

# Introduction to diffusion models for image generation and inpainting

---

Émile Pierret, supervised by Bruno Galerne

July, 5<sup>th</sup>

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

Introduction on generative models

Presentation of the diffusion models

- 2.1. The forward process
- 2.2. The backward process
- 2.3. The continuous framework
- 2.4. Conclusion

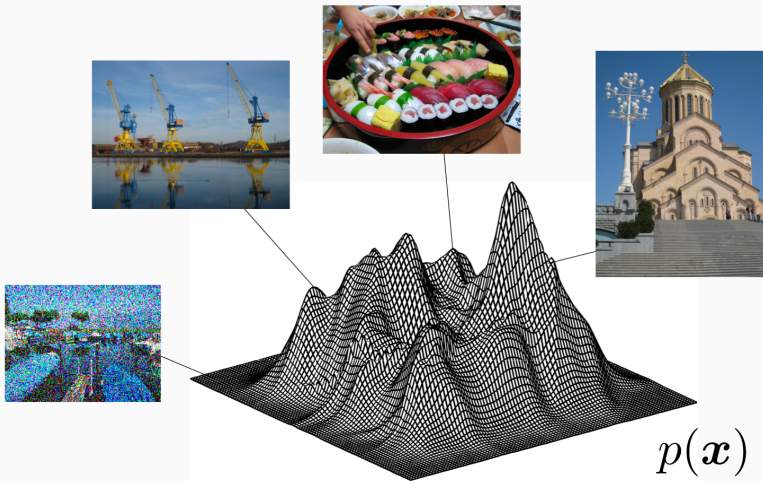
Presentation of RePaint

- 3.1. Conditional diffusion models
- 3.2. Presentation of the article

## Introduction on generative models

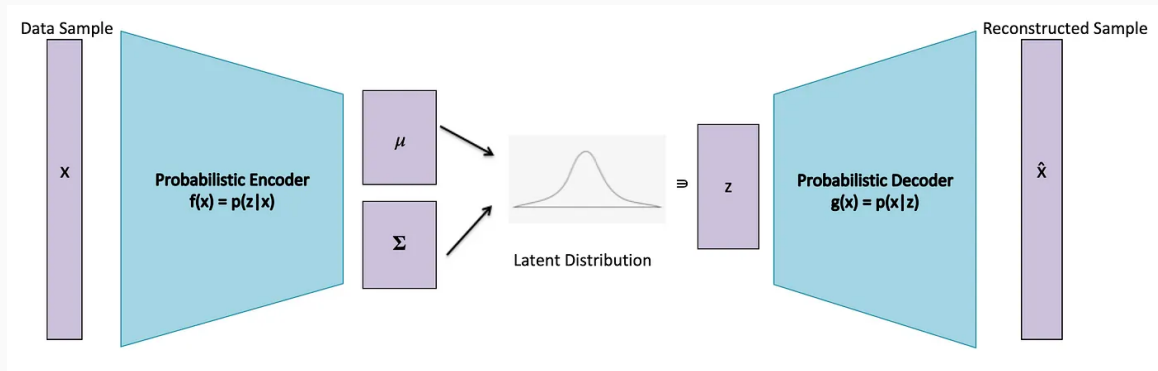
---

# What is a generative model ?



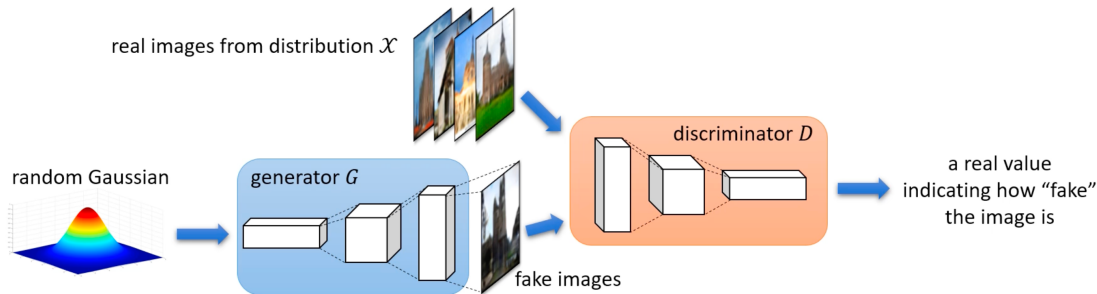


### 22/09 : Variational Auto-Encoder (VAE):



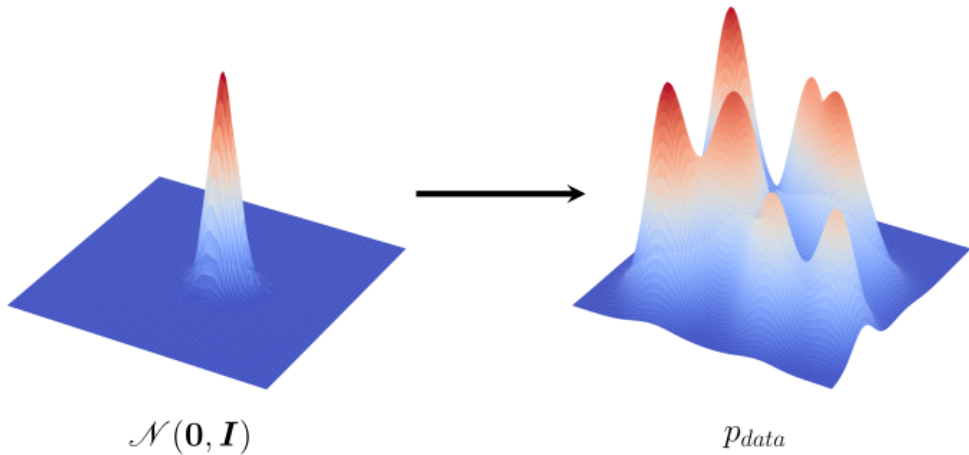
*Image extracted from <https://medium.com/@elzettevanrensburg/generating-the-intuition-behind-variational-auto-encoders-vaes-c7d2f8631a87>*

### 25/10: Generative Adversarial Netowrk (GAN):



*Image extracted from <https://www.microsoft.com/en-us/research/blog/how-can-generative-adversarial-networks-learn-real-life-distributions-easily/>*

# Main idea



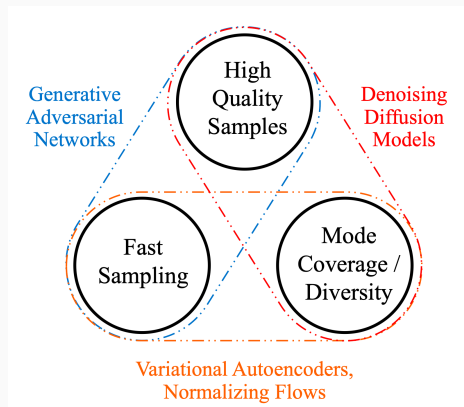


Image extrated from [Xiao et al., 2022]<sup>1</sup>

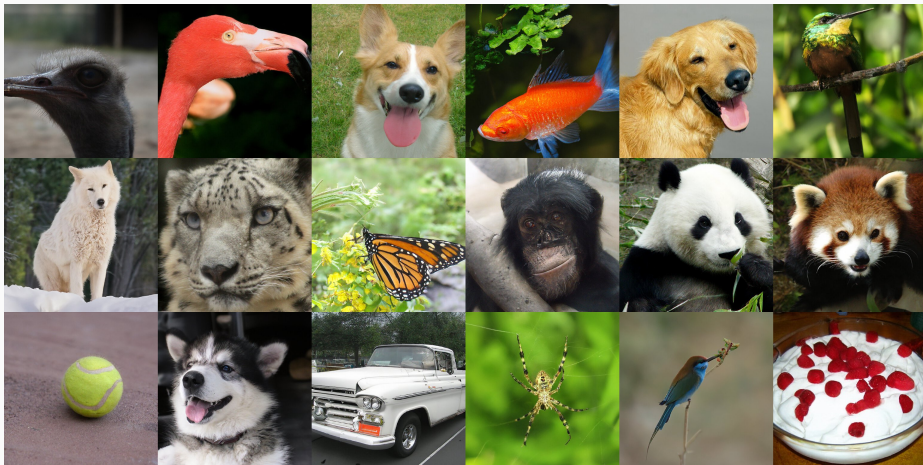
- Dhariwal, P., & Nichol, A. (2021). Diffusion models beat GANs on image synthesis. *Advances in Neural Information Processing Systems*

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

## Presentation of the diffusion models

---

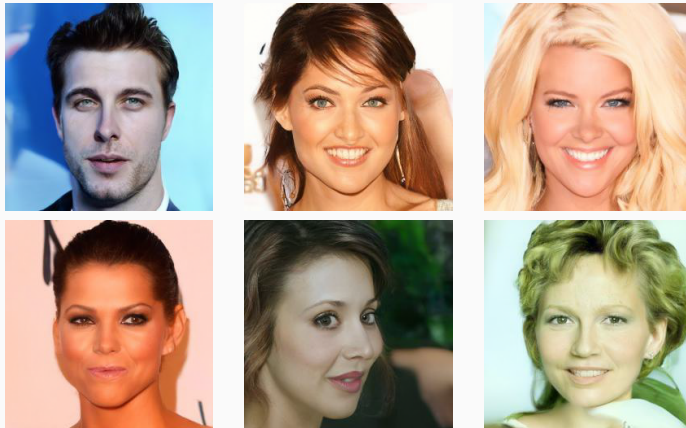
# Samples from diffusion models



Images extracted from [Dhariwal and Nichol, 2021]<sup>2</sup> from a model trained on ImageNet.

<sup>2</sup>Dhariwal, P., & Nichol, A. (2021). Diffusion models beat GANs on image synthesis. *Advances in Neural Information Processing Systems*

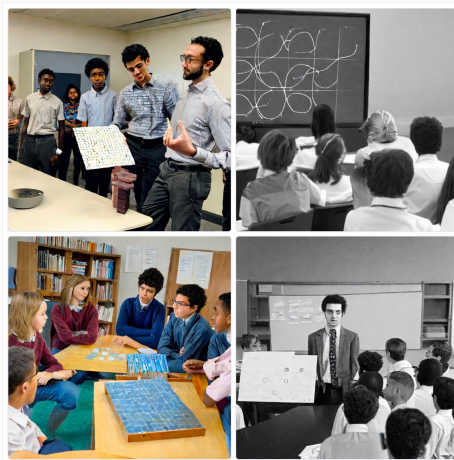
# Home-made with the code of Lugmayr et al., 2022<sup>3</sup>



**NB:** 100s for 10 samples.

<sup>3</sup>Lugmayr, A., Danelljan, M., Romero, A., Yu, F., Timofte, R., & Van Gool, L. (2022). RePaint: Inpainting using denoising diffusion probabilistic models. *RePaint*

# Stable diffusion<sup>4</sup>



*A young mathematician presenting diffusion models to his colleagues*

<sup>4</sup><https://huggingface.co/spaces/stabilityai/stable-diffusion>



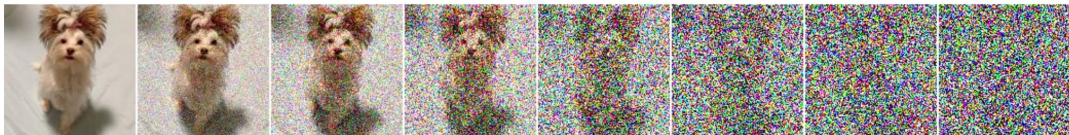
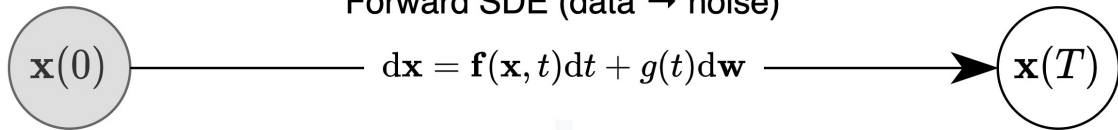
# Stable diffusion<sup>5</sup>



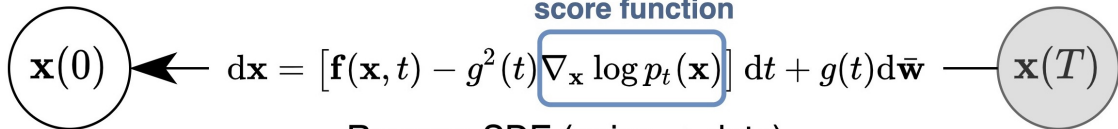
*A young mathematician sadly discovers that he should study stochastic differential equations for his thesis.*

<sup>5</sup><https://huggingface.co/spaces/stabilityai/stable-diffusion>

Forward SDE (data  $\rightarrow$  noise)



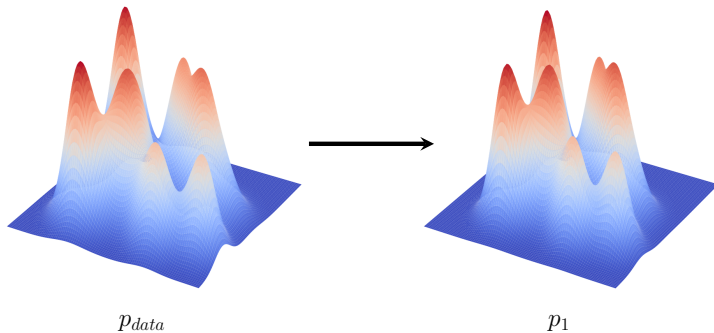
score function



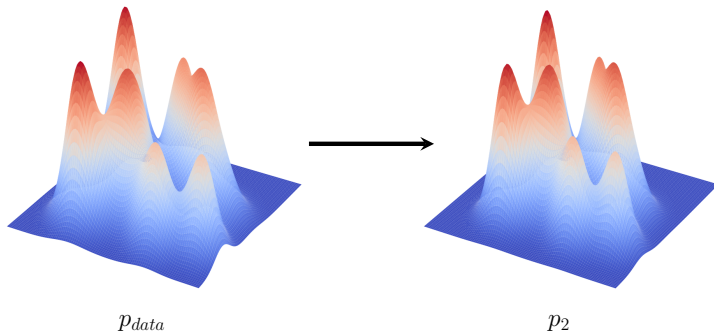
Reverse SDE (noise  $\rightarrow$  data)

Image extracted from [Y. Song et al., 2023]

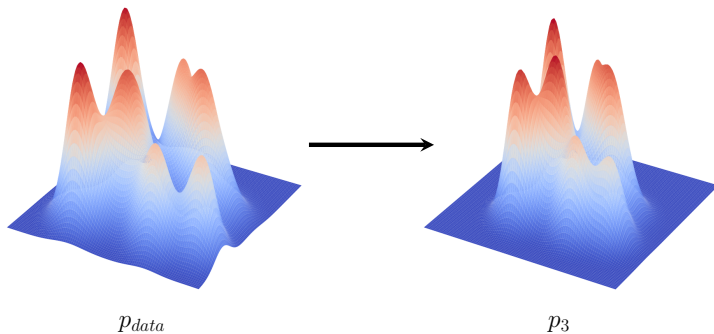
# The forward process



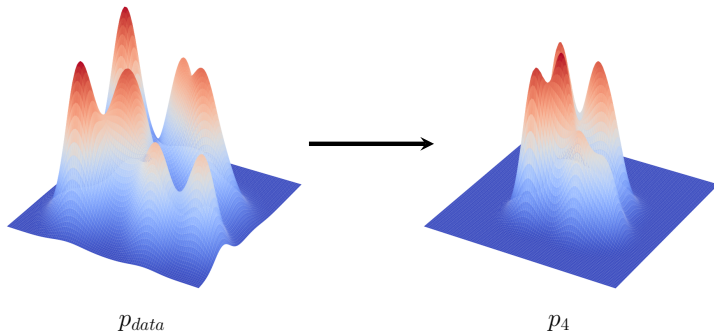
# The forward process



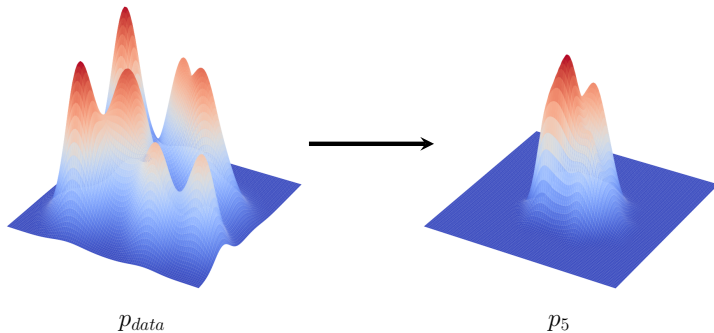
# The forward process



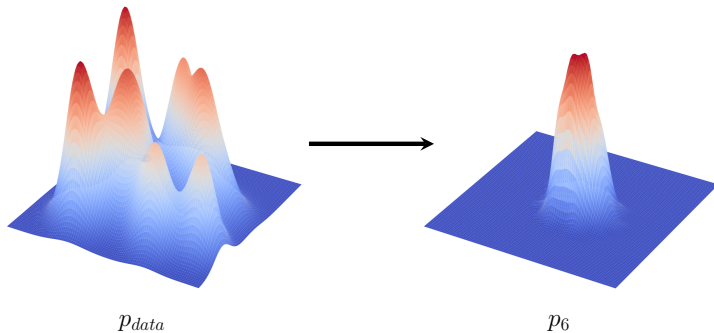
# The forward process



# The forward process

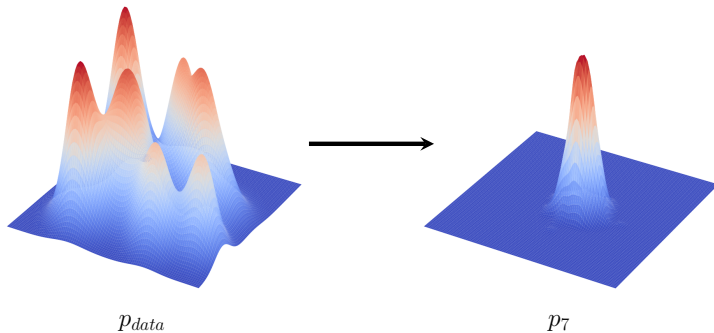


# The forward process





# The forward process

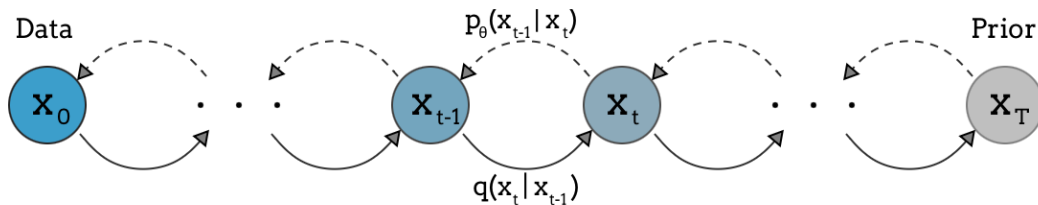


## 2.1. The forward process

# The forward process

Let  $\mathbf{x} \in \mathbb{R}^{3 \times M \times N}$  be an image, sample from  $p_{\text{data}}$ . We denote  $p_0 := p_{\text{data}}$ .

Our objective is to transform  $p_0$  into  $p_T \approx \mathcal{N}(\mathbf{0}, \mathbf{I})$  with  $T$  steps. I will present the process presented in DDPM [Ho et al., 2020]<sup>6</sup>.



# The forward process

Let  $\mathbf{x} \in \mathbb{R}^{3 \times M \times N}$  be an image, sample from  $p_{\text{data}}$ . We denote  $p_0 := p_{\text{data}}$ .

Our objective is to transform  $p_0$  into  $p_T \approx \mathcal{N}(\mathbf{0}, \mathbf{I})$  with  $T$  steps. I will present the process presented in DDPM [Ho et al., 2020]<sup>6</sup>.

Let  $(\beta_t)_{0 \leq t \leq T}$  be an increasing sequence of real numbers in  $[0, 1]$ . Let  $\alpha_t = 1 - \beta_t$ , we would like to construct  $(\mathbf{x}_t)_{0 \leq t \leq T}$ .

---

<sup>6</sup>Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*

# The forward process

Let  $\mathbf{x} \in \mathbb{R}^{3 \times M \times N}$  be an image, sample from  $p_{\text{data}}$ . We denote  $p_0 := p_{\text{data}}$ .

Our objective is to transform  $p_0$  into  $p_T \approx \mathcal{N}(\mathbf{0}, \mathbf{I})$  with  $T$  steps. I will present the process presented in DDPM [Ho et al., 2020]<sup>6</sup>.

Let  $(\beta_t)_{0 \leq t \leq T}$  be an increasing sequence of real numbers in  $[0, 1]$ . Let  $\alpha_t = 1 - \beta_t$ , we would like to construct  $(\mathbf{x}_t)_{0 \leq t \leq T}$ . For all  $t \geq 1$ ,

$$\mathbf{x}_t = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \sqrt{\beta_t} \boldsymbol{\varepsilon}_t, \text{ with } \boldsymbol{\varepsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

---

<sup>6</sup>Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*

## The forward process

Let  $\mathbf{x} \in \mathbb{R}^{3 \times M \times N}$  be an image, sample from  $p_{\text{data}}$ . We denote  $p_0 := p_{\text{data}}$ .

Our objective is to transform  $p_0$  into  $p_T \approx \mathcal{N}(\mathbf{0}, \mathbf{I})$  with  $T$  steps. I will present the process presented in DDPM [Ho et al., 2020]<sup>6</sup>.

Let  $(\beta_t)_{0 \leq t \leq T}$  be an increasing sequence of real numbers in  $[0, 1]$ . Let  $\alpha_t = 1 - \beta_t$ , we would like to construct  $(\mathbf{x}_t)_{0 \leq t \leq T}$ . For all  $t \geq 1$ ,

$$\mathbf{x}_t = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \sqrt{\beta_t} \boldsymbol{\varepsilon}_t, \text{ with } \boldsymbol{\varepsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Denoting  $\alpha_t = 1 - \beta_t$

$$\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\varepsilon}_t, \text{ with } \boldsymbol{\varepsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

---

<sup>6</sup>Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*

## The forward process

Let  $\mathbf{x} \in \mathbb{R}^{3 \times M \times N}$  be an image, sample from  $p_{\text{data}}$ . We denote  $p_0 := p_{\text{data}}$ .

Our objective is to transform  $p_0$  into  $p_T \approx \mathcal{N}(\mathbf{0}, \mathbf{I})$  with  $T$  steps. I will present the process presented in DDPM [Ho et al., 2020]<sup>6</sup>.

Let  $(\beta_t)_{0 \leq t \leq T}$  be an increasing sequence of real numbers in  $[0, 1]$ . Let  $\alpha_t = 1 - \beta_t$ , we would like to construct  $(\mathbf{x}_t)_{0 \leq t \leq T}$ . For all  $t \geq 1$ ,

$$\mathbf{x}_t = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \sqrt{\beta_t} \boldsymbol{\varepsilon}_t, \text{ with } \boldsymbol{\varepsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Denoting  $\alpha_t = 1 - \beta_t$

$$\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\varepsilon}_t, \text{ with } \boldsymbol{\varepsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Other formulation:

$$\mathbf{x}_t \sim \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \alpha_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

---

<sup>6</sup>Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*

## The forward process

Denoting  $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$ , For  $t \geq 1$ ,

$$\mathbf{x}_t \sim \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$$

In other words,  $\mathbf{x}_t$  can be expressed as

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \mathbf{z}_t, \text{ with } \mathbf{z}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

**Example:**

If  $\mathbf{x}_1 = \sqrt{\alpha_1} \mathbf{x}_0 + \sqrt{1 - \alpha_1} \boldsymbol{\varepsilon}_1$ ,

$$\begin{aligned} \mathbf{x}_2 &= \sqrt{\alpha_2} \mathbf{x}_1 + \sqrt{1 - \alpha_2} \boldsymbol{\varepsilon}_2 \\ &= \sqrt{\alpha_2} \sqrt{\alpha_1} \mathbf{x}_0 + \sqrt{\alpha_2} \sqrt{1 - \alpha_1} \boldsymbol{\varepsilon}_1 + \sqrt{1 - \alpha_2} \boldsymbol{\varepsilon}_2 \\ &= \sqrt{\bar{\alpha}_2} \mathbf{x}_0 + (\sqrt{1 - \alpha_1} \boldsymbol{\varepsilon}_1 + \sqrt{1 - \alpha_2} \boldsymbol{\varepsilon}_2) \end{aligned}$$

And:

$$(\sqrt{1 - \alpha_1} \boldsymbol{\varepsilon}_1 + \sqrt{1 - \alpha_2} \boldsymbol{\varepsilon}_2) \sim \mathcal{N}(\mathbf{0}, \alpha_2(1 - \alpha_1) + 1 - \alpha_2) = \mathcal{N}(\mathbf{0}, 1 - \bar{\alpha}_2)$$

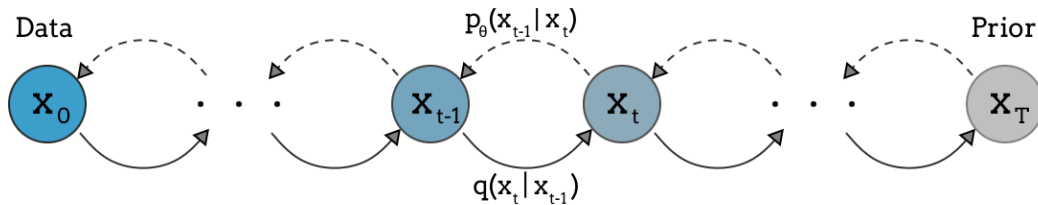


# The forward decomposition

The forward decomposition is:

$$q(\mathbf{x}_{0:T}) = q(\mathbf{x}_0) \prod_{t=1}^T q(\mathbf{x}_t \mid \mathbf{x}_{t-1})$$

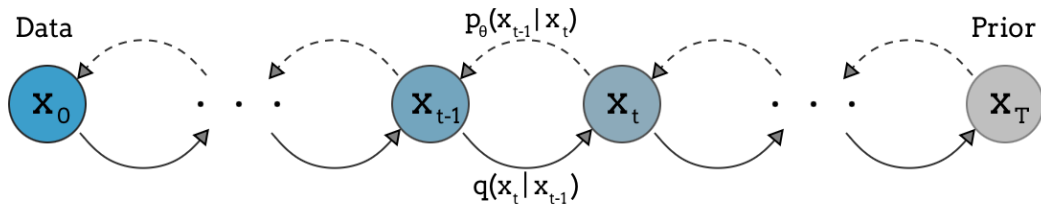
And we denote  $p_t$  the law of  $\mathbf{x}_t$ .



To keep in mind: Conditioning on  $\mathbf{x}_0$ , all is possible.

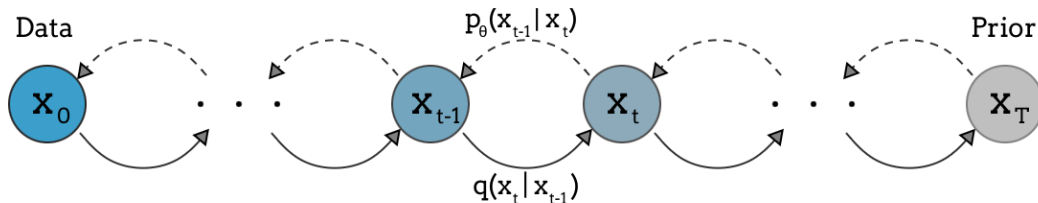
## 2.2. The backward process

## The backward process



<sup>7</sup>Bishop, C. M. (2006). *Pattern recognition and machine learning* (Vol. 4). Springer

# The backward process

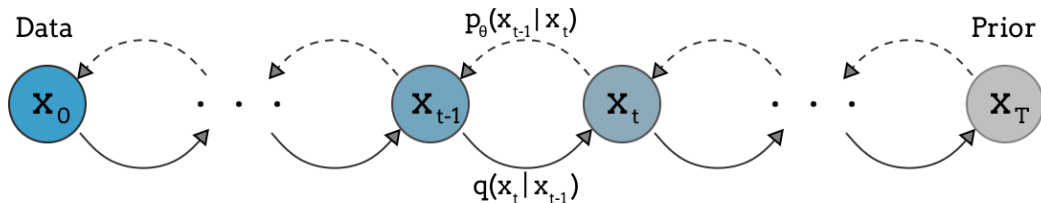


$q(\mathbf{x}_{t-1} | \mathbf{x}_0, \mathbf{x}_t)$  is tractable.

For example, see [Bishop, 2006]<sup>7</sup>

<sup>7</sup>Bishop, C. M. (2006). *Pattern recognition and machine learning* (Vol. 4). Springer

## The backward process



$q(\mathbf{x}_{t-1} | \mathbf{x}_0, \mathbf{x}_t)$  is tractable.

For example, see [Bishop, 2006]<sup>7</sup>

$$\mathbf{x}_{t-1} \sim \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \mathbf{I})$$

with

$$\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) = \frac{\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_t} \beta_t \mathbf{x}_0 + \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}_t$$
$$\tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$$

<sup>7</sup>Bishop, C. M. (2006). *Pattern recognition and machine learning* (Vol. 4). Springer

The backward probability  $p_\theta$

$$p_\theta(x_{0:T}) = p_\theta(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t)$$

We will suppose that  $p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t)$  is Gaussian. [Feller, 1949]<sup>8</sup>.

This can be also explained in the discrete case by [De Bortoli et al., 2021]<sup>9</sup>.

$$p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t) = \mathcal{N}(\mu_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t))$$

---

<sup>8</sup>Feller, W. (1949). On the theory of stochastic processes, with particular reference to applications. *Proceedings of the [First] Berkeley Symposium on Mathematical Statistics and Probability*

<sup>9</sup>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*

$$\begin{aligned}\log p_{\theta}(\mathbf{x}_0) &= \log \int p_{\theta}(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T} \\ &= \log \int \frac{p_{\theta}(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} q(\mathbf{x}_{1:T} | \mathbf{x}_0) d\mathbf{x}_{1:T} \\ &\geq \int \log \frac{p_{\theta}(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} q(\mathbf{x}_{1:T} | \mathbf{x}_0) d\mathbf{x}_{1:T} \text{ by Jensen's inequality}\end{aligned}$$

Consequently,

$$\begin{aligned}\mathbb{E}_{\mathbf{x}_0 \sim p_0} [\log p_{\theta}(\mathbf{x}_0)] &\geq \int \int \log \frac{p_{\theta}(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} q(\mathbf{x}_{1:T} | \mathbf{x}_0) q(\mathbf{x}_0) d\mathbf{x}_{1:T} d\mathbf{x}_0 \\ &= \int \log \frac{p_{\theta}(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} q(\mathbf{x}_{0:T}) d\mathbf{x}_{0:T} \\ &= \mathbb{E}_{\mathbf{x}_{0:T} \sim q} \left[ \log \frac{p_{\theta}(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \right]\end{aligned}$$

## Derivation of the ELBO [Sohl-Dickstein et al., 2015]

By the forward and the backward decompositions,

$$\begin{aligned}\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T} \mid \mathbf{x}_0)} &= \log \left[ p_\theta(\mathbf{x}_T) \prod_{t=1}^T \frac{p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t)}{q(\mathbf{x}_t \mid \mathbf{x}_{t-1})} \right] \\ &= \log p_\theta(\mathbf{x}_T) + \sum_{t=1}^T \log \frac{p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t)}{q(\mathbf{x}_t \mid \mathbf{x}_{t-1})}\end{aligned}$$

Consequently:

$$\begin{aligned}\mathbb{E}_{\mathbf{x}_0 \sim p_0} [\log p_\theta(\mathbf{x}_0)] \\ &\geq \mathbb{E}_{\mathbf{x}_{0:T} \sim q} [\log p_\theta(\mathbf{x}_T)] + \sum_{t=1}^T \mathbb{E}_{\mathbf{x}_{0:T} \sim q} \left[ \log \frac{p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t)}{q(\mathbf{x}_t \mid \mathbf{x}_{t-1})} \right] \\ &= \mathbb{E}_{\mathbf{x}_{0:T} \sim q} [\log p_\theta(\mathbf{x}_T)] + \sum_{t=2}^T \mathbb{E}_{\mathbf{x}_{0:T} \sim q} \left[ \log \frac{p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t)}{q(\mathbf{x}_t \mid \mathbf{x}_{t-1}, \mathbf{x}_0)} \frac{q(\mathbf{x}_{t-1} \mid \mathbf{x}_0)}{q(\mathbf{x}_t \mid \mathbf{x}_0)} \right] + \mathbb{E}_{\mathbf{x}_{0:T} \sim q} \left[ \log \frac{p_\theta(\mathbf{x}_0 \mid \mathbf{x}_1)}{q(\mathbf{x}_1 \mid \mathbf{x}_0)} \right] \\ &= \mathbb{E}_{\mathbf{x}_{0:T} \sim q} \left[ \log \frac{p_\theta(\mathbf{x}_T)}{q(\mathbf{x}_T \mid \mathbf{x}_0)} \right] + \sum_{t=2}^T \mathbb{E}_{\mathbf{x}_{0:T} \sim q} \left[ \log \frac{p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t)}{q(\mathbf{x}_t \mid \mathbf{x}_{t-1}, \mathbf{x}_0)} \right] + \mathbb{E}_{\mathbf{x}_{0:T} \sim q} [\log p_\theta(\mathbf{x}_0 \mid \mathbf{x}_1)]\end{aligned}$$



For  $2 \leq t \leq T$ ,

$$\begin{aligned}
 \mathbb{E}_{\mathbf{x}_{0:T}} \left[ \log \frac{p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)}{q(\mathbf{x}_t \mid \mathbf{x}_{t-1}, \mathbf{x}_0)} \right] &= - \int \log \left( \frac{q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)} \right) q(\mathbf{x}_{0:T}) d\mathbf{x}_{0:T} \\
 &= - \int \log \left( \frac{q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)} \right) q(\mathbf{x}_{0,t-1,t}) d\mathbf{x}_{0,t-1,t} \\
 &= - \int \log \left( \frac{q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)} \right) q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) q(\mathbf{x}_t, \mathbf{x}_0) d\mathbf{x}_{0,t-1,t} \\
 &= - \int \left( \int \log \left( \frac{q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0)}{p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)} \right) q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) d\mathbf{x}_{t-1} \right) q(\mathbf{x}_t, \mathbf{x}_0) d\mathbf{x}_0 \\
 &= - \mathbb{E}_{\mathbf{x}_t, \mathbf{x}_0 \sim q} [\text{D}_{KL}(q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t))]
 \end{aligned}$$

Finally,

$$\begin{aligned} & \mathbb{E}_{\mathbf{x}_0 \sim p_0} [\log p_\theta(\mathbf{x}_0)] \\ & \geq -\mathbb{E}_{\mathbf{x}_0} [\mathcal{D}_{KL}(q(\mathbf{x}_T | \mathbf{x}_0) \| p_\theta(\mathbf{x}_T))] - \sum_{t=2}^T \mathbb{E}_{\mathbf{x}_t, \mathbf{x}_0 \sim q} [\mathcal{D}_{KL}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t))] \\ & + \mathbb{E}_{\mathbf{x}_0, \mathbf{x}_1} [\log p_\theta(\mathbf{x}_0 | \mathbf{x}_1)] \end{aligned}$$

For  $2 \leq t \leq T$ , we are looking for  $D_{KL}(q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t))$ ,

By [Zhang et al., 2021]<sup>10</sup>,

$$D_{KL}(\mathcal{N}(\mu_1, \Sigma_1), \mathcal{N}(\mu_2, \Sigma_2)) = \frac{1}{2} \left[ \log \frac{|\Sigma_2|}{|\Sigma_1|} + \text{Tr}(\Sigma_2^{-1} \Sigma_1) + (\mu_2 - \mu_1)^T \Sigma_2^{-1} (\mu_2 - \mu_1) - n \right]$$

As a reminder,

$$\mathbf{x}_{t-1} \sim \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \mathbf{I})$$

If we suppose that  $p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t) = \mathcal{N}(\mu_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t))$  and  $\Sigma_\theta(\mathbf{x}_t, t) = \tilde{\beta}_t \mathbf{I}$ ,

$$D_{KL}(q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t)) = \frac{1}{2\tilde{\beta}_t^2} \|\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) - \mu_\theta(\mathbf{x}_t, t)\|^2$$

---

<sup>10</sup>Zhang, Y., Liu, W., Chen, Z., Li, K., & Wang, J. (2021). On the properties of kullback-leibler divergence between gaussians. *CoRR*, abs/2102.05485. <https://arxiv.org/abs/2102.05485>

For  $t = T$ ,  $\mathbb{E}_{\mathbf{x}_0} [\mathcal{D}_{KL}(q(\mathbf{x}_T | \mathbf{x}_0) \| p_\theta(\mathbf{x}_T))]$  is neglected because

$$q(\mathbf{x}_T | \mathbf{x}_0) = \mathcal{N}(\mathbb{E}(\mathbf{x}_0), \sqrt{1 - \bar{\alpha}_T} \text{Var}(\mathbf{x}_0)) \approx \mathcal{N}(\mathbf{0}, \mathbf{I})$$

For  $t = 0$ ,  $\mu_\theta(\mathbf{x}_1, 1)$  is penalized to obtain data images in  $[0, 1]$

The loss is finally,

$$L(\theta) = \sum_{t=2}^T \frac{1}{2\tilde{\beta}_t^2} \|\tilde{\mu}_t(\mathbf{x}_t, x_0) - \mu_\theta(\mathbf{x}_t, t)\|^2 - \mathbb{E}_{\mathbf{x}_0, \mathbf{x}_1} [\log p_\theta(\mathbf{x}_0 | \mathbf{x}_1)]$$

**Problem:** This loss is difficult to train, even for 2D distribution (see tutorials at this url<sup>11</sup>)

---

<sup>11</sup>[https://github.com/acids-ircam/diffusion\\_models](https://github.com/acids-ircam/diffusion_models)

## Solution:

Denoting  $\mathbf{x}_t = \sqrt{\alpha_t}\mathbf{x}_0 + \sqrt{1 - \alpha_t}\boldsymbol{\varepsilon}$  with  $\boldsymbol{\varepsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,

$$L_{t-1} = \mathbb{E}_{\mathbf{x}_0, \boldsymbol{\varepsilon}} \left[ \frac{1}{2\sigma_t^2} \left\| \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\varepsilon}) - \frac{\beta_t}{\sqrt{1 - \alpha_t}} \boldsymbol{\varepsilon} \right) - \mu_\theta(\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\varepsilon}), t) \right\|^2 \right].$$

## Solution:

Denoting  $\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\varepsilon}$  with  $\boldsymbol{\varepsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,

$$L_{t-1} = \mathbb{E}_{\mathbf{x}_0, \boldsymbol{\varepsilon}} \left[ \frac{1}{2\sigma_t^2} \left\| \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\varepsilon}) - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\varepsilon} \right) - \mu_\theta(\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\varepsilon}), t) \right\|^2 \right].$$

$$\mu_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\varepsilon}) - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t) \right)$$

$$\Longleftrightarrow$$

$$\boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t) = \frac{\sqrt{1 - \bar{\alpha}_t}}{\beta_t} (\mathbf{x}_t - \sqrt{\alpha_t} \mu_\theta(\mathbf{x}_t, t))$$

## Solution:

Denoting  $\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\varepsilon}$  with  $\boldsymbol{\varepsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,

$$L_{t-1} = \mathbb{E}_{\mathbf{x}_0, \boldsymbol{\varepsilon}} \left[ \frac{1}{2\sigma_t^2} \left\| \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\varepsilon}) - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\varepsilon} \right) - \mu_\theta(\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\varepsilon}), t) \right\|^2 \right].$$

$$\mu_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\varepsilon}) - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t) \right)$$

$$\Longleftrightarrow$$

$$\boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t) = \frac{\sqrt{1 - \bar{\alpha}_t}}{\beta_t} (\mathbf{x}_t - \sqrt{\alpha_t} \mu_\theta(\mathbf{x}_t, t))$$

$$L_{t-1} = \mathbb{E}_{\mathbf{x}_0, \boldsymbol{\varepsilon}} \left[ \frac{\beta_t^2}{2\sigma_t^2 \alpha_t (1 - \bar{\alpha}_t)} \left\| \boldsymbol{\varepsilon} - \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\varepsilon})) \right\|^2 \right].$$

## Solution:

Denoting  $\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\varepsilon}$  with  $\boldsymbol{\varepsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,

$$L_{t-1} = \mathbb{E}_{\mathbf{x}_0, \boldsymbol{\varepsilon}} \left[ \frac{1}{2\sigma_t^2} \left\| \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\varepsilon}) - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\varepsilon} \right) - \mu_\theta(\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\varepsilon}), t) \right\|^2 \right].$$

$$\mu_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\varepsilon}) - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t) \right)$$

$$\Longleftrightarrow$$

$$\boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t) = \frac{\sqrt{1 - \bar{\alpha}_t}}{\beta_t} (\mathbf{x}_t - \sqrt{\alpha_t} \mu_\theta(\mathbf{x}_t, t))$$

$$L_{t-1} = \mathbb{E}_{\mathbf{x}_0, \boldsymbol{\varepsilon}} \left[ \frac{\beta_t^2}{2\sigma_t^2 \alpha_t (1 - \bar{\alpha}_t)} \left\| \boldsymbol{\varepsilon} - \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\varepsilon}), t) \right\|^2 \right].$$

$$L_{\text{simple}}(\theta) = \mathbb{E}_{t, \mathbf{x}_0, \boldsymbol{\varepsilon}} \left[ \left\| \boldsymbol{\varepsilon} - \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\varepsilon}), t) \right\|^2 \right].$$

To convince yourself: Residual learning for denoiser.



The backward process becomes:

$$\begin{aligned} \mathbf{x}_T &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ \mathbf{x}_{t-1} &= \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t) \right) + \tilde{\beta}_t \mathbf{z}_t \text{ with } \mathbf{z}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \end{aligned}$$

The backward process becomes:

$$\begin{aligned}\mathbf{x}_T &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ \mathbf{x}_{t-1} &= \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t) \right) + \tilde{\beta}_t \mathbf{z}_t \text{ with } \mathbf{z}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})\end{aligned}$$

Ho et al., 2020 has the same results with:

$$\begin{aligned}\mathbf{x}_T &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ \mathbf{x}_{t-1} &= \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t) \right) + \beta_t \mathbf{z}_t \text{ with } \mathbf{z}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})\end{aligned}$$

As a reminder,  $\tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$

$$\begin{aligned}
L(\theta) &= \mathbb{E}_{t \sim [1, T], \mathbf{x}_0, \epsilon_t} \left[ \|\epsilon_t - \epsilon_\theta(\mathbf{x}_t, t)\|^2 \right] \\
&= \mathbb{E}_{t \sim [1, T], \mathbf{x}_0, \epsilon_t} \left[ \|\epsilon_t - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_t, t)\|^2 \right]
\end{aligned}$$

---

### Algorithm 1 Training

---

```

1: repeat
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$ 
4:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   Take gradient descent step on
        $\nabla_\theta \|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|^2$ 
6: until converged

```

---



---

### Algorithm 2 Sampling

---

```

1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 

```

---

## Some implementation tricks: Exponential Moving Average (EMA)

During the training, for each training step,

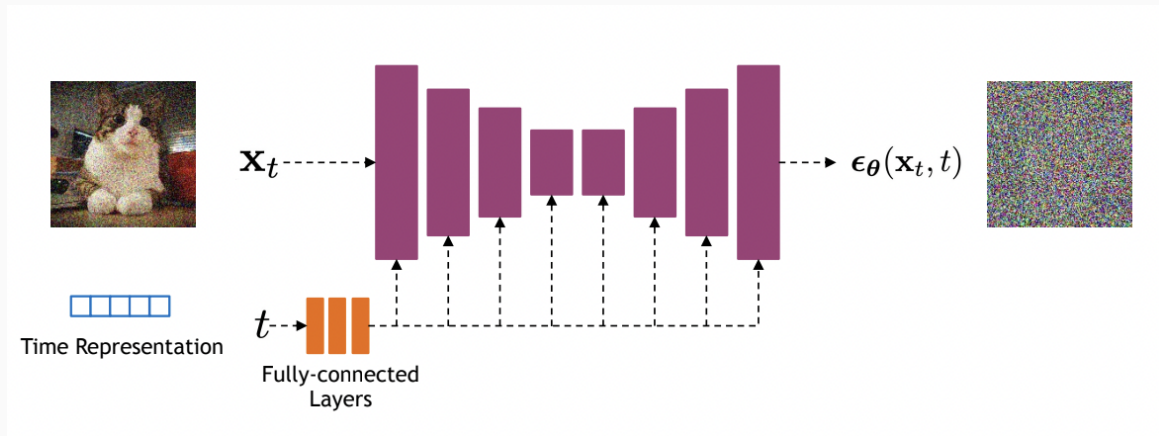
$$\tilde{\theta}_{n+1} = (1 - \mu)\theta_{n+1} + \mu\theta_n$$

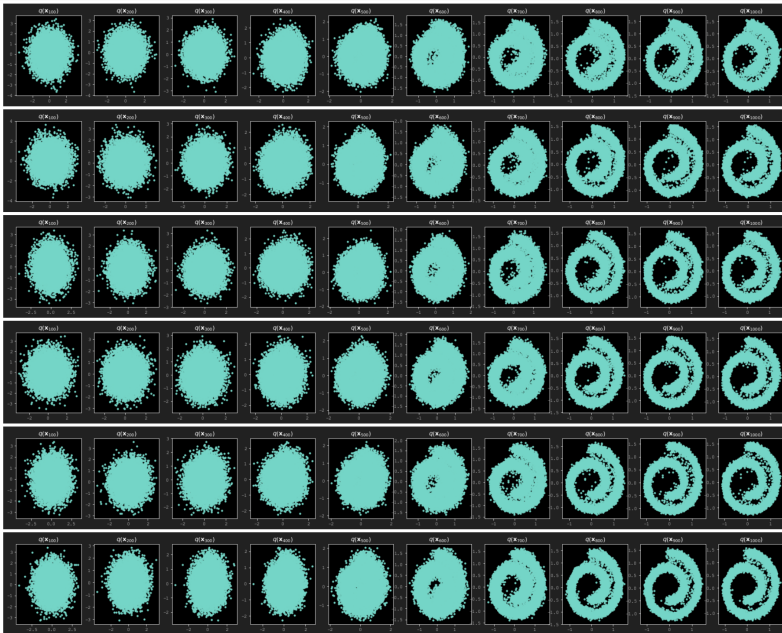
with  $\mu = 0.9$  for example.

# Time representation

Network architecture: UNET

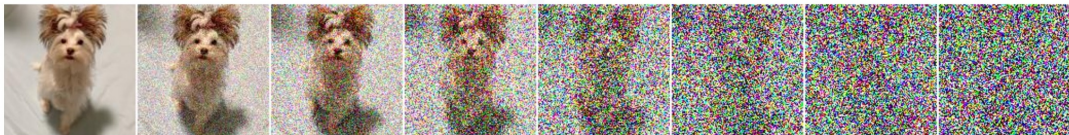
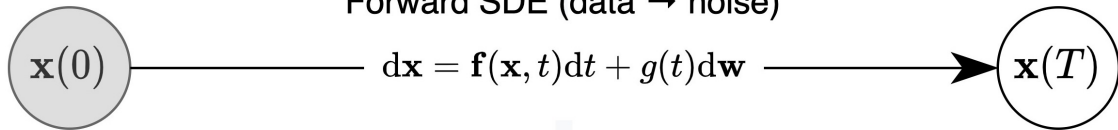
Time representation: sinusoidal positional embeddings or random Fourier features and the network is a U-NET.



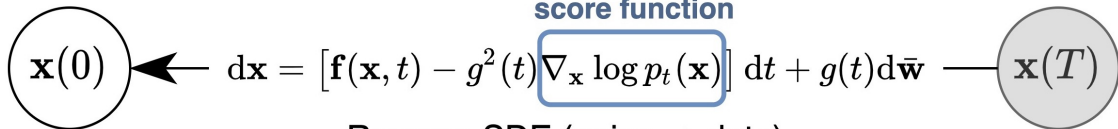


## 2.3. The continuous framework

Forward SDE (data  $\rightarrow$  noise)



score function



Reverse SDE (noise  $\rightarrow$  data)

Image extracted from [Y. Song et al., 2023]



## GROSSISSEMENT D'UNE FILTRATION

### ET RETOURNEMENT DU TEMPS D'UNE DIFFUSION

E. PARDOUX

#### 1. Introduction .

Soit  $\{X_t, 0 \leq t \leq 1\}$  un processus de diffusion dans  $\mathbb{R}^d$   
solution de l'E.D.S:

$$dX_t = b(t, X_t)dt + \sigma(t, X_t)dW_t$$

où  $\{W_t, 0 \leq t \leq 1\}$  est un mouvement brownien standard dans  $\mathbb{R}^{\ell}$  .

Continuant le travail de [2], nous nous posons la question suivante:  
existe-t-il un brownien standard dans  $\mathbb{R}^{\ell}$   $\{\bar{W}_t, 0 \leq t \leq 1\}$  et des coefficients  $\{\bar{b}(t, x), \bar{\sigma}(t, x); 0 \leq t \leq 1, x \in \mathbb{R}^d\}$  tels que le processus  $\bar{X}_t = X_{1-t}, 0 \leq t \leq 1$ , soit solution de :

$$d\bar{X}_t = \bar{b}(t, \bar{X}_t)dt + \bar{\sigma}(t, \bar{X}_t) d\bar{W}_t$$

Notre méthode consiste à identifier  $\{\bar{W}_t\}$ , en résolvant un problème de grossissement de filtration . On pourrait probablement déduire le résultat ci-dessous de ceux de Jeulin [4] et de Jacod

<sup>12</sup>Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., & Poole, B. (2023). Score-based generative modeling through stochastic differential equations. *International Conference on Learning Representations*

What means  $d\mathbf{x}_t = b(t, \mathbf{x}_t)dt + \sigma(t)d\mathbf{w}_t$  ?

What means  $d\mathbf{x}_t = b(t, \mathbf{x}_t)dt + \sigma(t)d\mathbf{w}_t$  ?

$$\mathbf{x}_t = \mathbf{x}_0 + \int_0^t b(s, \mathbf{x}_s)ds + \int_0^t \sigma(s)d\mathbf{w}_s$$

What means  $d\mathbf{x}_t = b(t, \mathbf{x}_t)dt + \sigma(t)d\mathbf{w}_t$  ?

$$\mathbf{x}_t = \mathbf{x}_0 + \int_0^t b(s, \mathbf{x}_s)ds + \int_0^t \sigma(s)d\mathbf{w}_s$$

Discretely,

$$\mathbf{x}(t + \Delta t) - \mathbf{x}(t) \approx \Delta t b(t, \mathbf{x}(t)) + \sigma(t)\sqrt{\Delta t}\xi_t, \text{ with } \xi_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

**Forward process:** Let  $(\boldsymbol{x}_t)_{0 \leq t \leq 1}$  be a diffusion process verifying the SDE:

$$d\boldsymbol{x}_t = b(t, \boldsymbol{x}_t)dt + \sigma(t)d\boldsymbol{w}_t$$

where  $(\boldsymbol{w}_t)_{0 \leq t \leq 1}$  is a standard Brownian motion.

**Forward process:** Let  $(\boldsymbol{x}_t)_{0 \leq t \leq 1}$  be a diffusion process verifying the SDE:

$$d\boldsymbol{x}_t = b(t, \boldsymbol{x}_t)dt + \sigma(t)d\boldsymbol{w}_t$$

where  $(\boldsymbol{w}_t)_{0 \leq t \leq 1}$  is a standard Brownian motion. **Example:**

$$d\boldsymbol{x}_t = -\frac{1}{2}\beta(t)\boldsymbol{x}_tdt + \sqrt{\beta(t)}d\boldsymbol{w}_t$$

**Forward process:** Let  $(\mathbf{x}_t)_{0 \leq t \leq 1}$  be a diffusion process verifying the SDE:

$$d\mathbf{x}_t = b(t, \mathbf{x}_t)dt + \sigma(t)d\mathbf{w}_t$$

where  $(\mathbf{w}_t)_{0 \leq t \leq 1}$  is a standard Brownian motion. **Example:**

$$d\mathbf{x}_t = -\frac{1}{2}\beta(t)\mathbf{x}_tdt + \sqrt{\beta(t)}d\mathbf{w}_t$$

$$\mathbf{x}(t + \Delta t) = \mathbf{x}_t - \frac{1}{2}\beta(t)\Delta t\mathbf{x}_t + \sqrt{\beta(t)\Delta t}\boldsymbol{\varepsilon}_t$$

**Forward process:** Let  $(\mathbf{x}_t)_{0 \leq t \leq 1}$  be a diffusion process verifying the SDE:

$$d\mathbf{x}_t = b(t, \mathbf{x}_t)dt + \sigma(t)d\mathbf{w}_t$$

where  $(\mathbf{w}_t)_{0 \leq t \leq 1}$  is a standard Brownian motion. **Example:**

$$d\mathbf{x}_t = -\frac{1}{2}\beta(t)\mathbf{x}_tdt + \sqrt{\beta(t)}d\mathbf{w}_t$$

$$\begin{aligned}\mathbf{x}(t + \Delta t) &= \mathbf{x}_t - \frac{1}{2}\beta(t)\Delta t\mathbf{x}_t + \sqrt{\beta(t)\Delta t}\boldsymbol{\varepsilon}_t \\ &= (1 - \frac{1}{2}\beta(t)\Delta t)\mathbf{x}_t + \sqrt{\beta(t)\Delta t}\boldsymbol{\varepsilon}_t\end{aligned}$$



**Forward process:** Let  $(\mathbf{x}_t)_{0 \leq t \leq 1}$  be a diffusion process verifying the SDE:

$$d\mathbf{x}_t = b(t, \mathbf{x}_t)dt + \sigma(t)d\mathbf{w}_t$$

where  $(\mathbf{w}_t)_{0 \leq t \leq 1}$  is a standard Brownian motion. **Example:**

$$d\mathbf{x}_t = -\frac{1}{2}\beta(t)\mathbf{x}_tdt + \sqrt{\beta(t)}d\mathbf{w}_t$$

$$\begin{aligned}\mathbf{x}(t + \Delta t) &= \mathbf{x}_t - \frac{1}{2}\beta(t)\Delta t\mathbf{x}_t + \sqrt{\beta(t)\Delta t}\boldsymbol{\varepsilon}_t \\ &= (1 - \frac{1}{2}\beta(t)\Delta t)\mathbf{x}_t + \sqrt{\beta(t)\Delta t}\boldsymbol{\varepsilon}_t \\ &\approx \sqrt{1 - \beta(t)\Delta t}\mathbf{x}_t + \sqrt{\beta(t)\Delta t}\boldsymbol{\varepsilon}_t\end{aligned}$$

**Forward process:** Let  $(\mathbf{x}_t)_{0 \leq t \leq 1}$  be a diffusion process verifying the SDE:

$$d\mathbf{x}_t = b(t, \mathbf{x}_t)dt + \sigma(t)d\mathbf{w}_t$$

where  $(\mathbf{w}_t)_{0 \leq t \leq 1}$  is a standard Brownian motion. **Example:**

$$d\mathbf{x}_t = -\frac{1}{2}\beta(t)\mathbf{x}_tdt + \sqrt{\beta(t)}d\mathbf{w}_t$$

$$\begin{aligned}\mathbf{x}(t + \Delta t) &= \mathbf{x}_t - \frac{1}{2}\beta(t)\Delta t\mathbf{x}_t + \sqrt{\beta(t)\Delta t}\boldsymbol{\varepsilon}_t \\ &= (1 - \frac{1}{2}\beta(t)\Delta t)\mathbf{x}_t + \sqrt{\beta(t)\Delta t}\boldsymbol{\varepsilon}_t \\ &\approx \sqrt{1 - \beta(t)\Delta t}\mathbf{x}_t + \sqrt{\beta(t)\Delta t}\boldsymbol{\varepsilon}_t \\ &\approx \sqrt{1 - \beta_t}\mathbf{x}_t + \sqrt{\beta_t}\mathbf{z}_t\end{aligned}$$

**Forward process:** Let  $(\mathbf{x}_t)_{0 \leq t \leq 1}$  be a diffusion process verifying the SDE:

$$d\mathbf{x}_t = b(t, \mathbf{x}_t)dt + \sigma(t)d\mathbf{w}_t$$

where  $(\mathbf{w}_t)_{0 \leq t \leq 1}$  is a standard Brownian motion. **Example:**

$$d\mathbf{x}_t = -\frac{1}{2}\beta(t)\mathbf{x}_tdt + \sqrt{\beta(t)}d\mathbf{w}_t$$

$$\begin{aligned}\mathbf{x}(t + \Delta t) &= \mathbf{x}_t - \frac{1}{2}\beta(t)\Delta t\mathbf{x}_t + \sqrt{\beta(t)\Delta t}\boldsymbol{\varepsilon}_t \\ &= (1 - \frac{1}{2}\beta(t)\Delta t)\mathbf{x}_t + \sqrt{\beta(t)\Delta t}\boldsymbol{\varepsilon}_t \\ &\approx \sqrt{1 - \beta(t)\Delta t}\mathbf{x}_t + \sqrt{\beta(t)\Delta t}\boldsymbol{\varepsilon}_t \\ &\approx \sqrt{1 - \beta_t}\mathbf{x}_t + \sqrt{\beta_t}\mathbf{z}_t\end{aligned}$$

It is DDPM ! With  $\beta_t = \Delta_t\beta\left(\frac{t}{T}\right)$

Question from [Pardoux, 1986]<sup>13</sup>: Is there  $\bar{b}, \bar{\sigma}$  such that  $\bar{x}_t = x_{1-t}$  is solution of the reverse SDE:

$$d\bar{x}_t = \bar{b}(t, \bar{x}_t)dt + \bar{\sigma}(t)d\bar{w}_t$$

---

<sup>13</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

Question from [Pardoux, 1986]<sup>13</sup>: Is there  $\bar{b}, \bar{\sigma}$  such that  $\bar{\mathbf{x}}_t = \mathbf{x}_{1-t}$  is solution of the reverse SDE:

$$d\bar{\mathbf{x}}_t = \bar{b}(t, \bar{\mathbf{x}}_t)dt + \bar{\sigma}(t)d\bar{\mathbf{w}}_t$$

**Backward process:** Under certain assumptions,  $\mathbf{x}_t$  is a weak solution of the following SDE, backward in time:

$$d\mathbf{x}_t = [b(t, \mathbf{x}_t) - \sigma(t)^2 \nabla_{\mathbf{x}} \log p_t(\bar{\mathbf{x}}_t)] dt + \sigma(t) \otimes d\hat{\mathbf{w}}_t$$

with  $\hat{\mathbf{w}}_t = \mathbf{w}_t - \mathbf{w}_1 - \int_t^1 \nabla \log p_s(\mathbf{x}_s) ds$  which is a backward Brownian motion adapted to the filtration  $\sigma(\mathbf{x}_t \cup \{\mathbf{w}_s - \mathbf{w}_t, t \leq s \leq 1\})$

---

<sup>13</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

Question from [Pardoux, 1986]<sup>13</sup>: Is there  $\bar{b}, \bar{\sigma}$  such that  $\bar{x}_t = x_{1-t}$  is solution of the reverse SDE:

$$d\bar{x}_t = \bar{b}(t, \bar{x}_t)dt + \bar{\sigma}(t)d\bar{w}_t$$

**Backward process:** Under certain assumptions,  $x_t$  is a weak solution of the following SDE, backward in time:

$$dx_t = [b(t, x_t) - \sigma(t)^2 \nabla_x \log p_t(\bar{x}_t)] dt + \sigma(t) \otimes d\hat{w}_t$$

with  $\hat{w}_t = w_t - w_1 - \int_t^1 \nabla \log p_s(x_s) ds$  which is a backward Brownian motion adapted to the filtration  $\sigma(x_t \cup \{w_s - w_t, t \leq s \leq 1\})$

That means:

$$\bar{b}(t, \bar{x}_t) = -b(1-t, \bar{x}_t) - \sigma(1-t)^2 \nabla_x \log p_{1-t}(\bar{x})$$

$$\bar{\sigma}(t) = -\sigma(1-t)$$

$$\bar{w}_t = w_{1-t} - w_1 - \int_{1-t}^1 \nabla \log p_s(x_s) ds$$

---

<sup>13</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

$\mathbf{x}_t \approx \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)$  and:

$$q(\mathbf{x}_t \mid \mathbf{x}_0) \propto \exp \left( -\frac{\|\mathbf{x}_t - \sqrt{\bar{\alpha}_t} \mathbf{x}_0\|^2}{2(1 - \bar{\alpha}_t)} \right)$$

---

<sup>14</sup>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)*

$\mathbf{x}_t \approx \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)$  and:

$$q(\mathbf{x}_t \mid \mathbf{x}_0) \propto \exp \left( -\frac{\|\mathbf{x}_t - \sqrt{\bar{\alpha}_t} \mathbf{x}_0\|^2}{2(1 - \bar{\alpha}_t)} \right)$$

$$\nabla \log q(\mathbf{x}_t \mid \mathbf{x}_0) = -\frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t} \mathbf{x}_0}{1 - \bar{\alpha}_t} \approx -\frac{\boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)}{1 - \bar{\alpha}_t}$$

---

<sup>14</sup>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)*



$\mathbf{x}_t \approx \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)$  and:

$$q(\mathbf{x}_t \mid \mathbf{x}_0) \propto \exp \left( -\frac{\|\mathbf{x}_t - \sqrt{\bar{\alpha}_t} \mathbf{x}_0\|^2}{2(1 - \bar{\alpha}_t)} \right)$$

$$\nabla \log q(\mathbf{x}_t \mid \mathbf{x}_0) = -\frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t} \mathbf{x}_0}{1 - \bar{\alpha}_t} \approx -\frac{\boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)}{1 - \bar{\alpha}_t}$$

Consequently,  $-\frac{\boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)}{1 - \bar{\alpha}_t}$  minimizes  $\mathbb{E}_{\mathbf{x}_0} (\|\nabla \log q(\mathbf{x}_t \mid \mathbf{x}_0) - X\|^2)$

---

<sup>14</sup>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)*

$\mathbf{x}_t \approx \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)$  and:

$$q(\mathbf{x}_t \mid \mathbf{x}_0) \propto \exp \left( -\frac{\|\mathbf{x}_t - \sqrt{\bar{\alpha}_t} \mathbf{x}_0\|^2}{2(1 - \bar{\alpha}_t)} \right)$$

$$\nabla \log q(\mathbf{x}_t \mid \mathbf{x}_0) = -\frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t} \mathbf{x}_0}{1 - \bar{\alpha}_t} \approx -\frac{\boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)}{1 - \bar{\alpha}_t}$$

Consequently,  $-\frac{\boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)}{1 - \bar{\alpha}_t}$  minimizes  $\mathbb{E}_{\mathbf{x}_0} (\|\nabla \log q(\mathbf{x}_t \mid \mathbf{x}_0) - X\|^2)$

and,  $s_\theta(\mathbf{x}_t, t) = -\frac{\boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)}{1 - \bar{\alpha}_t} \approx \nabla \log q(\mathbf{x}_t)$

---

<sup>14</sup>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)*

$\mathbf{x}_t \approx \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)$  and:

$$q(\mathbf{x}_t | \mathbf{x}_0) \propto \exp \left( -\frac{\|\mathbf{x}_t - \sqrt{\bar{\alpha}_t} \mathbf{x}_0\|^2}{2(1 - \bar{\alpha}_t)} \right)$$

$$\nabla \log q(\mathbf{x}_t | \mathbf{x}_0) = -\frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t} \mathbf{x}_0}{1 - \bar{\alpha}_t} \approx -\frac{\boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)}{1 - \bar{\alpha}_t}$$

Consequently,  $-\frac{\boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)}{1 - \bar{\alpha}_t}$  minimizes  $\mathbb{E}_{\mathbf{x}_0} (\|\nabla \log q(\mathbf{x}_t | \mathbf{x}_0) - X\|^2)$

and,  $s_\theta(\mathbf{x}_t, t) = -\frac{\boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)}{1 - \bar{\alpha}_t} \approx \nabla \log q(\mathbf{x}_t)$

$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} (\mathbf{x}_t + (1 - \alpha_t) s_\theta(\mathbf{x}_t, t)) + \tilde{\beta}_t \mathbf{z}_t \text{ with } \mathbf{z}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

This is well explained in [Chung et al., 2022]<sup>14</sup>

---

<sup>14</sup>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)*

With  $b(t, \mathbf{x}) = -\frac{1}{2}\beta(t)\mathbf{x}$ ,  $\sigma(t) = \sqrt{\beta(t)}$ ,

$$\bar{\mathbf{x}}_t = \left[ b(t, \bar{\mathbf{x}}_t) - \sigma(t)^2 \nabla_{\mathbf{x}} \log p_t(\bar{\mathbf{x}}) \right] dt + \sigma(t) d\bar{\mathbf{w}}_t$$

becomes:

$$\begin{aligned} d\bar{\mathbf{x}}_t &= \frac{1}{2}\beta(t)\mathbf{x}dt - \beta(t)\nabla_{\mathbf{x}} \log p_t(\bar{\mathbf{x}})dt + \sqrt{\beta(t)}d\bar{\mathbf{w}}_t \\ \mathbf{x}_{t-1} &= \frac{1}{\sqrt{1-\beta_t}}\mathbf{x}_t + \frac{\beta_t}{\sqrt{1-\beta_t}}s_{\theta}(\mathbf{x}_t, t) + \tilde{\beta}_t\mathbf{z}_t \end{aligned}$$

And  $\frac{1}{\sqrt{1-\beta_t}} \approx 1 + \frac{1}{2}\beta_t$ ,  $\frac{\beta_t}{\sqrt{1-\beta_t}} \approx \beta_t(1 + \frac{1}{2}\beta_t) \approx \beta_t$

Can be applied to other SDE as the model SMLD [Y. Song and Ermon, 2019]<sup>15</sup> where  $d\mathbf{x}_t = \sqrt{\frac{d[\sigma^2(t)]}{dt}} d\mathbf{w}_t$

---

<sup>15</sup>Song, Y., & Ermon, S. (2019). Generative modeling by estimating gradients of the data distribution. *Advances in Neural Information Processing Systems*

## To other approaches

To a SDE  $dx_t = f(x_t, t)dt + g(t)dw_t$  can be associated an ODE:  $dx = [f(x, t) - \frac{1}{2}g(t)^2 \nabla_x \log p_t(x)] dt$

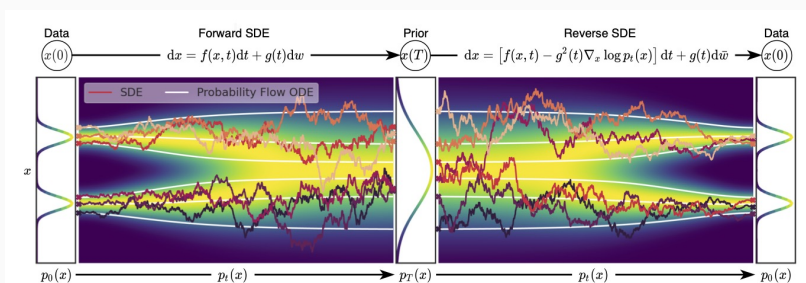


Image extracted from [Y. Song et al., 2023]<sup>16</sup>

<sup>16</sup>Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., & Poole, B. (2023). Score-based generative modeling through stochastic differential equations. *International Conference on Learning Representations*

## 2.4. Conclusion

- Speed-up the convergence. [De Bortoli et al., 2021]<sup>17</sup>, [Dhariwal and Nichol, 2021]<sup>18</sup>
- The data distribution does not have density Bortoli, 2022<sup>19</sup>
- Not presented here: DDIM [J. Song et al., 2021]<sup>20</sup>.
- Interesting tutorials: Song's code: [https://github.com/yang-song/score\\_sde](https://github.com/yang-song/score_sde), 2D models: [https://github.com/acids-ircam/diffusion\\_models/tree/main](https://github.com/acids-ircam/diffusion_models/tree/main)
- Reading recommendation : [Ho et al., 2020]<sup>21</sup>, [Y. Song et al., 2023]<sup>22</sup>

---

<sup>17</sup>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*

<sup>18</sup>Dhariwal, P., & Nichol, A. (2021). Diffusion models beat GANs on image synthesis. *Advances in Neural Information Processing Systems*

<sup>19</sup>Bortoli, V. D. (2022). Convergence of denoising diffusion models under the manifold hypothesis. *Transactions on Machine Learning Research*

<sup>20</sup>Song, J., Meng, C., & Ermon, S. (2021). Denoising diffusion implicit models. *International Conference on Learning Representations*

<sup>21</sup>Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*

<sup>22</sup>Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., & Poole, B. (2023). Score-based generative modeling through stochastic differential equations. *International Conference on Learning Representations*



## Presentation of RePaint

---

Various approaches:

- Re-train a score-based model. [Saharia et al., 2022]<sup>23</sup>
- Classifier guidance method. [Y. Song et al., 2023]<sup>24</sup>
- The replacement method, presented today. [Lugmayr et al., 2022]<sup>25</sup>, [Chung et al., 2022]<sup>26</sup>
- Pseudo-inverse reasoning [Choi et al., 2021]<sup>27</sup>

---

<sup>23</sup>Saharia, C., Ho, J., Chan, W., Salimans, T., Fleet, D. J., & Norouzi, M. (2022). Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*

<sup>24</sup>Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., & Poole, B. (2023). Score-based generative modeling through stochastic differential equations. *International Conference on Learning Representations*

<sup>25</sup>Lugmayr, A., Danelljan, M., Romero, A., Yu, F., Timofte, R., & Van Gool, L. (2022). RePaint: Inpainting using denoising diffusion probabilistic models. *RePaint*

<sup>26</sup>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)*

<sup>27</sup>Choi, J., Kim, S., Jeong, Y., Gwon, Y., & Yoon, S. (2021). ILVR: Conditioning method for denoising diffusion probabilistic models. *ILVR*

## Diverse inpainting

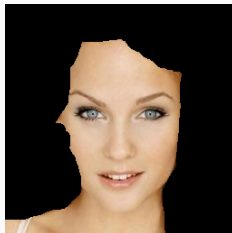
Let suppose that  $x \sim p_{\text{data}}$ , we suppose that there exists a mask  $m$  such that  $mx$  is **known**.

**Objective:** Sample  $x$  conditioned on  $mx = y$ .

Example:



$x$



$y = mx$

## Diverse inpainting

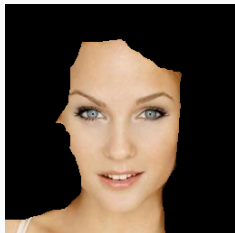
Let suppose that  $x \sim p_{\text{data}}$ , we suppose that there exists a mask  $m$  such that  $mx$  is **known**.

**Objective:** Sample  $x$  conditioned on  $mx = y$ .

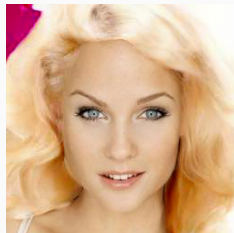
Example:



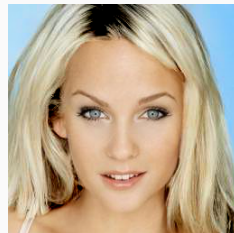
$x$



$y = mx$



Sample 1



Sample 2

Let denote  $\overline{m} = 1 - m$

**Objective:** Sample  $x_0$  conditioned on  $mx_0 = y$ .

Let suppose that we know a diffusion model for  $x_0 \sim p_{\text{data}}$ .

We would like to model the evolution of  $z_t = \overline{m}x_t$  for the forward and the backward processes.

Let denote  $\overline{m} = 1 - m$

**Objective:** Sample  $x_0$  conditioned on  $mx_0 = y$ .

Let suppose that we know a diffusion model for  $x_0 \sim p_{\text{data}}$ .

We would like to model the evolution of  $z_t = \overline{m}x_t$  for the forward and the backward processes.

The forward process is:

$$dz_t = -\frac{1}{2}\beta(t)z_t dt + \sqrt{\beta(t)}dw_t$$

Let denote  $\overline{m} = 1 - m$

**Objective:** Sample  $x_0$  conditioned on  $mx_0 = y$ .

Let suppose that we know a diffusion model for  $x_0 \sim p_{\text{data}}$ .

We would like to model the evolution of  $z_t = \overline{m}x_t$  for the forward and the backward processes.

The forward process is:

$$dz_t = -\frac{1}{2}\beta(t)z_t dt + \sqrt{\beta(t)}dw_t$$

The reversed SDE is conditioned on  $mx_0 = y$  and becomes:

$$dz_t = \left[ -\frac{1}{2}\beta(t)z_t - \beta(t)\nabla_z \log p(z_t \mid mx_0 = y) \right] dt + \sqrt{\beta(t)}d\overline{w}_t$$

Let denote  $\overline{m} = 1 - m$

**Objective:** Sample  $x_0$  conditioned on  $mx_0 = y$ .

Let suppose that we know a diffusion model for  $x_0 \sim p_{\text{data}}$ .

We would like to model the evolution of  $z_t = \overline{m}x_t$  for the forward and the backward processes.

The forward process is:

$$dz_t = -\frac{1}{2}\beta(t)z_t dt + \sqrt{\beta(t)}dw_t$$

The reversed SDE is conditioned on  $mx_0 = y$  and becomes:

$$dz_t = \left[ -\frac{1}{2}\beta(t)z_t - \beta(t)\nabla_z \log p(z_t \mid mx_0 = y) \right] dt + \sqrt{\beta(t)}d\overline{w}_t$$

**Problem:**  $\nabla_z \log p_t(z_t \mid mx_0 = y)$  is not tractable.



Let denote  $\overline{m} = 1 - m$

$$p_t [\mathbf{z}_t \mid m\mathbf{x}_0 = \mathbf{y}]$$

Let denote  $\overline{m} = 1 - m$

$$p_t [z_t \mid mx_0 = y] = p_t [\overline{m}x_t \mid mx_0 = y]$$

Let denote  $\overline{m} = 1 - m$

$$\begin{aligned} p_t [z_t \mid mx_0 = y] &= p_t [\overline{m}x_t \mid mx_0 = y] \\ &= \int p_t [\overline{m}x_t \mid mx_0 = y, mx_t] p_t [mx_t \mid mz_0 = y] d(mx_t) \end{aligned}$$

Let denote  $\overline{m} = 1 - m$

$$\begin{aligned} p_t [z_t \mid mx_0 = y] &= p_t [\overline{m}x_t \mid mx_0 = y] \\ &= \int p_t [\overline{m}x_t \mid mx_0 = y, mx_t] p_t [mx_t \mid mz_0 = y] d(mx_t) \\ &= \mathbb{E}_{p_t[mx_t \mid mx_0 = y]} \left( p_t [\overline{m}x_t \mid mx_0 = y, mx_t] \right) \end{aligned}$$

Let denote  $\overline{m} = 1 - m$

$$\begin{aligned} p_t [z_t \mid mx_0 = y] &= p_t [\overline{m}x_t \mid mx_0 = y] \\ &= \int p_t [\overline{m}x_t \mid mx_0 = y, mx_t] p_t [mx_t \mid mz_0 = y] d(mx_t) \\ &= \mathbb{E}_{p_t[mx_t \mid mx_0 = y]} \left( p_t [\overline{m}x_t \mid mx_0 = y, mx_t] \right) \\ &\approx \mathbb{E}_{p_t[mx_t \mid mx_0 = y]} \left( p_t [\overline{m}x_t \mid mx_t] \right) \end{aligned}$$

Let denote  $\overline{m} = 1 - m$

$$\begin{aligned} p_t [z_t \mid mx_0 = y] &= p_t [\overline{m}x_t \mid mx_0 = y] \\ &= \int p_t [\overline{m}x_t \mid mx_0 = y, mx_t] p_t [mx_t \mid mz_0 = y] d(mx_t) \\ &= \mathbb{E}_{p_t[mx_t \mid mx_0 = y]} \left( p_t [\overline{m}x_t \mid mx_0 = y, mx_t] \right) \\ &\approx \mathbb{E}_{p_t[mx_t \mid mx_0 = y]} \left( p_t [\overline{m}x_t \mid mx_t] \right) \\ &\approx p_t [\overline{m}x_t \mid mx_t] \end{aligned}$$

Let denote  $\overline{m} = 1 - m$

$$\begin{aligned} p_t [z_t \mid mx_0 = y] &= p_t [\overline{m}x_t \mid mx_0 = y] \\ &= \int p_t [\overline{m}x_t \mid mx_0 = y, mx_t] p_t [mx_t \mid mz_0 = y] d(mx_t) \\ &= \mathbb{E}_{p_t[mx_t \mid mx_0 = y]} \left( p_t [\overline{m}x_t \mid mx_0 = y, mx_t] \right) \\ &\approx \mathbb{E}_{p_t[mx_t \mid mx_0 = y]} \left( p_t [\overline{m}x_t \mid mx_t] \right) \\ &\approx p_t [\overline{m}x_t \mid mx_t] \end{aligned}$$

Now,

$$\log p_t (x_t) = \log p_t (\overline{m}x_t, x_t)$$

Let denote  $\overline{m} = 1 - m$

$$\begin{aligned} p_t [z_t \mid mx_0 = y] &= p_t [\overline{m}x_t \mid mx_0 = y] \\ &= \int p_t [\overline{m}x_t \mid mx_0 = y, mx_t] p_t [mx_t \mid mz_0 = y] d(mx_t) \\ &= \mathbb{E}_{p_t[mx_t \mid mx_0 = y]} \left( p_t [\overline{m}x_t \mid mx_0 = y, mx_t] \right) \\ &\approx \mathbb{E}_{p_t[mx_t \mid mx_0 = y]} \left( p_t [\overline{m}x_t \mid mx_t] \right) \\ &\approx p_t [\overline{m}x_t \mid mx_t] \end{aligned}$$

Now,

$$\log p_t (x_t) = \log p_t (\overline{m}x_t, x_t) = \log p_t (\overline{m}x_t \mid mx_t) + \log p_t (mx_t)$$



Let denote  $\overline{m} = 1 - m$

$$\begin{aligned} p_t [z_t \mid mx_0 = y] &= p_t [\overline{m}x_t \mid mx_0 = y] \\ &= \int p_t [\overline{m}x_t \mid mx_0 = y, mx_t] p_t [mx_t \mid mz_0 = y] d(mx_t) \\ &= \mathbb{E}_{p_t[mx_t \mid mx_0 = y]} \left( p_t [\overline{m}x_t \mid mx_0 = y, mx_t] \right) \\ &\approx \mathbb{E}_{p_t[mx_t \mid mx_0 = y]} \left( p_t [\overline{m}x_t \mid mx_t] \right) \\ &\approx p_t [\overline{m}x_t \mid mx_t] \end{aligned}$$

Now,

$$\log p_t (x_t) = \log p_t (\overline{m}x_t, x_t) = \log p_t (\overline{m}x_t \mid mx_t) + \log p_t (mx_t)$$

Thus,

$$\nabla_z \log p_t (\overline{m}x_t \mid mx_t) = \nabla_{\overline{m}x} \log p_t (x_t)$$

Let denote  $\overline{m} = 1 - m$

$$\begin{aligned} p_t [z_t \mid mx_0 = y] &= p_t [\overline{m}x_t \mid mx_0 = y] \\ &= \int p_t [\overline{m}x_t \mid mx_0 = y, mx_t] p_t [mx_t \mid mz_0 = y] d(mx_t) \\ &= \mathbb{E}_{p_t[mx_t \mid mx_0 = y]} \left( p_t [\overline{m}x_t \mid mx_0 = y, mx_t] \right) \\ &\approx \mathbb{E}_{p_t[mx_t \mid mx_0 = y]} \left( p_t [\overline{m}x_t \mid mx_t] \right) \\ &\approx p_t [\overline{m}x_t \mid mx_t] \end{aligned}$$

Now,

$$\log p_t (x_t) = \log p_t (\overline{m}x_t, x_t) = \log p_t (\overline{m}x_t \mid mx_t) + \log p_t (mx_t)$$

Thus,

$$\nabla_z \log p_t (\overline{m}x_t \mid mx_t) = \nabla_{\overline{m}x} \log p_t (x_t) = \overline{m} \nabla_x \log p_t (x_t)$$

Let denote  $\overline{m} = 1 - m$

$$\begin{aligned}
p_t [z_t \mid mx_0 = y] &= p_t [\overline{m}x_t \mid mx_0 = y] \\
&= \int p_t [\overline{m}x_t \mid mx_0 = y, mx_t] p_t [mx_t \mid mz_0 = y] d(mx_t) \\
&= \mathbb{E}_{p_t[mx_t \mid mx_0 = y]} \left( p_t [\overline{m}x_t \mid mx_0 = y, mx_t] \right) \\
&\approx \mathbb{E}_{p_t[mx_t \mid mx_0 = y]} \left( p_t [\overline{m}x_t \mid mx_t] \right) \\
&\approx p_t [\overline{m}x_t \mid mx_t]
\end{aligned}$$

Now,

$$\log p_t (x_t) = \log p_t (\overline{m}x_t, x_t) = \log p_t (\overline{m}x_t \mid mx_t) + \log p_t (mx_t)$$

Thus,

$$\nabla_z \log p_t (\overline{m}x_t \mid mx_t) = \nabla_{\overline{m}x} \log p_t (x_t) = \overline{m} \nabla_x \log p_t (x_t)$$

Consequently, the approached backward SDE becomes:

$$dz_t = \left[ -\frac{1}{2} \overline{m} \beta(t) x_t - \overline{m} \beta(t) \nabla_x \log p_t (x_t) \right] dt + \sqrt{\beta(t)} d\overline{w}_t$$

## How to build a conditional diffusion model ?

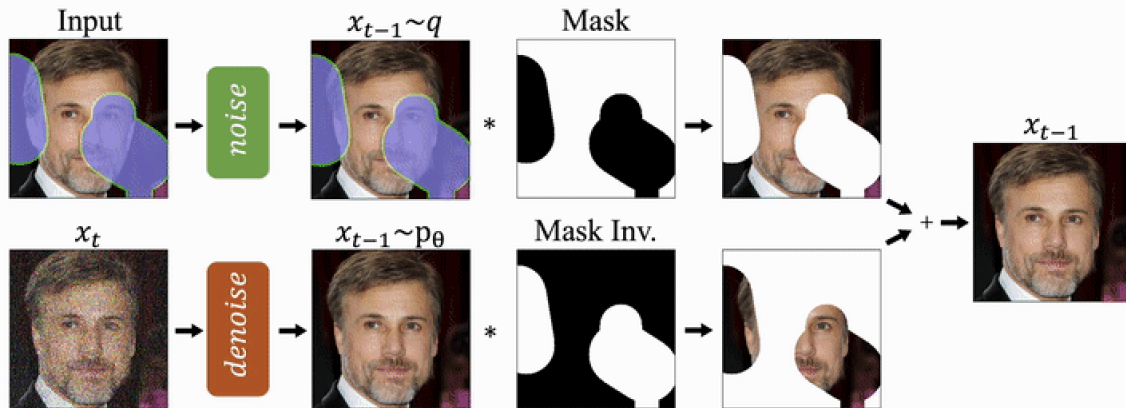


Image extracted from [Lugmayr et al., 2022]<sup>28</sup>

<sup>28</sup>Lugmayr, A., Danelljan, M., Romero, A., Yu, F., Timofte, R., & Van Gool, L. (2022). RePaint: Inpainting using denoising diffusion probabilistic models. *RePaint*

1. **Input:** A masked image  $m\mathbf{x} = \mathbf{y}$
2. Sample  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
3. For  $t = T, \dots, 1$
4.   Sample  $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
5.    $\mathbf{x}_{t-1}^{\text{known}} = \sqrt{\bar{\alpha}_t} \mathbf{y} + (1 - \bar{\alpha}_t) \boldsymbol{\varepsilon}$
6.   Sample  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$  else  $\mathbf{z} = \mathbf{0}$
7.    $\mathbf{x}_{t-1}^{\text{unknown}} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
8.    $\mathbf{x}_{t-1} = m \odot \mathbf{x}_{t-1}^{\text{known}} + (1 - m) \odot \mathbf{x}_{t-1}^{\text{unknown}}$
9. **Output:**  $\mathbf{x}_0$

- $\mathbf{x}_{t-1}^{\text{known}}$  is resampled at each step knowing:

$$\mathbf{x}_{t-1} \sim \mathcal{N}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$$

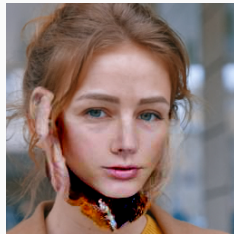
## Let's try it



## Let's try it

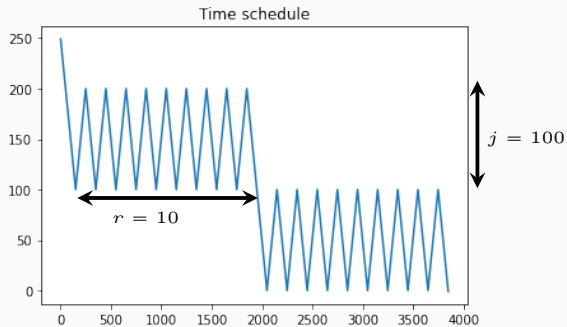


## Let's try it





# Time schedule



1. For each group of times of size  $j$
2.     Do  $r$  times:
3.         For  $s = T, \dots, T - j$
4.             Do the backward from  $s$  to  $s - 1$
5.         For  $s = T - j, \dots, T$
6.             Do the forward from  $s$  to  $s + 1$
7.         For  $s = T, \dots, T - j$
8.             Do the backward from  $s$  to  $s - 1$

**Problem:** From 1000 steps to 18820 steps with  $j = r = 10$

In DDPM, there are two choices for the variance schedule of the backward process:  $\beta_t$  or  $\tilde{\beta}_t$ . Nichol and Dhariwal, 2021 proposes to learn the variance as

$$\Sigma_{\theta}(\mathbf{x}, t) = \exp \left[ v \log(\beta_t) + (1 - v) \log(\tilde{\beta}_t) \right]$$

This reduces the number of steps from 1000 steps to 250 steps.

**Now:** From 250 steps to 4570 steps with  $j = r = 10$ .

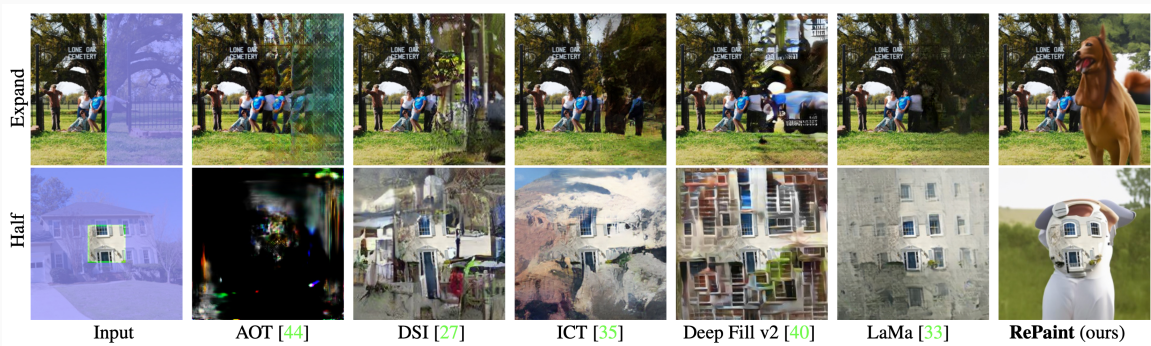
## "Justification" of the algorithm

From the article: "However, the sampling of the known pixels using is performed without considering the generated parts of the image, which introduces disharmony. Although the model tries to harmonize the image again in every step, it can never fully converge because the same issue occurs in the next step. Moreover, in each reverse step, the maximum change to an image declines due to the variance schedule of  $\beta_t$ . Thus, the method cannot correct mistakes that lead to disharmonious boundaries in the subsequent steps due to restricted flexibility. **As a consequence, the model needs more time to harmonize the conditional information  $x^{\text{known}}$  with the generated information  $x^{\text{unknown}}$  in one step before advancing to the next denoising step.**"

## Let's try it



# Fails



- Free learning method is an interesting approach.
- Probably far from the true conditional distribution [Trippe et al., 2023]<sup>29</sup> .
- An other approach with manifold constraints [Chung et al., 2022]<sup>30</sup>.

---

<sup>29</sup>Trippe, B. L., Yim, J., Tischer, D., Baker, D., Broderick, T., Barzilay, R., & Jaakkola, T. S. (2023). Diffusion probabilistic modeling of protein backbones in 3d for the motif-scaffolding problem. *International Conference on Learning Representations*

<sup>30</sup>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)*

Thank you for your attention !

## References

---

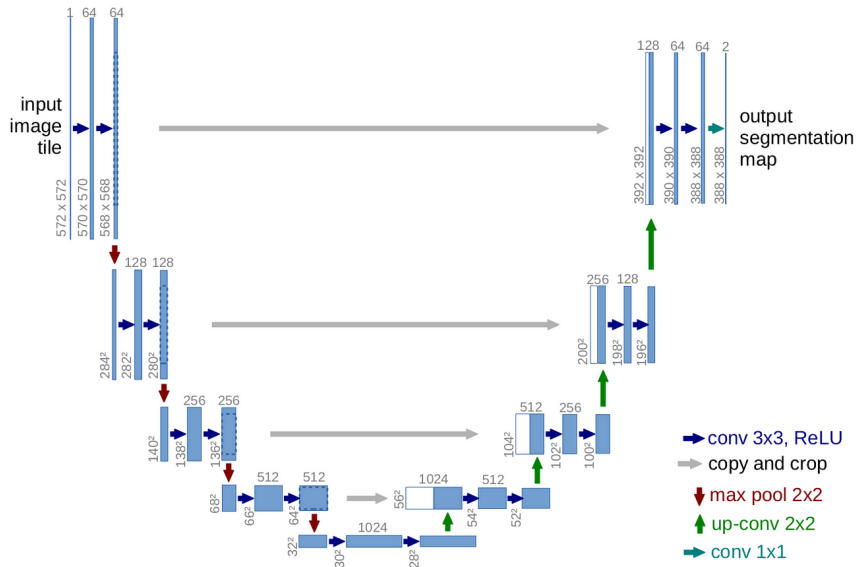
- Bishop, C. M. (2006). *Pattern recognition and machine learning* (Vol. 4). Springer.
- Bortoli, V. D. (2022). Convergence of denoising diffusion models under the manifold hypothesis. *Transactions on Machine Learning Research*.
- Choi, J., Kim, S., Jeong, Y., Gwon, Y., & Yoon, S. (2021). ILVR: Conditioning method for denoising diffusion probabilistic models. *ILVR*.
- 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., Thornton, J., Heng, J., & Doucet, A. (2021). Diffusion schrödinger bridge with applications to score-based generative modeling. *Advances in Neural Information Processing Systems*.
- Dhariwal, P., & Nichol, A. (2021). Diffusion models beat GANs on image synthesis. *Advances in Neural Information Processing Systems*.
- Feller, W. (1949). On the theory of stochastic processes, with particular reference to applications. *Proceedings of the [First] Berkeley Symposium on Mathematical Statistics and Probability*.
- Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*.



- Lugmayr, A., Danelljan, M., Romero, A., Yu, F., Timofte, R., & Van Gool, L. (2022). RePaint: Inpainting using denoising diffusion probabilistic models. *RePaint*.
- Nichol, A. Q., & Dhariwal, P. (2021). Improved denoising diffusion probabilistic models. In M. Meila & T. Zhang (Eds.), *Proceedings of the 38th international conference on machine learning* (pp. 8162–8171). PMLR. <https://proceedings.mlr.press/v139/nichol21a.html>
- 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.
- Saharia, C., Ho, J., Chan, W., Salimans, T., Fleet, D. J., & Norouzi, M. (2022). Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., & Ganguli, S. (2015). Deep unsupervised learning using nonequilibrium thermodynamics [ISSN: 1938-7228]. *Proceedings of the 32nd International Conference on Machine Learning*, 2256–2265. Retrieved 2022-12-01, from <https://proceedings.mlr.press/v37/sohl-dickstein15.html>
- Song, J., Meng, C., & Ermon, S. (2021). Denoising diffusion implicit models. *International Conference on Learning Representations*.
- Song, Y., & Ermon, S. (2019). Generative modeling by estimating gradients of the data distribution. *Advances in Neural Information Processing Systems*.

- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., & Poole, B. (2023). Score-based generative modeling through stochastic differential equations. *International Conference on Learning Representations*.
- Trippe, B. L., Yim, J., Tischer, D., Baker, D., Broderick, T., Barzilay, R., & Jaakkola, T. S. (2023). Diffusion probabilistic modeling of protein backbones in 3d for the motif-scaffolding problem. *International Conference on Learning Representations*.
- Xiao, Z., Kreis, K., & Vahdat, A. (2022). Tackling the generative learning trilemma with denoising diffusion GANs. *International Conference on Learning Representations*.
- Zhang, Y., Liu, W., Chen, Z., Li, K., & Wang, J. (2021). On the properties of kullback-leibler divergence between gaussians. *CoRR*, abs/2102.05485. <https://arxiv.org/abs/2102.05485>

- Architecture of U-NET: 63
- $\beta$  or  $\tilde{\beta}$ : 63
- Sketch of proof Pardoux, 1986 : 64



We can read in Ho et al., 2020, "Experimentally, both  $\sigma_t^2 = \beta_t$  and  $\sigma_t^2 = \tilde{\beta}_t$  had similar results. The first choice is optimal for  $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , and the second is optimal for  $\mathbf{x}_0$  deterministically set to one point. These are the two extreme choices corresponding to upper and lower bounds on reverse process entropy for data with coordinatewise unit variance Sohl-Dickstein et al., 2015."

1.  $\hat{\boldsymbol{w}}_t = \boldsymbol{w}_t - \boldsymbol{w}_1 - \int_t^1 \nabla \log p_s(\boldsymbol{x}_s) ds$  is a backward Brownian motion adapted to the filtration  $\sigma(\boldsymbol{x}_t \cup \{\boldsymbol{w}_s - \boldsymbol{w}_t, t \leq s \leq 1\})$  because it is a local backward martingale.
2. Rewriting of the SDE in the Stratonovich sense.
3. Insertion of  $\hat{\boldsymbol{w}}_t$  in the SDE.
4. Change of variable.