

# Introduction to diffusion models and study of their restriction to the Gaussian case

---

Émile Pierret, supervised by Bruno Galerne

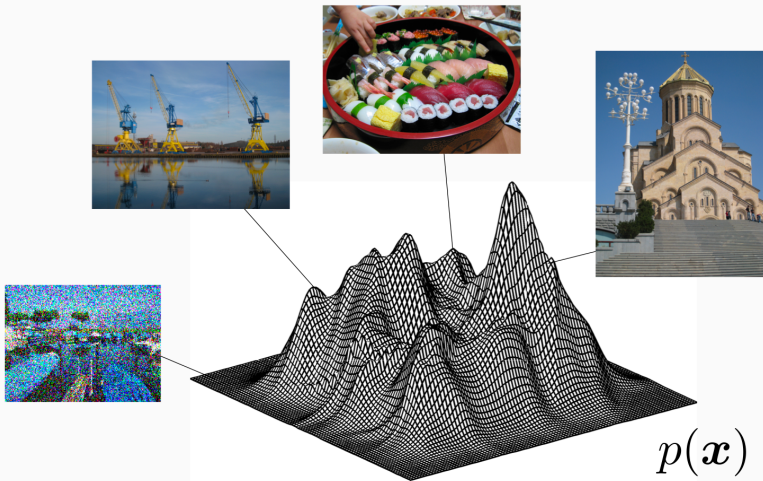
May, 27<sup>th</sup>, CANUM 2024

Institut Denis Poisson – Université d'Orléans, Université de Tours, CNRS

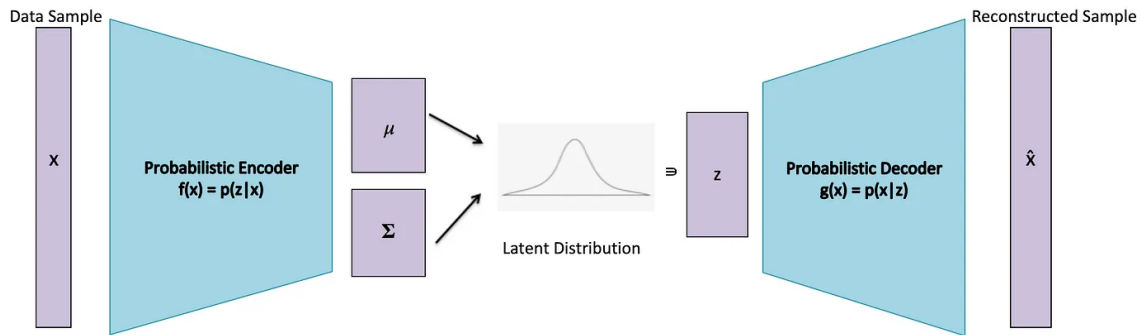
## Introduction to generative models

---

# What is a generative model ?



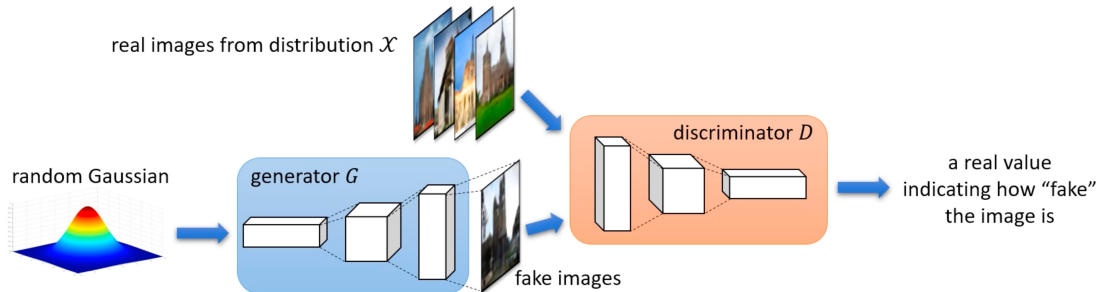
# Variational Auto-Encoder (VAE)



*Image extracted from this url*



# Generative Adversarial Network (GAN)



*Image extracted from this url*

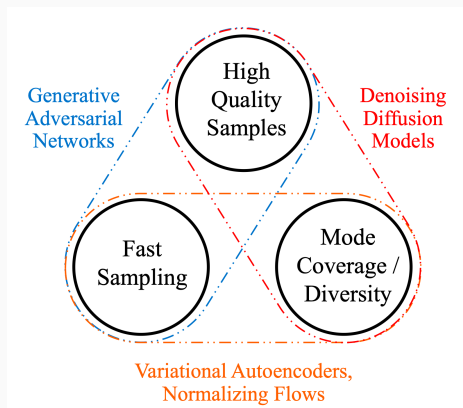
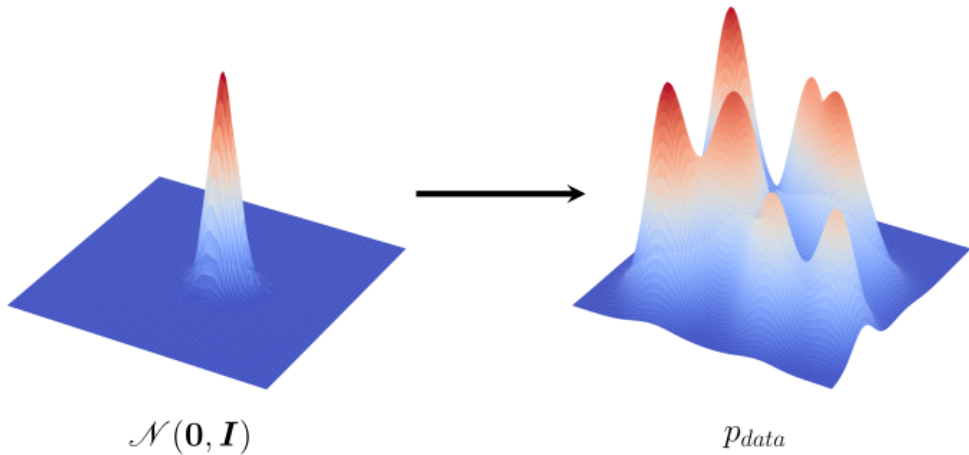


Image extrated from [Xiao et al., 2022]<sup>1</sup>

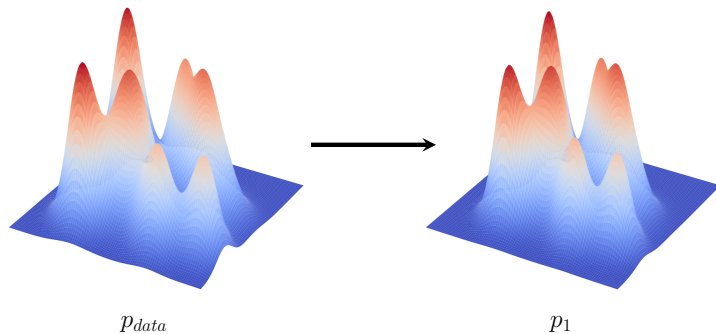
- Dhariwal, P., & Nichol, A. (2021). Diffusion models beat GANs on image synthesis. *Advances in Neural Information Processing Systems*

<sup>1</sup>Xiao, Z., Kreis, K., & Vahdat, A. (2022). Tackling the generative learning trilemma with denoising diffusion GANs. *International Conference on Learning Representations*

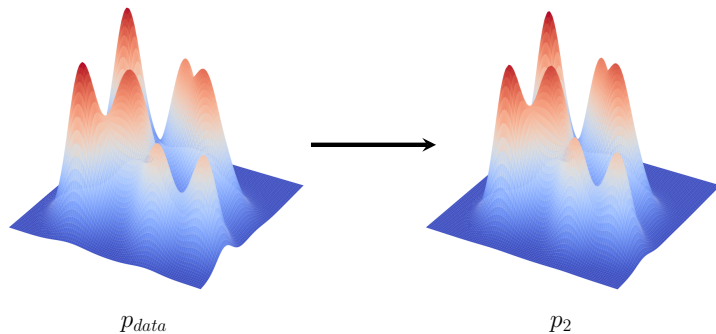
# Main idea



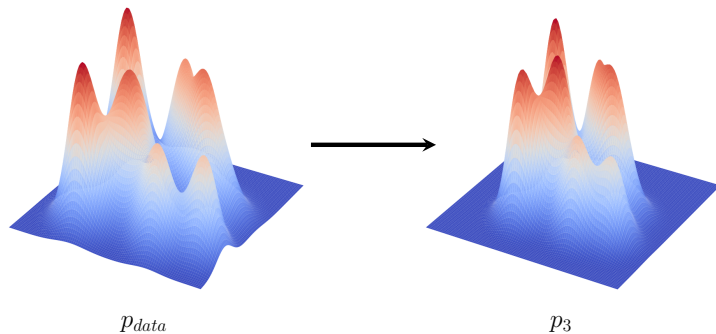
# The forward process of Denoising diffusion models



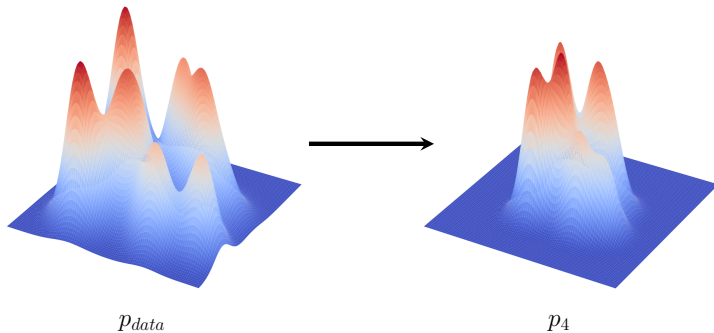
# The forward process of Denoising diffusion models



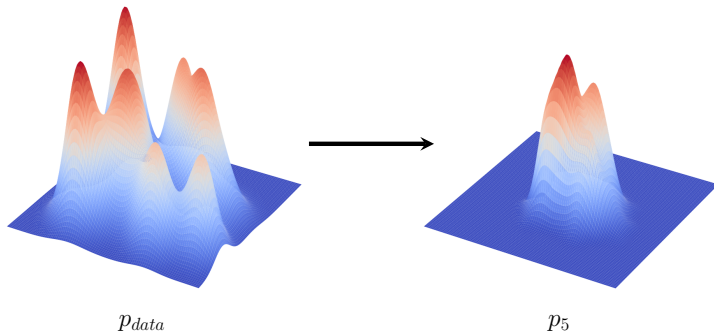
# The forward process of Denoising diffusion models



# The forward process of Denoising diffusion models

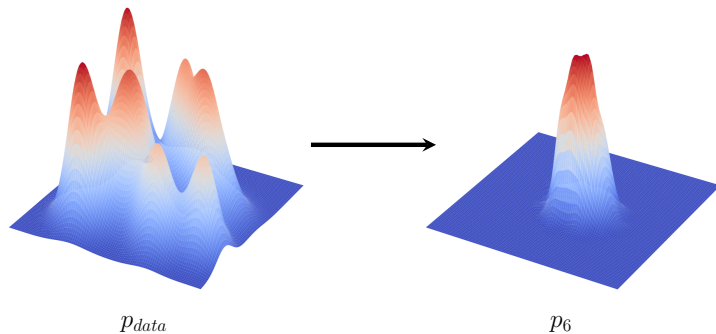


# The forward process of Denoising diffusion models

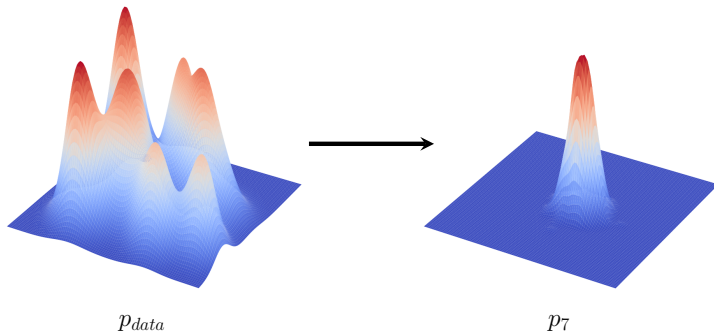




# The forward process of Denoising diffusion models



# The forward process of Denoising diffusion models



## Diffusion models through SDE

---

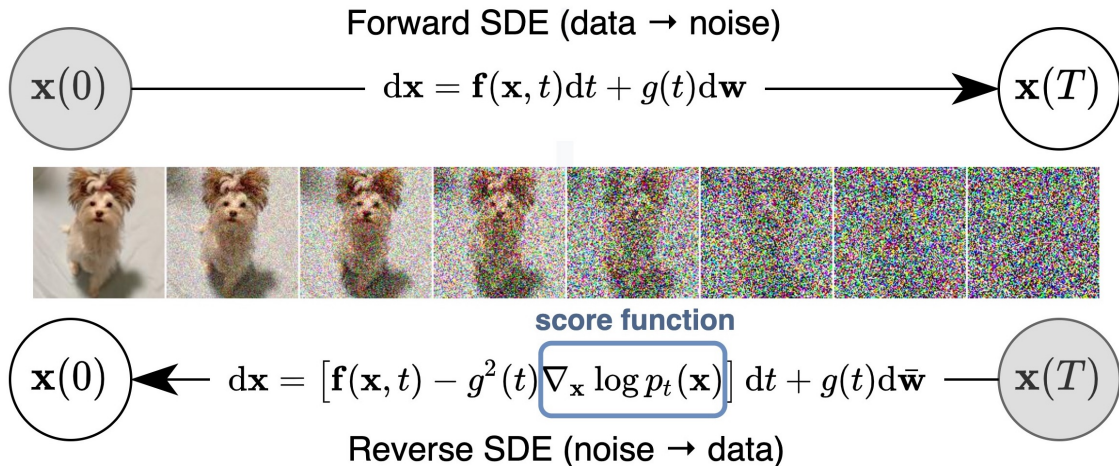


Image extracted from [Song et al., 2021]

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

where  $\beta_t$  is an affine non-decreasing function. We denote  $(p_t)_{0 \leq t \leq T}$  the density of  $\mathbf{x}_t$ .

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

where  $\beta_t$  is an affine non-decreasing function. We denote  $(p_t)_{0 \leq t \leq T}$  the density of  $\mathbf{x}_t$ .

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

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

where  $\beta_t$  is an affine non-decreasing function. We denote  $(p_t)_{0 \leq t \leq T}$  the density of  $\mathbf{x}_t$ .

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

$$\begin{aligned} d\mathbf{z}_t &= \beta_t e^{B_t} \mathbf{x}_t + e^{B_t} d\mathbf{x}_t \\ &= \sqrt{2\beta_t} e^{B_t} d\mathbf{w}_t. \end{aligned}$$

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

where  $\beta_t$  is an affine non-decreasing function. We denote  $(p_t)_{0 \leq t \leq T}$  the density of  $\mathbf{x}_t$ .

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

$$\begin{aligned} d\mathbf{z}_t &= \beta_t e^{B_t} \mathbf{x}_t + e^{B_t} d\mathbf{x}_t \\ &= \sqrt{2\beta_t} e^{B_t} d\mathbf{w}_t. \end{aligned}$$

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

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

with  $\boldsymbol{\eta}_t \sim \mathcal{N}(\mathbf{0}, (1 - e^{-2B_t}) \mathbf{I})$ . In particular,  $\boldsymbol{\Sigma}_t := \text{Cov}(\mathbf{x}_t) = e^{-2B_t} \text{Cov}(\mathbf{x}_0) + (1 - e^{-2B_t}) \mathbf{I}$ .



## The forward process

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

where  $\beta_t$  is an affine non-decreasing function. We denote  $(p_t)_{0 \leq t \leq T}$  the density of  $\mathbf{x}_t$ .

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

$$\begin{aligned} d\mathbf{z}_t &= \beta_t e^{B_t} \mathbf{x}_t + e^{B_t} d\mathbf{x}_t \\ &= \sqrt{2\beta_t} e^{B_t} d\mathbf{w}_t. \end{aligned}$$

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

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

with  $\boldsymbol{\eta}_t \sim \mathcal{N}(\mathbf{0}, (1 - e^{-2B_t}) \mathbf{I})$ . In particular,  $\boldsymbol{\Sigma}_t := \text{Cov}(\mathbf{x}_t) = e^{-2B_t} \text{Cov}(\mathbf{x}_0) + (1 - e^{-2B_t}) \mathbf{I}$ .

Consequently, if  $t \rightarrow +\infty$ ,  $\mathbf{x}_\infty \sim \mathcal{N}_0$

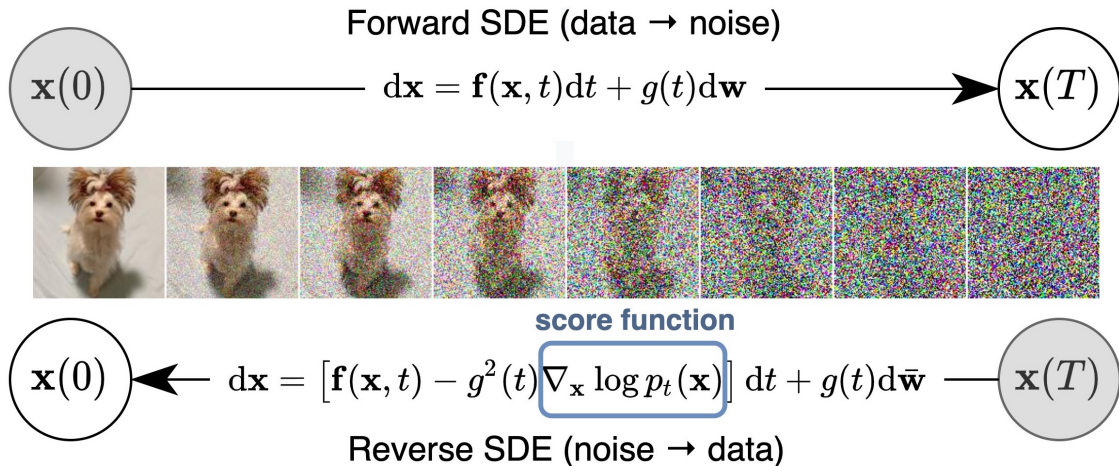


Image extracted from [Song et al., 2021]

Under some assumptions on the distribution  $p_{\text{data}}$  [Pardoux, 1986]<sup>2</sup>, the backward process  $(x_{T-t})_{0 \leq t \leq T}$  verifies the backward SDE

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

---

<sup>2</sup>Pardoux, E. (1986). Grossissement d'une filtration et retournement du temps d'une diffusion. In J. Azéma & M. Yor (Eds.), *Séminaire de probabilités xx 1984/85* (pp. 48–55). Springer Berlin Heidelberg

Under some assumptions on the distribution  $p_{\text{data}}$  [Pardoux, 1986]<sup>2</sup>, the backward process  $(x_{T-t})_{0 \leq t \leq T}$  verifies the backward SDE

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

- $\nabla \log p_{T-t}$  is called the score function.
- The backward Brownian motion  $\bar{w}$  is not defined on the same filtration than the forward  $w$
- We are unable to derive the score function.

---

<sup>2</sup>Pardoux, E. (1986). Grossissement d'une filtration et retournement du temps d'une diffusion. In J. Azéma & M. Yor (Eds.), *Séminaire de probabilités xx 1984/85* (pp. 48–55). Springer Berlin Heidelberg

# How to sample $p_{\text{data}}$ ?

1. Learn the score function  $s_{\theta}(\mathbf{x}, t) \approx \nabla \log p_t(\mathbf{x})$  by applying the forward process to data and minimizing

$$\mathbb{E}_t \left\{ \mathbb{E}_t \lambda(t) \mathbb{E}_{\mathbf{x}_0} \mathbb{E}_{(\mathbf{x}_t | \mathbf{x}_0)} [\|s_{\theta}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t | \mathbf{x}_0)\|^2] \right\}. \quad (4)$$

2. Discretize the backward SDE

- $\mathbf{y}_0 \sim \mathcal{N}_0$  (and not  $p_T$ )
- By Euler Maruyama's scheme,

$$d\mathbf{y}_t = \beta_{T-t}(\mathbf{y}_t + 2\nabla \log p_{T-t}(\mathbf{y}_t))dt + \sqrt{2\beta_{T-t}}d\bar{\mathbf{w}}_t, \quad 0 \leq t < T, \quad \mathbf{y}_0 \sim p_T. \quad (5)$$

becomes:

$$\tilde{\mathbf{y}}_{k+1}^{\Delta, \text{EM}} = \tilde{\mathbf{y}}_k^{\Delta, \text{EM}} + \Delta_t \beta_{T-t_k} \left( \tilde{\mathbf{y}}_k^{\Delta, \text{EM}} - 2\Sigma_{T-t_k}^{-1} \tilde{\mathbf{y}}_k^{\Delta, \text{EM}} \right) + \sqrt{2\Delta_t \beta_{T-t_k}} \mathbf{z}_k, \quad \mathbf{z}_k \sim \mathcal{N}_0 \quad (6)$$

# The flow ODE

With a SDE can be associated an ODE

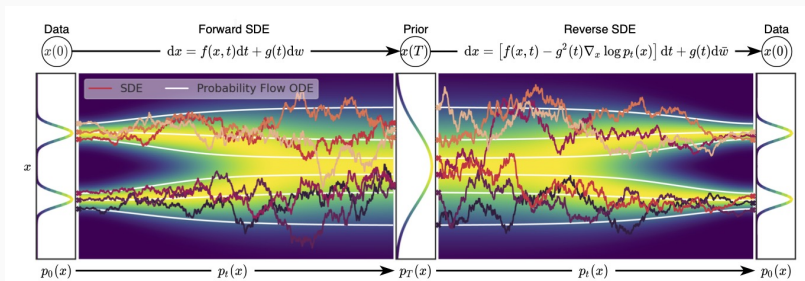


Image extracted from [Song et al., 2021]<sup>3</sup>

<sup>3</sup>Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., & Poole, B. (2021). Score-based generative modeling through stochastic differential equations. *International Conference on Learning Representations*. <https://openreview.net/forum?id=PXTIG12RRHS>

As a reminder, the forward process is.

$$d\mathbf{x}_t = -\beta_t \mathbf{x}_t dt + \sqrt{2\beta_t} d\mathbf{w}_t, \quad 0 \leq t \leq T, \quad \mathbf{x}_0 \sim p_{\text{data}}. \quad (7)$$

With Fokker-Planck equation, we can introduce the associated flow ODE

$$d\mathbf{x}_t = [-\beta_t \mathbf{x}_t - \beta_t \nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t)] dt, \quad 0 < t \leq T, \quad \mathbf{x}_0 \sim p_{\text{data}} \quad (8)$$

such that: if  $\mathbf{y}_0 \sim p_T$  and verifies Equation (9) then for all  $t$ ,  $\mathbf{y}_t \sim p_t$ .

$$d\mathbf{y}_t = [\beta_{T-t} \mathbf{y}_t + \beta_{T-t} \nabla_{\mathbf{y}} \log p_{T-t}(\mathbf{y}_t)] dt, \quad 0 \leq t < T. \quad (9)$$

## Two techniques to sample

1. Learn the score function  $s_\theta(\mathbf{x}, t) \approx \nabla \log p_t(\mathbf{x})$  by applying the forward process.

2. Discretize the backward SDE

- $\mathbf{y}_0 \sim \mathcal{N}_0$  (and not  $p_T$ )
- By Euler Maruyama's scheme,

$$d\mathbf{y}_t = \beta_{T-t}(\mathbf{y}_t + 2\nabla \log p_{T-t}(\mathbf{y}_t))dt + \sqrt{2\beta_{T-t}}d\bar{\mathbf{w}}_t$$

becomes

$$\begin{aligned}\tilde{\mathbf{y}}_{k+1}^{\Delta, \text{EM}} &= \tilde{\mathbf{y}}_k^{\Delta, \text{EM}} + \Delta_t \beta_{T-t_k} \left( \tilde{\mathbf{y}}_k^{\Delta, \text{EM}} - 2\Sigma_{T-t_k}^{-1} \tilde{\mathbf{y}}_k^{\Delta, \text{EM}} \right) \\ &\quad + \sqrt{2\Delta_t \beta_{T-t_k}} \mathbf{z}_k, \quad \mathbf{z}_k \sim \mathcal{N}_0\end{aligned}$$

2. Discretize the flow ODE in reverse-time

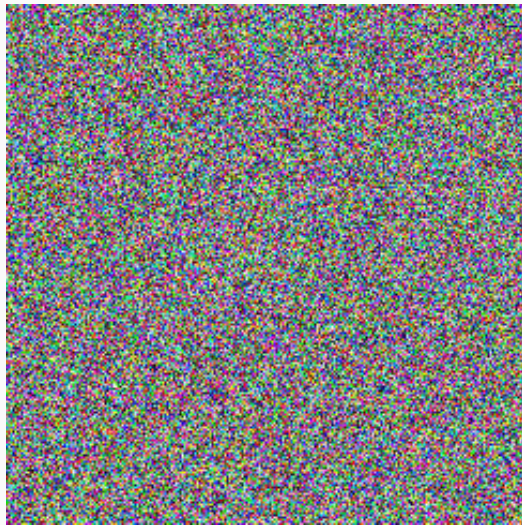
- $\mathbf{y}_0 \sim \mathcal{N}_0$  (and not  $p_T$ )
- By Euler's scheme,

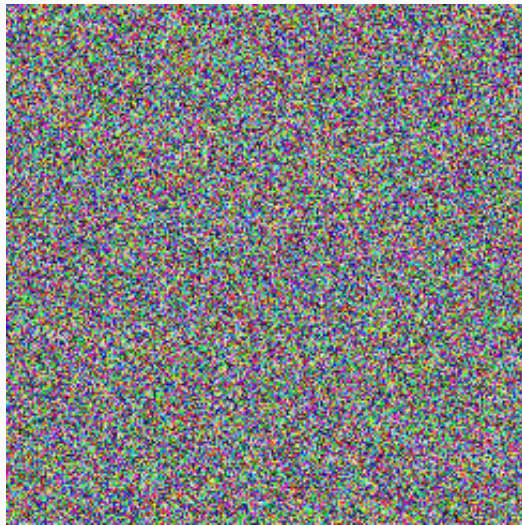
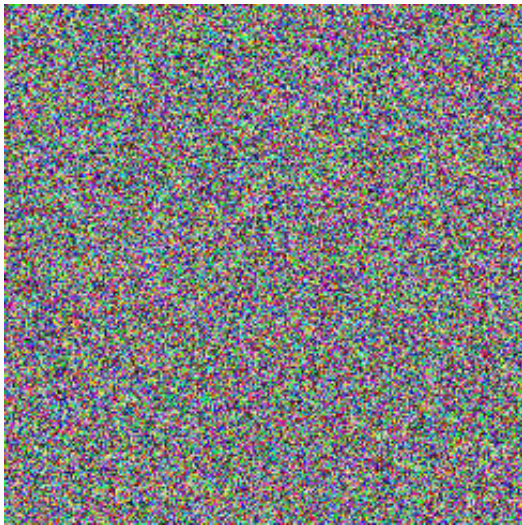
$$d\mathbf{y}_t = [\beta_{T-t}\mathbf{y}_t + \beta_{T-t}\nabla_{\mathbf{y}} \log p_{T-t}(\mathbf{y}_t)] dt$$

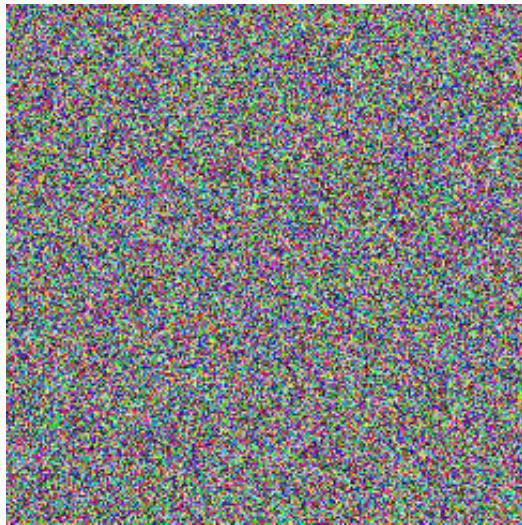
becomes

$$\begin{aligned}\hat{\mathbf{y}}_{k+1}^{\Delta, \text{Euler}} &= \hat{\mathbf{y}}_k^{\Delta, \text{Euler}} + \Delta_t f(t_k, \hat{\mathbf{y}}_k^{\Delta, \text{Euler}}) \\ \text{with } f(t, \mathbf{y}) &= \beta_{T-t}\mathbf{y} - \beta_{T-t}\Sigma_{T-t}^{-1}\mathbf{y}\end{aligned}$$



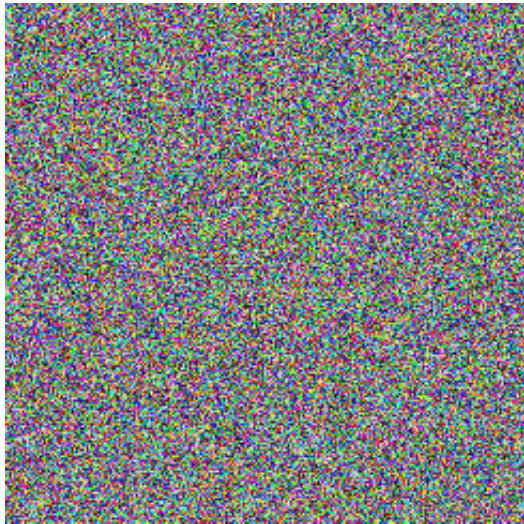


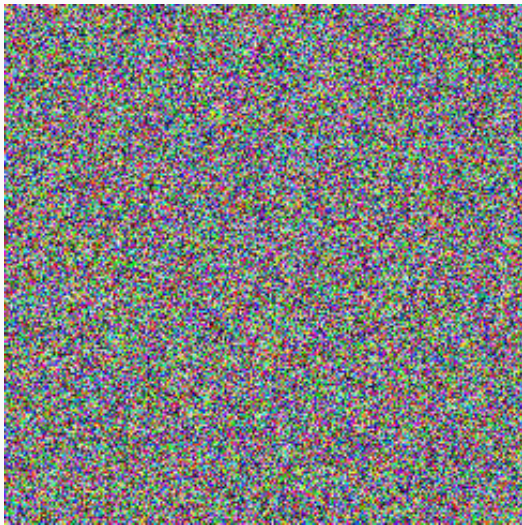










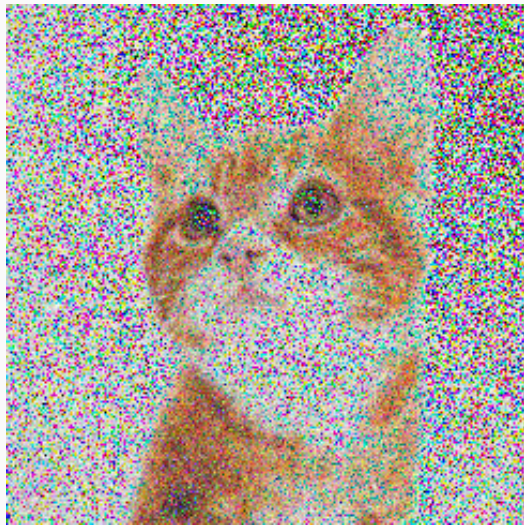
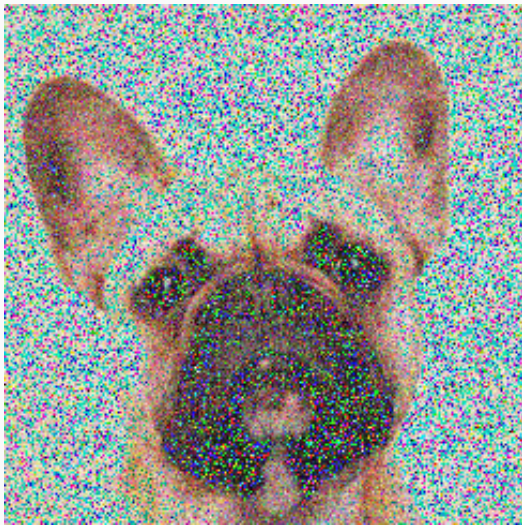


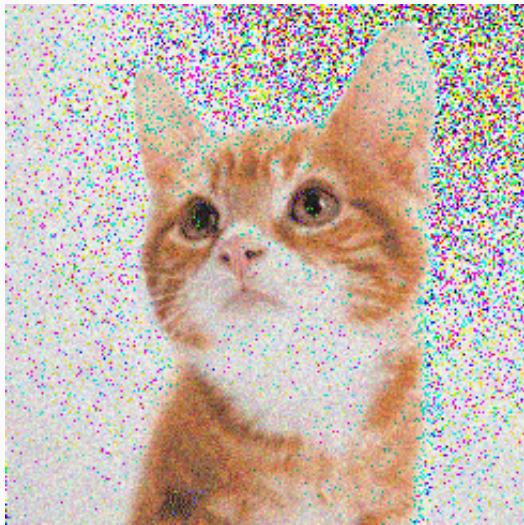
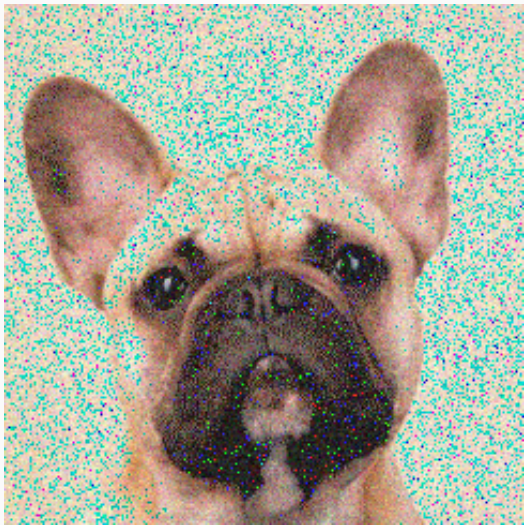














# To image restoration

The sampling of  $p_{\text{data}}$  provides a prior knowledge on data to achieve restoration tasks on images (inpainting, super-resolution, deblurring,...) [Song et al., 2021],[Lugmayr et al., 2022],[Chung et al., 2022], Pseudo-inverse reasoning [Choi et al., 2021]



## Study of the convergence

---

- Experimental study: [S. Chen, Chewi, Lee, et al., 2023; Franzese et al., 2023; Karras et al., 2022]
- Theoretical study: [Benton et al., 2024; S. Chen, Chewi, Li, et al., 2023; De Bortoli et al., 2021; Lee et al., 2022, 2024]
- Under manifold assumption: [M. Chen et al., 2023; De Bortoli, 2022; Wenliang and Moran, 2022]
- 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
- In practice, the convergence of diffusion models is observed using the Frechet Inception Distance (FID) which is 2-Wasserstein distance between Gaussians fitted to datasets.



There are four types of error:

- The initialization error
- The discretization error
- The truncation error
- The score approximation error

$$d\mathbf{y}_t = \beta_{T-t}(\mathbf{y}_t + 2\nabla \log p_{T-t}(\mathbf{y}_t))dt + \sqrt{2\beta_{T-t}}d\mathbf{w}_t, \quad 0 \leq t < T, \quad \mathbf{y}_0 \sim p_T. \quad (10)$$

is replaced by:

$$d\mathbf{y}_t = \beta_{T-t}(\mathbf{y}_t + 2\nabla \log p_{T-t}(\mathbf{y}_t))dt + \sqrt{2\beta_{T-t}}d\mathbf{w}_t, \quad 0 \leq t < T, \quad \mathbf{y}_0 \sim \mathcal{N}_0. \quad (11)$$

- The result is: if  $\mathbf{y}_t$  verifies Equation (14),  $\mathbf{y}_T \sim p_{T-t}$
- Equation (15) produces another stochastic process.
- This holds also for the ODE.



Several choice for the discretization:

SDE schemes	Euler-Maruyama (EM)	$\begin{cases} \tilde{\mathbf{y}}_0^{\Delta, \text{EM}} & \sim \mathcal{N}_0 \\ \tilde{\mathbf{y}}_{k+1}^{\Delta, \text{EM}} & = \tilde{\mathbf{y}}_k^{\Delta, \text{EM}} + \Delta_t \beta_{T-t_k} \left( \tilde{\mathbf{y}}_k^{\Delta, \text{EM}} - 2 \Sigma_{T-t_k}^{-1} \tilde{\mathbf{y}}_k^{\Delta, \text{EM}} \right) + \sqrt{2 \Delta_t \beta_{T-t_k}} \mathbf{z}_k, \mathbf{z}_k \sim \mathcal{N}_0 \end{cases} \quad (12)$	
	Exponential integrator (EI)	$\begin{cases} \tilde{\mathbf{y}}_0^{\Delta, \text{EI}} & \sim \mathcal{N}_0 \\ \tilde{\mathbf{y}}_{k+1}^{\Delta, \text{EI}} & = \tilde{\mathbf{y}}_k^{\Delta, \text{EI}} + \gamma_{1,k} \left( \tilde{\mathbf{y}}_k^{\Delta, \text{EI}} - 2 \Sigma_{T-t_k}^{-1} \tilde{\mathbf{y}}_k^{\Delta, \text{EI}} \right) + \sqrt{2 \gamma_{2,k}} \mathbf{z}_k, \mathbf{z}_k \sim \mathcal{N}_0 \end{cases} \quad (13)$ <p>where <math>\gamma_{1,k} = \exp(B_{T-t_k} - B_{T-t_{k+1}}) - 1</math> and <math>\gamma_{2,k} = \frac{1}{2}(\exp(2B_{T-t_k} - 2B_{T-t_{k+1}}) - 1)</math></p>	
ODE schemes	Explicit Euler	$\begin{cases} \hat{\mathbf{y}}_0^{\Delta, \text{Euler}} & \sim \mathcal{N}_0 \\ \hat{\mathbf{y}}_{k+1}^{\Delta, \text{Euler}} & = \hat{\mathbf{y}}_k^{\Delta, \text{Euler}} + \Delta_t f(t_k, \hat{\mathbf{y}}_k^{\Delta, \text{Euler}}) \quad \text{with } f(t, \mathbf{y}) = \beta_{T-t} \mathbf{y} - \beta_{T-t} \Sigma_{T-t}^{-1} \mathbf{y} \end{cases} \quad (14)$	
	Heun's method	$\begin{cases} \hat{\mathbf{y}}_0^{\Delta, \text{Heun}} & \sim \mathcal{N}_0 \\ \hat{\mathbf{y}}_{k+1/2}^{\Delta, \text{Heun}} & = \hat{\mathbf{y}}_k^{\Delta, \text{Heun}} + \Delta_t f(t_k, \hat{\mathbf{y}}_k^{\Delta, \text{Heun}}) \quad \text{with } f(t, \mathbf{y}) = \beta_{T-t} \mathbf{y} - \beta_{T-t} \Sigma_{T-t}^{-1} \mathbf{y} \\ \hat{\mathbf{y}}_{k+1}^{\Delta, \text{Heun}} & = \hat{\mathbf{y}}_k^{\Delta, \text{Heun}} + \frac{\Delta_t}{2} \left( f(t_k, \hat{\mathbf{y}}_k^{\Delta, \text{Heun}}) + f(t_{k+1}, \hat{\mathbf{y}}_{k+1/2}^{\Delta, \text{Heun}}) \right) \end{cases} \quad (15)$	

$$d\mathbf{y}_t = \beta_{T-t}(\mathbf{y}_t + 2\nabla \log p_{T-t}(\mathbf{y}_t))dt + \sqrt{2\beta_{T-t}}d\mathbf{w}_t, \quad 0 \leq t < T, \quad \mathbf{y}_0 \sim p_T. \quad (16)$$

- At time 0,  $p_0$  does not necessary exists.
- It is preferable to solve Equation (20) from 0 to  $T - \varepsilon$ .
- In general,  $\varepsilon = 10^{-3}$  (Karras et al., 2022; Song et al., 2021)

$$d\mathbf{y}_t = \beta_{T-t}(\mathbf{y}_t + 2\nabla \log p_{T-t}(\mathbf{y}_t))dt + \sqrt{2\beta_{T-t}}d\mathbf{w}_t, \quad 0 \leq t < T, \quad \mathbf{y}_0 \sim p_T. \quad (17)$$

$$d\mathbf{y}_t = \beta_{T-t}(\mathbf{y}_t + 2s_\theta(T-t, \mathbf{y}_t))dt + \sqrt{2\beta_{T-t}}d\mathbf{w}_t, \quad 0 \leq t < T, \quad \mathbf{y}_0 \sim p_T. \quad (18)$$

where  $s_\theta$  is a neural network.

1. The most difficult to estimate theoretically.
2. In general, bounds on the  $L^2$  norm.

## Restriction to the Gaussian case

---

**Gaussian assumption:**  $p_{\text{data}}$  is a centered Gaussian distribution  $\mathcal{N}(\mathbf{0}, \Sigma)$ . ( $\Sigma$  is not necessarily invertible)

- $p_t \sim \mathcal{N}(\mathbf{0}, \Sigma_t)$ , with  $\Sigma_t = e^{-2Bt} \text{Cov}(\mathbf{x}_0) + (1 - e^{-2Bt}) \mathbf{I}$
- $\nabla \log p_t(\mathbf{x}) = -\Sigma_t^{-1} \mathbf{x}$
- Also known if  $p_{\text{data}}$  is a Gaussian mixture [Shah et al., 2023; Zach et al., 2024; Zach et al., 2023].

Note that  $\nabla \log p_t$  is linear.

**Gaussian assumption:**  $p_{\text{data}}$  is a centered Gaussian distribution  $\mathcal{N}(\mathbf{0}, \Sigma)$ . ( $\Sigma$  is not necessarily invertible)

- $p_t \sim \mathcal{N}(\mathbf{0}, \Sigma_t)$ , with  $\Sigma_t = e^{-2Bt} \text{Cov}(\mathbf{x}_0) + (1 - e^{-2Bt}) \mathbf{I}$
- $\nabla \log p_t(\mathbf{x}) = -\Sigma_t^{-1} \mathbf{x}$
- Also known if  $p_{\text{data}}$  is a Gaussian mixture [Shah et al., 2023; Zach et al., 2024; Zach et al., 2023].

Note that  $\nabla \log p_t$  is linear.

### Proposition 2: Characterization of Gaussian distributions through diffusion models

The three following propositions are equivalent:

- (i)  $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \Sigma)$  for some covariance  $\Sigma$ .
- (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  $t > 0$ ,  $\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) = -\Sigma_t^{-1} \mathbf{x}$ .

## Proposition 3: Solution of the backward SDE under Gaussian assumption

Under Gaussian assumption, the strong solution to Equation (3) can be written as:

$$\mathbf{y}_t = e^{-(B_T - B_{T-t})} \Sigma_{T-t} \Sigma_T^{-1} \mathbf{y}_0 + \boldsymbol{\xi}_t, \quad 0 \leq t \leq T \quad (19)$$

where  $\boldsymbol{\xi}_t$  is a Gaussian process. Finally:

$$\text{Cov}(\mathbf{y}_t) = \Sigma_{T-t} + e^{-2(B_T - B_{T-t})} \Sigma_{T-t}^2 \Sigma_T^{-1} (\Sigma_{T-t}^{-1} \text{Cov}(\mathbf{y}_0) \Sigma_T^{-1} \Sigma_{T-t} - \mathbf{I}), \quad (20)$$

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

$$\text{Cov}(\mathbf{y}_t) = \Sigma_{T-t} + e^{-2(B_T - B_{T-t})} \Sigma_{T-t}^2 \Sigma_T^{-1} [\Sigma_T^{-1} \text{Cov}(\mathbf{y}_0) - \mathbf{I}] \quad (21)$$

- $\mathbf{y}_0$  can follow any law.

## Proposition 4: Solution of the ODE probability flow under Gaussian assumption

The solution to the probability flow ODE (8) under Gaussian assumption corresponds to the optimal transport map between  $p_T$  and  $p_{\text{data}}$ . 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 (9). Consequently, the covariance matrix  $\text{Cov}(\mathbf{y}_t)$  verifies

$$\text{Cov}(\mathbf{y}_t) = \Sigma_T^{-1/2} \Sigma_{T-t}^{1/2} \text{Cov}(\mathbf{y}_0) \Sigma_{T-t}^{1/2} \Sigma_T^{-1/2}, \quad 0 \leq t \leq T, \quad (22)$$

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

$$\text{Cov}(\mathbf{y}_t) = \Sigma_T^{-1} \Sigma_{T-t} \text{Cov}(\mathbf{y}_0), \quad 0 \leq t \leq T. \quad (23)$$

- The relation between optimal transport and probability flow ODE (also called Fokker-Planck ODE) has been discussed in Khrlukov et al., 2023; Lavenant and Santambrogio, 2022<sup>4</sup> in the asymptotic case where  $T \mapsto +\infty$ .

<sup>4</sup>Lavenant, H., & Santambrogio, F. (2022). The flow map of the fokker-planck equation does not provide optimal transport. *Applied Mathematics Letters*, 133, 108225.  
<https://doi.org/https://doi.org/10.1016/j.aml.2022.108225>



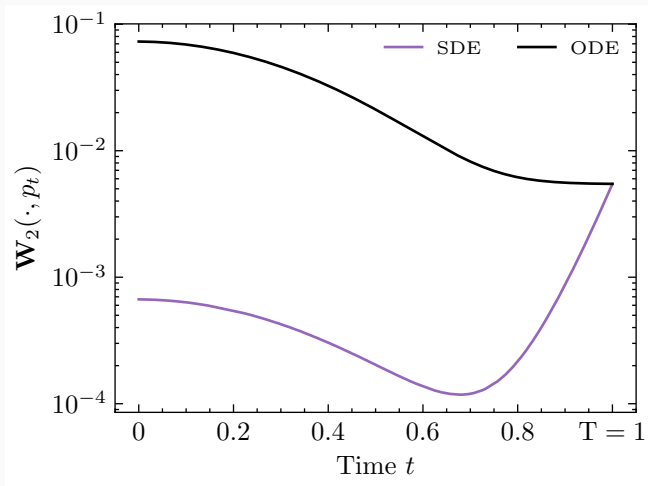
## Proposition 5: Marginals of the generative processes under Gaussian assumption

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  $0 \leq t \leq T$ ,

$$\mathbf{W}_2(\tilde{p}_t, p_t) \leq \mathbf{W}_2(\hat{p}_t, p_t) \quad (24)$$

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

## Initialization error



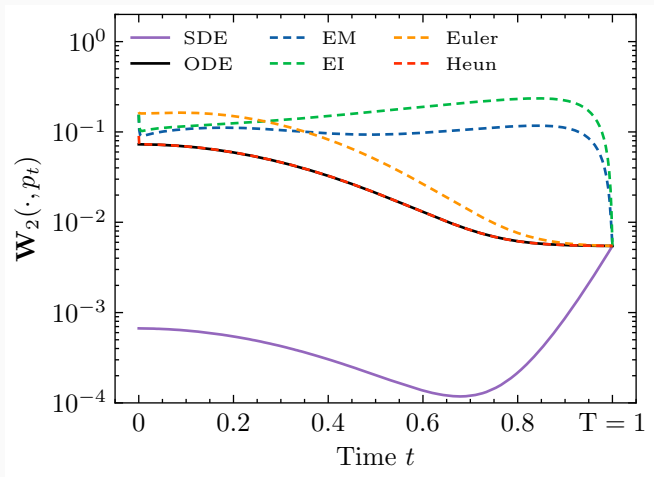
From

$$\tilde{\mathbf{y}}_{k+1}^{\Delta, \text{EM}} = \tilde{\mathbf{y}}_k^{\Delta, \text{EM}} + \Delta_t \beta_{T-t_k} \left( \tilde{\mathbf{y}}_k^{\Delta, \text{EM}} - 2 \Sigma_{T-t_k}^{-1} \tilde{\mathbf{y}}_k^{\Delta, \text{EM}} \right) + \sqrt{2 \Delta_t \beta_{T-t_k}} \mathbf{z}_k, \quad \mathbf{z}_k \sim \mathcal{N}_0$$

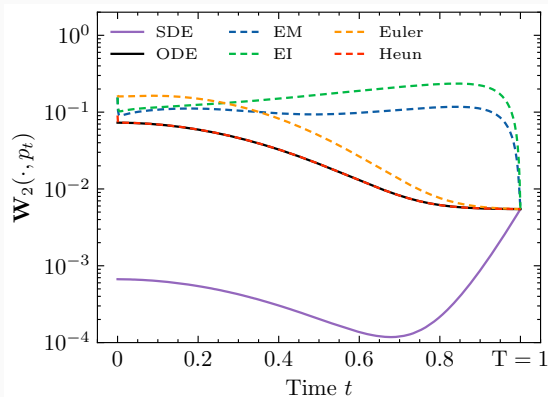
we have:

$$\lambda_i^{\text{EM}, k+1} = \left( 1 + \Delta_t \beta_{T-t_k} \left( 1 - \frac{2}{\lambda_i^{T-t_k}} \right) \right)^2 \lambda_i^{\text{EM}, k} + 2 \Delta_t \beta_{T-t_k}, \quad 1 \leq i \leq d, 0 \leq k \leq N-2 \quad (25)$$

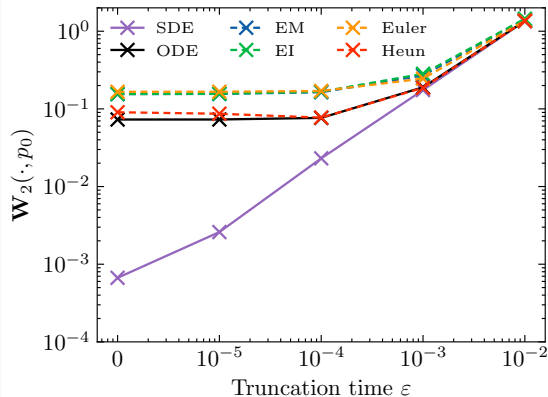
## Discretization error



# Truncation error



(a) Initialization error along the integration time

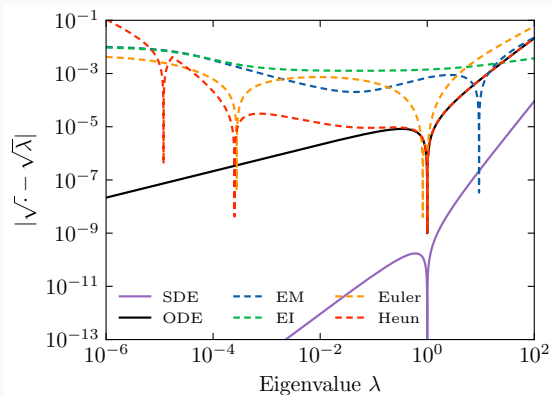


(b) Truncation error for different truncation time  $\varepsilon$

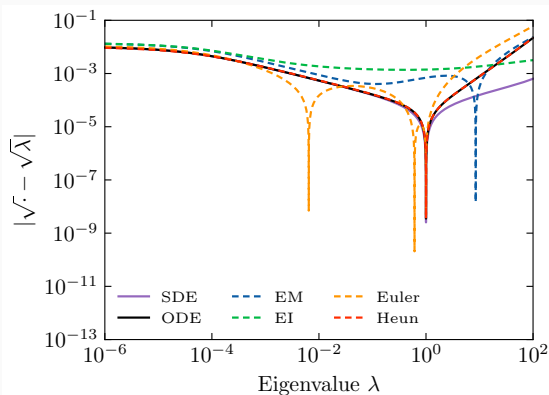
# Ablation study

		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

## Data dependent errors



(a) Initialization error at final time



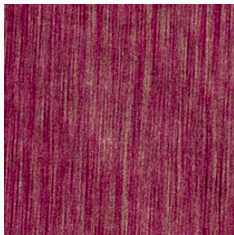
(b) Truncation error at final time for  $\varepsilon = 10^{-3}$

# Score approximation

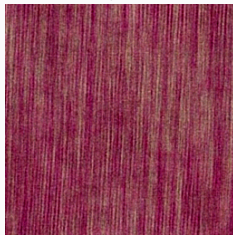
Gaussian sample



$p_{\theta}^{\text{EM}}$



$p_{\theta}^{\text{Heun}}$



$p_{\theta}^{\text{Heun}}$





# Score approximation

$p$	Exact score distribution			Learned score distribution	
	$\mathbf{W}_2(p, p_{\text{data}}) \downarrow$	$\mathbf{W}_2^{\text{emp.}}(p^{\text{emp.}}, p_{\text{data}}) \downarrow$	$\text{FID}(p^{\text{emp.}}, p_{\text{data}}^{\text{emp.}}) \downarrow$	$\mathbf{W}_2^{\text{emp.}}(p_{\theta}^{\text{emp.}}, p_{\text{data}}^{\text{emp.}}) \downarrow$	$\text{FID}(p_{\theta}^{\text{emp.}}, p_{\text{data}}^{\text{emp.}}) \downarrow$
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

- Heun's method fails.
- EM discretization more resilient to score approximation.

- This theoretical analysis led to conclude that Heun's scheme is the best numerical solution, in accordance with empirical previous work [Karras et al., 2022].
- We conducted an empirical analysis with a learned score function using standard architecture which showed the most important one in practice.
- This suggests that assessing the quality of learned score functions is an important research direction for future work.

Thank you for your attention !

Preprint : Diffusion models for Gaussian distributions: Exact solutions and Wasserstein errors, E. Pierret, B. Galerne, 2024, hal, Arxiv

## References

---

- Benton, J., Bortoli, V. D., Doucet, A., & Deligiannidis, G. (2024). Nearly  $\mathcal{O}(1)$ -linear convergence bounds for diffusion models via stochastic localization. *The Twelfth International Conference on Learning Representations*. <https://openreview.net/forum?id=r5njV3BsuD>
- Chen, M., Huang, K., Zhao, T., & Wang, M. (2023). Score approximation, estimation and distribution recovery of diffusion models on low-dimensional data. In A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, & J. Scarlett (Eds.), *Proceedings of the 40th international conference on machine learning* (pp. 4672–4712). PMLR. <https://proceedings.mlr.press/v202/chen23o.html>
- Chen, S., Chewi, S., Lee, H., Li, Y., Lu, J., & Salim, A. (2023). The probability flow ODE is provably fast. *Thirty-seventh Conference on Neural Information Processing Systems*. <https://openreview.net/forum?id=KD6MFeWSAd>
- Chen, S., Chewi, S., Li, J., Li, Y., Salim, A., & Zhang, A. (2023). Sampling is as easy as learning the score: Theory for diffusion models with minimal data assumptions. *The Eleventh International Conference on Learning Representations*. [https://openreview.net/forum?id=zyLVMgsZ0U\\_](https://openreview.net/forum?id=zyLVMgsZ0U_)

- Choi, J., Kim, S., Jeong, Y., Gwon, Y., & Yoon, S. (2021). ILVR: Conditioning method for denoising diffusion probabilistic models. *ILVR*, 14367–14376. Retrieved 2022-11-28, from [https://openaccess.thecvf.com/content/ICCV2021/html/Choi\\_ILVR\\_Conditioning\\_Method\\_for\\_Denoising\\_Diffusion\\_Probabilistic\\_Models\\_ICCV\\_2021\\_paper.html](https://openaccess.thecvf.com/content/ICCV2021/html/Choi_ILVR_Conditioning_Method_for_Denoising_Diffusion_Probabilistic_Models_ICCV_2021_paper.html)
- Chung, H., Sim, B., Ryu, D., & Ye, J. C. (2022). Improving diffusion models for inverse problems using manifold constraints. *Advances in Neural Information Processing Systems (NeurIPS)*.
- De Bortoli, V. (2022). Convergence of denoising diffusion models under the manifold hypothesis. *Transactions on Machine Learning Research*. <https://openreview.net/forum?id=MhK5aXo3gB>
- De Bortoli, V., Thornton, J., Heng, J., & Doucet, A. (2021). Diffusion schrödinger bridge with applications to score-based generative modeling. *Advances in Neural Information Processing Systems*, 34, 17695–17709. Retrieved 2022-11-08, from <https://papers.nips.cc/paper/2021/hash/940392f5f32a7ade1cc201767cf83e31-Abstract.html>
- Dhariwal, P., & Nichol, A. (2021). Diffusion models beat GANs on image synthesis. *Advances in Neural Information Processing Systems*.
- Franzese, G., Rossi, S., Yang, L., Finamore, A., Rossi, D., Filippone, M., & Michiardi, P. (2023). How much is enough? a study on diffusion times in score-based generative models. *Entropy*, 25(4). <https://doi.org/10.3390/e25040633>

- Karras, T., Aittala, M., Aila, T., & Laine, S. (2022). Elucidating the design space of diffusion-based generative models. *Proc. NeurIPS*.
- Khrulkov, V., Ryzhakov, G., Chertkov, A., & Oseledets, I. (2023). Understanding DDPM latent codes through optimal transport. *The Eleventh International Conference on Learning Representations*.  
<https://openreview.net/forum?id=6PIrhAx1j4i>
- Lavenant, H., & Santambrogio, F. (2022). The flow map of the fokker–planck equation does not provide optimal transport. *Applied Mathematics Letters*, 133, 108225.  
<https://doi.org/https://doi.org/10.1016/j.aml.2022.108225>
- Lee, H., Lu, J., & Tan, Y. (2022). Convergence of score-based generative modeling for general data distributions. *NeurIPS 2022 Workshop on Score-Based Methods*.
- Lee, H., Lu, J., & Tan, Y. (2024). Convergence for score-based generative modeling with polynomial complexity. *Proceedings of the 36th International Conference on Neural Information Processing Systems*.
- Lugmayr, A., Danelljan, M., Romero, A., Yu, F., Timofte, R., & Gool, L. V. (2022). Repaint: Inpainting using denoising diffusion probabilistic models. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 11451–11461. <https://api.semanticscholar.org/CorpusID:246240274>
- Pardoux, E. (1986). Grossissement d'une filtration et retournement du temps d'une diffusion. In J. Azéma & M. Yor (Eds.), *Séminaire de probabilités xx 1984/85* (pp. 48–55). Springer Berlin Heidelberg.

- Shah, K., Chen, S., & Klivans, A. (2023). Learning mixtures of gaussians using the DDPM objective. *Thirty-seventh Conference on Neural Information Processing Systems*.  
<https://openreview.net/forum?id=aig7sgdRfl>
- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., & Poole, B. (2021). Score-based generative modeling through stochastic differential equations. *International Conference on Learning Representations*. <https://openreview.net/forum?id=PXTIG12RRHS>
- Wenliang, L. K., & Moran, B. (2022). Score-based generative model learn manifold-like structures with constrained mixing. *NeurIPS 2022 Workshop on Score-Based Methods*.  
<https://openreview.net/forum?id=eSZqalrDLZR>
- Xiao, Z., Kreis, K., & Vahdat, A. (2022). Tackling the generative learning trilemma with denoising diffusion GANs. *International Conference on Learning Representations*.
- Zach, M., Kobler, E., Chambolle, A., & Pock, T. (2024). Product of gaussian mixture diffusion models. *Journal of Mathematical Imaging and Vision*. <https://doi.org/10.1007/s10851-024-01180-3>
- Zach, M., Pock, T., Kobler, E., & Chambolle, A. (2023). Explicit diffusion of gaussian mixture model based image priors. In L. Calatroni, M. Donatelli, S. Morigi, M. Prato, & M. Santacesaria (Eds.), *Scale space and variational methods in computer vision* (pp. 3–15). Springer International Publishing.