier domain such that:

$$\widehat{\mathbf{w}} = \sum_{\xi \in \mathbb{R}^{M \times N}} \left(\left[\widehat{\mathbf{w}}^T \widehat{\mathbf{v}}^{1,\xi}\right] \widehat{\mathbf{v}}^{1,\xi} + \left[\widehat{\mathbf{w}}^T \widehat{\mathbf{v}}^{2,\xi}\right] \widehat{\mathbf{v}}^{2,\xi} + \left[\widehat{\mathbf{w}}^T \widehat{\mathbf{v}}^{3,\xi}\right] \widehat{\mathbf{v}}^{3,\xi} \right) \quad (93)$$

$$(94)$$

A sample drawn from $\mathcal{N}(\mathbf{0}, \Gamma)$ has the same distribution as \mathbf{Y} given by

$$\widehat{\mathbf{Y}} = \sum_{\xi \in \mathbb{R}^{M \times N}} \sqrt{\lambda_1^\xi} \left[\widehat{\mathbf{w}}^T \widehat{\mathbf{v}}^{1,\xi}\right] \widehat{\mathbf{v}}^{1,\xi} + \sum_{\xi \in \mathbb{R}^{M \times N}} \sqrt{\lambda_2^\xi} \left[\widehat{\mathbf{w}}^T \widehat{\mathbf{v}}^{2,\xi}\right] \widehat{\mathbf{v}}^{2,\xi} + \sum_{\xi \in \mathbb{R}^{M \times N}} \sqrt{\lambda_3^\xi} \left[\widehat{\mathbf{w}}^T \widehat{\mathbf{v}}^{3,\xi}\right] \widehat{\mathbf{v}}^{3,\xi}. \quad (95)$$

Note that the three channels of \mathbf{w} are independent. Furthermore, for $1 \leq j \leq 3$

$$\left(\widehat{\mathbf{v}}^{j,\xi}\right)^T \widehat{\mathbf{Y}} = \sqrt{\lambda_1^\xi} \left[\widehat{\mathbf{w}}^T \widehat{\mathbf{v}}^{j,\xi}\right] \left\|\widehat{\mathbf{v}}^{j,\xi}\right\|^2 = \sqrt{\lambda_1^\xi} \left[\widehat{\mathbf{w}}^T \widehat{\mathbf{v}}^{j,\xi}\right] \quad (96)$$

$$\left|\left(\widehat{\mathbf{v}}^{j,\xi}\right)^T \widehat{\mathbf{Y}}\right|^2 = \lambda_j^\xi \left|\widehat{\mathbf{w}}^T \widehat{\mathbf{v}}^{j,\xi}\right|^2 \quad (97)$$

$$\mathbb{E} \left[\left| \left(\widehat{\mathbf{v}}^{j,\xi}\right)^T \widehat{\mathbf{Y}} \right|^2 \right] = \lambda_j^\xi \mathbb{E} \left[\left| \widehat{\mathbf{w}}^T \widehat{\mathbf{v}}^{j,\xi} \right|^2 \right] \quad (98)$$

$$\mathbb{E} \left[\left| \widehat{\mathbf{w}}^T \widehat{\mathbf{v}}^{j,\xi} \right|^2 \right] = \sum_{1 \leq c_1, c_2 \leq 3} \mathbb{E} \left[\widehat{\mathbf{w}}_{c_1}(\xi) \widehat{\mathbf{w}}_{c_2}(\xi) \right] \widehat{\mathbf{v}}_{c_1}^{j,\xi}(\xi) \widehat{\mathbf{v}}_{c_2}(\xi) \text{ by Equation (90)} \quad (99)$$

$$= \sum_{1 \leq c \leq 3} \mathbb{E} \left[\left| \widehat{\mathbf{w}}_c(\xi) \right|^2 \right] \left| \widehat{\mathbf{v}}_c^{j,\xi}(\xi) \right|^2 \text{ because the channels are independent} \quad (100)$$

$$= 3MN \sum_{1 \leq c \leq 3} \left| \widehat{\mathbf{v}}_c^{j,\xi}(\xi) \right|^2 \text{ because } \mathbb{E} \left[\left| \widehat{\mathbf{w}}_c(\xi) \right|^2 \right] = MN \quad (101)$$

$$= 3MN \text{ by Equation (91)}. \quad (102)$$

Finally,

$$\mathbb{E} \left[\left| \left(\widehat{\mathbf{v}}^{j,\xi}\right)^T \widehat{\mathbf{Y}} \right|^2 \right] = 3MN \lambda_1^\xi \quad (103)$$

Finally, for a given sampling $(\mathbf{Y}_k)_{1 \leq k \leq N_{\text{samples}}}$ following the distribution $\mathcal{N}(\mathbf{0}, \Gamma)$, an estimator of λ_j^ξ is:

$$\lambda_j^{\xi, \text{emp.}} = \frac{1}{3N_{\text{samples}}MN} \sum_{k=1}^{N_{\text{samples}}} \left| \left(\widehat{\mathbf{v}}^{j,\xi}\right)^T \widehat{\mathbf{Y}}_k \right|^2. \quad (104)$$

The empirical Wasserstein distance between the Gaussian distribution $\mathcal{N}(\mathbf{0}, \Gamma)$ and the ADSN model with texton \mathbf{t} is:

$$\mathbf{W}_2^{\text{emp.}}(\mathcal{N}^{\text{emp.}}(\mathbf{0}, \Gamma), \text{ADSN}(\mathbf{u})) = \sqrt{\sum_{\xi \in \mathbb{R}^{M \times N}} \left(\left(\sqrt{\lambda_1^{\xi, \text{emp.}}} - \sqrt{\lambda_1^{\xi, \text{ADSN}}} \right)^2 + \lambda_2^{\xi, \text{emp.}} + \lambda_3^{\xi, \text{emp.}} \right)} \quad (105)$$

with $\lambda_1^{\xi, \text{ADSN}} = |\widehat{\mathbf{t}}_1(\xi)|^2 + |\widehat{\mathbf{t}}_2(\xi)|^2 + |\widehat{\mathbf{t}}_3(\xi)|^2$ for $\xi \in \mathbb{R}^{M \times N}$.

Furthermore, the computations can be vectorized by componentwise products in the Fourier domain.