

# Generative Modeling with Diffusion Models

Juho Lee

Kim Jaechul Graduate School of AI, KAIST

*Foundation and Future of Generative Models: Mathematics, Algorithms, and Applications @ CityU*

# Table of Contents

## Motivation

Denoising Diffusion Probabilistic Models

From discrete to continuous time

Working in Latent Spaces

Diffusion Models for Learned Samplers

Summary & Conclusion

# Generative models

Generative modeling, density estimation, . . . ,

$$x_1, \dots, x_n \stackrel{\text{i.i.d.}}{\sim} p_{\text{data}}, \quad p_{\theta} \approx p_{\text{data}}? \quad (1)$$

Maximum likelihood training:

$$\begin{aligned} \theta^* &= \arg \min_{\theta} D_{\text{KL}}[p_{\text{data}} \| p_{\theta}], \\ \nabla_{\theta} D_{\text{KL}}[p_{\text{data}} \| p_{\theta}] &= -\mathbb{E}_{p_{\text{data}}}[\nabla_{\theta} \log p_{\theta}(x)]. \end{aligned} \quad (2)$$

# Generative models

How to build a model distribution  $p_{\theta}(x)$ ?

- ▶ Variational AutoEncoder (VAE) ([Kingma et al., 2014](#)).
- ▶ Autoregressive models ([Uria et al., 2016](#)).
- ▶ Flow-based models ([Rezende et al., 2015](#)).
- ▶ Implicit models (Generative Adversarial Network (GAN) ([Goodfellow et al., 2014](#))).

# Variational autoencoders

The Evidence Lower Bound (ELBO) of VAE:

$$\log p_{\theta}(x) \geq \mathbb{E}_{q_{\phi}(z|x)}[\log p_{\theta}(x|z)] - D_{\text{KL}}[q_{\phi}(z|x)||p_{\theta}(z)]. \quad (3)$$

- ▶ An **encoder**  $q_{\phi}(z|x)$  compresses a data point  $x$  into a code  $z$ .
- ▶ A **decoder**  $p_{\theta}(x|z)$  reconstructs  $x$  from  $z$ .
- ▶ The code  $z$  should be simple in a sense that  $q_{\phi}(z|x)$  is similar to  $p(z)$ .

# Autoregressive models

Let  $x := [x_1, \dots, x_d]^\top$  and  $x_{1:j} := [x_1, \dots, x_j]^\top$ .

$$p_\theta(x) = \prod_{j=1}^d p_\theta(x_j | x_{1:j-1}), \quad (4)$$

$$p_\theta(x_j | x_{1:j-1}) = \mathcal{N}(x_j | f_\mu(x_{1:j-1}), f_\sigma^2(x_{1:j-1})).$$

- ▶ Exact likelihood computation.
- ▶ Slow generation and order dependency. Rarely used for images.

# Flow-based models

Let  $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^d$  be a bijection with an easily computable Jacobian.

$$\begin{aligned} z &\sim p_0(z), \quad x = f_\theta(z), \\ p_\theta(x) &= p_0(f_\theta^{-1}(x)) \left| \det \frac{df_\theta^{-1}(x)}{dx} \right| \end{aligned} \tag{5}$$

- ▶ Exact likelihood computation.
- ▶ Fast generation.
- ▶ Limited flexibility due to the bijection constraint on  $f_\theta$ .

# Implicit models

$f_\theta : \mathbb{R}^h \rightarrow \mathbb{R}^d$  is an unconstrained function.

$$z \sim p_0(z), \quad x = f_\theta(z). \quad (6)$$

- ▶ Cannot compute the likelihood  $p_\theta(x)$ .
- ▶ Requires adversarial training, which can be unstable.
- ▶ Hard to get the code  $z$  corresponding to a given  $x$ .

# A better generative model?

- ▶ Minimal constraints on  $f_\theta$ .
- ▶ No adversarial training.
- ▶ Moderately fast generation.
- ▶ Handy manipulation of the code  $x$ .

# Table of Contents

Motivation

Denoising Diffusion Probabilistic Models

From discrete to continuous time

Working in Latent Spaces

Diffusion Models for Learned Samplers

Summary & Conclusion

# Challenges in VAE

- ▶ The tension between the encoder and decoder:
  - ▶ The encoder  $q_\phi(z|x)$  compresses  $x$  into a noise-like code  $z$ .
  - ▶ The decoder  $p_\theta(x|z)$  reconstructs  $x$  from  $z$ .
- ▶ It is generally hard to learn powerful mappings that achieve both good compression and sample fidelity.

# An idea: solve easier problems

- ▶ For the encoder, an easy way is to simply corrupt the input with noise.

$$z = \alpha x + \sigma \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, I) \Rightarrow q(z|x) = \mathcal{N}(z|\alpha x, \sigma^2 I). \quad (7)$$

- ▶ In terms of code simplicity, we want to choose  $\alpha = 0$  and  $\sigma = 1$ , but in that case, the decoder cannot reconstruct  $x$  from  $z$ .
- ▶ The main idea behind diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020) is to break down this problem of corruption and reconstruction as a **sequence of easier problems**.

## An idea: solve easier problems

Let  $(\alpha_k, \sigma_k)_{k=1}^K$  be a sequence such that the log Signal-to-Noise Ratio (SNR)  $\rho_k := \log(\alpha_k^2/\sigma_k^2)$  is non-increasing. Starting from  $x_0 := x$ , we decompose the corruption into smaller corruptions:

$$\begin{aligned}x_1 &= \alpha_1 x_0 + \sigma_1 \varepsilon_1, \\x_2 &= \frac{\alpha_2}{\alpha_1} x_1 + \sqrt{\alpha_2^2 (e^{-\rho_2} - e^{-\rho_1})} \varepsilon_2, \\&\vdots \\x_K &= \frac{\alpha_K}{\alpha_{K-1}} x_{K-1} + \sqrt{\alpha_K^2 (e^{-\rho_K} - e^{-\rho_{K-1}})} \varepsilon_K,\end{aligned}\tag{8}$$

leading to  $x_k \sim \mathcal{N}(\alpha_k x_0, \sigma_k^2 I)$  for  $k = 1, \dots, K$ .

# An idea: solve easier problems

The corruption sequence forms a **forward process** with a joint density,

$$q(x_{1:K} | x_0) = \prod_{k=1}^K q(x_k | x_{k-1}), \quad (9)$$
$$q(x_k | x_{k-1}) = \mathcal{N}(x_k | (\alpha_k / \alpha_{k-1})x_{k-1}, \alpha_k^2(e^{-\rho_k} - e^{-\rho_{k-1}})I).$$

To denoise (decode) the corrupted input, we build a sequence of denoisers,

$$p_\theta(x_{K-1} | x_K), p_\theta(x_{K-2} | x_{K-1}), \dots, p_\theta(x_1 | x_2), p_\theta(x | x_1), \quad (10)$$

defining a **reverse process** with a joint density,

$$p_\theta(x_{0:K}) = p_{\text{noise}}(x_K) \prod_{k=1}^K p_\theta(x_{k-1} | x_k), \quad (11)$$

where  $p_{\text{noise}}(x_K)$  is a prior for the terminal noise.

# Denosing diffusion probabilistic models

Denosing Diffusion Probabilistic Models (DDPM) (Ho et al., 2020), in a nutshell, is a VAE whose code  $z$  corresponds to the sequence of corrupted inputs  $z = x_{1:K}$ . It can thus be trained by maximizing the ELBO,

$$\log p_{\theta}(x) \geq \mathbb{E}_{q(x_{1:K}|x)} \left[ \log \frac{p_{\theta}(x, x_{1:K})}{q(x_{1:K}|x)} \right] := -\mathcal{L}(\theta, x). \quad (12)$$

The negative ELBO can further be decomposed as,

$$\mathcal{L}(x, \theta) = \mathbb{E}_q \left[ \log \frac{q(x_K|x)}{p_{\text{noise}}(x_K)} + \sum_{k=1}^K \log \frac{q(x_{k-1}|x_k, x)}{p_{\theta}(x_{k-1}|x_k)} \right]. \quad (13)$$

# Denosing diffusion probabilistic models

The roles of the terms in the objective:

$$\mathcal{L}(x, \theta) = \mathbb{E}_q \left[ \log \frac{q(x_K | x)}{p_{\text{noise}}(x_K)} + \sum_{k=1}^K \log \frac{q(x_{k-1} | x_k, x)}{p_{\theta}(x_{k-1} | x_k)} \right]. \quad (14)$$

- ▶ The first term measures the discrepancy between the prior noise distribution  $p_{\text{noise}}(x_K)$  and the terminal noise distribution  $q(x_K | x)$ .
- ▶ The second term measures the divergence between the oracle decoder  $q(x_{k-1} | x_k, x)$  and the intermediate decoder  $p_{\theta}(x_{k-1} | x_k)$ .

# Table of Contents

Motivation

Denoising Diffusion Probabilistic Models

From discrete to continuous time

Working in Latent Spaces

Diffusion Models for Learned Samplers

Summary & Conclusion

## What if $K \rightarrow \infty$ ?

- ▶ The key idea was to decompose  $x \rightarrow z$  into small-sized steps.
- ▶ Can we take more steps with smaller sizes?
- ▶  $K \rightarrow \infty$  with infinitesimal step sizes?

# A crash course on SDE

A Stochastic Differential Equation (SDE) on  $\mathbb{R}^d$  has the following form,

$$dX_t = f(X_t, t)dt + g(t)dW_t, \quad t \in [0, T]. \quad (15)$$

- ▶  $f : \mathbb{R}^d \times [0, T] \rightarrow \mathbb{R}^d$  is called a **drift** coefficient.
- ▶  $g : [0, T] \rightarrow \mathbb{R}$  is called a **diffusion** coefficient.
- ▶  $W_t$  is called a **Wiener process**, satisfying
  - ▶  $W_t$  has independent increments.
  - ▶  $W_t - W_s \sim \mathcal{N}(0, t - s)$  for all  $0 \leq s \leq t \leq T$ .

## A crash course on SDE

Just like the Ordinary Differential Equation (ODE), we can do something like,

$$X_t = X_0 + \underbrace{\int_0^t f(X_s, s) ds}_{:=A} + \underbrace{\int_0^t g(s) dW_s}_{:=B}. \quad (16)$$

- ▶ You can treat  $A$  as an ordinary integral, but technically it is a random variable (depending on  $\omega$  from the underlying probability space).

$$A(\omega) = \int_0^t f(X_s(\omega), s) ds, \quad \omega \in \Omega. \quad (17)$$

## A crash course on SDE

Just like the ODE, we can do something like,

$$X_t = X_0 + \underbrace{\int_0^t f(X_s, s) ds}_{:=A} + \underbrace{\int_0^t g(s) dW_s}_{:=B}. \quad (18)$$

- ▶  $B$  term requires an elaboration. You can roughly understand it as,

$$B(\omega) = \lim_{\Delta \rightarrow 0} \sum_{j \geq 0} g(t_j) [W_{t_{j+1}} - W_{t_j}](\omega), \quad (19)$$

where  $\Delta$  is the largest length among the segments  $\{[t_j, t_{j+1}]\}_{j \geq 0}$ .

# A crash course on SDE

Without rigour, we see that  $B$  is Gaussian with parameters

$$\begin{aligned}\mathbb{E}[B] &= \lim_{\Delta \rightarrow 0} \sum_{j \geq 0} g(t_j) \mathbb{E}[W_{t_{j+1}} - W_{t_j}] = 0, \\ \text{Var}[B] &= \lim_{\Delta \rightarrow 0} \sum_{j \geq 0} g^2(t_j) \text{Var}[W_{t_{j+1}} - W_{t_j}] \\ &= \lim_{\Delta \rightarrow 0} \sum_{j \geq 0} g^2(t_j) (t_{j+1} - t_j) \\ &= \int_0^t g^2(s) ds.\end{aligned}\tag{20}$$

# A crash course on SDE

The solution  $X_t$  of Eq. (15) is generally intractable. To obtain a numerical approximation, divide  $[0, T]$  into  $K$  bins and define for  $k = 0, \dots, K$ ,

$$t_k = k \cdot \delta \text{ where } \delta := T/K. \quad (21)$$

If we had only the drift term, that is,  $dX_t = f(X_t, t)dt$ , we consider a linear approximation of  $X_T$  around  $X_{T-\delta} = X_{t_{K-1}}$ :

$$X_T \approx X_{t_{K-1}} + \delta f(X_{t_{K-1}}, t_{K-1}). \quad (22)$$

Applying this approximation recursively, we get an iterative method,

$$X_{t_k} \approx X_{t_{k-1}} + \delta f(X_{t_{k-1}}, t_{k-1}) \text{ for } k = 1, \dots, K. \quad (23)$$

# A crash course on SDE

As we also have the Wiener process term. We know that

$$\int_{t_{k-1}}^{t_k} g(s) dW_s \sim \mathcal{N}\left(0, \int_{t_{k-1}}^{t_k} g^2(s) ds\right), \quad (24)$$

hence

$$\int_{t_{k-1}}^{t_k} g(s) dW_s \approx \sqrt{\delta} g(t_{k-1}) \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, I). \quad (25)$$

So the iterative approximation method for SDE become

$$X_{t_k} \approx X_{t_{k-1}} + \delta f(X_{t_{k-1}}, t_{k-1}) + \sqrt{\delta} g(t_{k-1}) \varepsilon. \quad (26)$$

# A crash course on SDE

- ▶ The law  $\mathbb{P}$  of the process  $X$  is called a **path measure** of  $X$ .
- ▶ We assume that the time-marginals  $\mathbb{P}_t$  admit a density  $p_t$  w.r.t. the Lebesgue measure.
- ▶ We also define joint and conditional densities, for instance,

$$\begin{aligned}\mathbb{P}(X_t \in A, X_s \in B) &= \iint_{A \times B} p_{t,s}(x, y) dx dy, \\ \mathbb{P}(X_t \in A | X_s = x) &= \int_A p_{t|s}(y|x) dy.\end{aligned}\tag{27}$$

# A crash course on SDE

A key question: how does  $p_t$  evolve over time?

## Theorem 1 (Fokker-Planck Equation)

Let  $(p_t)_{t \in [0, T]}$  be the time-marginal densities of  $X$  solving (15). Then  $p_t$  solves the *Fokker-Planck Equation (FPE)*,

$$\partial_t p_t(x) = -\nabla \cdot (f(x, t)p_t(x)) + \frac{g^2(t)}{2} \Delta p_t(x). \quad (28)$$

# A crash course on SDE

Are there simple cases where  $p_t$  or at least  $p_{t|s}$  is tractable?

## Proposition 1 (Affine drifts)

Assume  $f(x, t) = a(t)x$  in Eq. (15). Then for any  $0 \leq s \leq t \leq T$ ,

$$X_t | X_s = x \sim \mathcal{N}(\mu_{t|s}(x), \sigma_{t|s}^2 I), \quad (29)$$

where

$$\begin{aligned} \mu_{t|s}(x) &= \exp\left(\int_s^t a(u) du\right) x, \\ \sigma_{t|s}^2 &= \int_s^t \exp\left(2 \int_u^t a(v) dv\right) g^2(u) du. \end{aligned} \quad (30)$$

# From discrete to continuous time

Coming back to the discrete-time diffusion models,

$$q(x_k | x_{k-1}) = \mathcal{N}(x_k | (\alpha_k / \alpha_{k-1}) x_{k-1}, \alpha_k^2 (e^{-\rho_k} - e^{-\rho_{k-1}}) I). \quad (31)$$

Map each step  $k$  to  $t_k := k \cdot T/K$  in  $[0, T]$ . Let  $\alpha : [0, T] \rightarrow \mathbb{R}$  and  $\sigma : [0, T] \rightarrow [0, \infty)$  be differentiable functions where

$$\begin{aligned} \alpha_k &= \alpha(t_k), & \sigma_k &= \sigma(t_k) \text{ for } k = 0, \dots, K, \\ \alpha(0) &= 1, & \sigma(0) &= 0. \end{aligned} \quad (32)$$

Then

$$\frac{\alpha_k}{\alpha_{k-1}} = \delta \frac{\alpha(t_{k-1} + \delta) - \alpha(t_{k-1})}{\delta \alpha(t_{k-1})} + 1 \approx \delta \frac{\dot{\alpha}(t_{k-1})}{\alpha(t_{k-1})} + 1, \quad (33)$$

where  $\dot{f} = \frac{df}{dt}$  denotes the time derivative of a function and  $\delta := T/K$ .

# From discrete to continuous time

Similarly,

$$\begin{aligned}\alpha_k^2(e^{-\rho_k} - e^{-\rho_{k-1}}) &= \frac{\alpha_{k-1}^2(\sigma_k^2 - \sigma_{k-1}^2) - \sigma_{k-1}^2(\alpha_k^2 - \alpha_{k-1}^2)}{\alpha_{k-1}^2} \\ &\approx \delta\dot{\sigma}^2(t_{k-1}) - \delta\frac{\sigma^2(t_{k-1})}{\alpha^2(t_{k-1})}\dot{\alpha}^2(t_{k-1}) \\ &= -\delta\sigma^2(t_{k-1})\dot{\rho}(t_{k-1}).\end{aligned}\tag{34}$$

Here,  $\alpha(t)$  and  $\sigma(t)$  are properly chosen such that the log SNR  $\rho(t) = \log(\alpha^2(t)/\sigma^2(t))$  is a non-increasing function.

# From discrete to continuous time

Then the corruption density  $q(x_k | x_{k-1})$  can be rewritten as,

$$X_{t_k} = X_{t_{k-1}} + \delta \frac{\dot{\alpha}(t_{k-1})}{\alpha(t_{k-1})} X_{t_{k-1}} + \sqrt{-\delta \sigma^2(t_{k-1}) \dot{\rho}(t_{k-1})} \cdot \varepsilon, \quad (35)$$

where  $X_{t_k} = x_k$  and  $X_{t_{k-1}} = x_{k-1}$ . This is an approximate update step of the following SDE,

$$dX_t = \underbrace{\frac{\dot{\alpha}(t)}{\alpha(t)} X_t}_{:=f(X_t, t)} dt + \underbrace{\sqrt{-\sigma^2(t) \dot{\rho}(t)}}_{:=g(t)} dW_t, \quad X_0 \sim p_{\text{data}}. \quad (36)$$

By [Proposition 1](#), we recover the corruption equation we started from.

$$X_t | X_0 = x \sim \mathcal{N}(\alpha(t)x, \sigma^2(t)I), \quad (37)$$

# The generative process via time reversal

- ▶ The forward process corrupts data into noise. For generation, we must reverse this process.

$$dX_t = f(X_t, t)dt + g(t)dW_t, \quad X_0 \sim p_{\text{data}}. \quad (38)$$

- ▶ Let  $\mathbb{P}^*$  be the path measure of the forward process with time marginals  $(p_t^*)_{t \in [0, T]}$  and  $p_0^* = p_{\text{data}}$ . Find the **reverse-time SDE** that starts from  $p_T^*$  and transform it back to  $p_0^*$ .

$$\tilde{X}_0 = X_T \sim p_T^* \xrightarrow{\text{SDE}} \tilde{X}_T = X_0 \sim p_0^* = p_{\text{data}}, \quad (39)$$

where  $\tilde{f} := f(T - t)$  denotes the time-reversal.

# The generative process via time reversal

Starting from the FPE, using the identity  $\nabla p_t^* = p_t^* \nabla \log p_t^*$ ,

$$\begin{aligned}\partial_t p_t^* &= -\nabla \cdot (f p_t^*) + \frac{g^2}{2} \nabla \cdot (p_t^* \nabla \log p_t^*) \\ &= -\nabla \cdot ((f - g^2 \nabla \log p_t^*) p_t^*) - \frac{g^2}{2} \Delta p_t^*.\end{aligned}\tag{40}$$

Taking the change of variable  $t \leftarrow T - t$ ,

$$\partial_t \tilde{p}_t^* = \nabla \cdot ((\tilde{f} - \tilde{g}^2 \nabla \log \tilde{p}_t^*) \tilde{p}_t^*) + \frac{\tilde{g}^2}{2} \Delta \tilde{p}_t^*.\tag{41}$$

From [Theorem 1](#), we realize that  $\tilde{p}_t^*$  as the time-marginal of the SDE defining the reverse process ([Anderson, 1982](#); [Song, Sohl-Dickstein, et al., 2021](#)):

$$d\tilde{X}_t = -(\tilde{f} - \tilde{g}^2 \nabla \log \tilde{p}_t^*)(\tilde{X}_t, t) dt + \tilde{g}(t) d\tilde{W}_t, \quad \tilde{X}_0 \sim p_T^*.\tag{42}$$

# Score-based generative models

To implement the reverse SDE as a generative model, two components are required.

- ▶ The initial distribution  $p_T^*$ : this is generally unknown, but we can approximate it with a simple noise distribution  $p_{\text{noise}}$ .
- ▶ The **score function**  $\nabla \log \tilde{p}_t^*$ : this is generally intractable, so approximated via a neural network  $s_\theta(x, t)$ .

How to learn the parameter  $\theta$  of the score network  $s_\theta$ ?

# Score-based generative models

Based on two choices, we construct a **score-based generative model**,

$$d\tilde{X}_t^\theta = -(\tilde{f} - \tilde{g}^2\tilde{\varsigma}_\theta)(\tilde{X}_t^\theta, t)dt + \tilde{g}(t)d\tilde{W}_t^\theta, \quad \tilde{X}_0^\theta \sim p_{\text{noise}}. \quad (43)$$

Let  $\tilde{\mathbb{P}}^\theta$  be the path measure associated with  $\tilde{X}^\theta$ . We want to minimize the KL divergence between them,

$$D_{\text{KL}}[\tilde{\mathbb{P}}^\star \|\tilde{\mathbb{P}}^\theta] = \mathbb{E}_{\tilde{\mathbb{P}}^\star} \left[ \log \frac{d\tilde{\mathbb{P}}^\star}{d\tilde{\mathbb{P}}^\theta}(\tilde{X}) \right]. \quad (44)$$

How to compute the KL divergence between the path measures?

# Computing divergence between path measures

The KL divergence between two path measures can be decomposed into the KL divergence between the initial distributions and the expected KL divergence between the paths afterwards ([Léonard, 2014](#))

$$D_{\text{KL}}[\tilde{\mathbb{P}}^* \|\tilde{\mathbb{P}}^\theta] = D_{\text{KL}}[p_T^* \|\mathit{p}_{\text{noise}}] + \mathbb{E}_{p_T^*} \left[ D_{\text{KL}}[\tilde{\mathbb{P}}^*(\cdot | \tilde{X}_0) \|\tilde{\mathbb{P}}^\theta(\cdot | \tilde{X}_0)] \right]. \quad (45)$$

The first term is irreducible, so we must focus on the second term.

# Score matching for generative models

## Theorem 2 (Girsanov Theorem, Informal)

Consider two SDEs,

$$\begin{aligned}dX_t^a &= a(X_t, t)dt + \sigma(t)dW_t^a, \\dX_t^b &= b(X_t, t)dt + \sigma(t)dW_t^b,\end{aligned}\tag{46}$$

and let  $\mathbb{P}^a$  and  $\mathbb{P}^b$  denote the path measures of  $X^a$  and  $X^b$ , respectively.

Assume that  $h := (b - a)/\sigma$  satisfies  $\mathbb{E}_{\mathbb{P}^a}[\exp(\frac{1}{2} \int_0^T h_t^2 dt)] < \infty$ . Then the Radon-Nikodym Derivative (RND) between the path measures is given as,

$$\log \frac{d\mathbb{P}^b}{d\mathbb{P}^a}(X) = \int_0^T h(X_t, t) \cdot dW_t^a - \frac{1}{2} \int_0^T \|h\|^2(X_t, t) dt.\tag{47}$$

# Score matching for generative models

Applying [Theorem 2](#) to  $\tilde{\mathbb{P}}$  and  $\tilde{\mathbb{P}}^\theta$ ,

$$\mathbb{E}_{\tilde{\mathbb{P}}} \left[ \log \frac{d\tilde{\mathbb{P}}}{d\tilde{\mathbb{P}}^\theta}(\tilde{X}) \right] = \frac{1}{2} \mathbb{E}_{\tilde{\mathbb{P}}} \left[ \int_0^T g^2(t) \|\tilde{\mathbf{s}}_\theta(\tilde{X}_t, t) - \nabla \log \tilde{p}_t^*(\tilde{X}_t)\|^2 dt \right]. \quad (48)$$

That is, minimizing the divergence between the path measure amounts to matching the score network  $\mathbf{s}_\theta$  to the true score function  $\nabla \log p_t$ . The integral along the time axis can be written as an expectation, so we write

$$\mathcal{L}_{\text{SM}}(\theta; g) := \frac{1}{2} \mathbb{E}_{\mathcal{U}(0, T), p_t^*} \left[ g^2(t) \|\mathbf{s}_\theta(X_t, t) - \nabla \log p_t^*(X_t)\|^2 \right], \quad (49)$$

where  $\mathcal{U}(0, T)$  denotes a uniform distribution over  $[0, T]$ .

# Score matching for generative models

- ▶ The weight  $g$  can be replaced by an arbitrary weight  $w(t)$  without changing the optimal solution (i.e.,  $\mathbf{s}_\theta(X_t, t) = \nabla \log p_t^*(X_t)$  for all  $t \in [0, T]$ ).

$$\mathcal{L}_{\text{SM}}(\theta; w) = \frac{1}{2} \mathbb{E}_{\mathcal{U}(0, T), p_t^*} [w^2(t) \|\mathbf{s}_\theta(X_t, t) - \nabla \log p_t^*(X_t)\|^2]. \quad (50)$$

- ▶ For instance, [Ho et al. \(2020\)](#) suggested to set  $w(t) = 1$  for all  $t \in [0, T]$ , leading to better quality for image generative models.
- ▶ However, [Eq. \(50\)](#) is still intractable due to the score  $\nabla \log p_t^*(x)$ .

# Conditional score matching

The marginal score is an expectation of the conditional score:

$$\nabla \log p_t^*(x) = \mathbb{E}_{p_{0|t}^*} \left[ \nabla \log p_{t|0}^*(x | X_0) \right]. \quad (51)$$

Based on this observation,

$$\begin{aligned} & \mathbb{E}_{p_t^*} \left[ \|\mathbf{s}_\theta(X_t, t) - \nabla \log p_t^*(X_t)\|^2 \right] \\ & \stackrel{c}{=} \int p_t^*(x_t) \left( \|\mathbf{s}_\theta(x_t, t)\|^2 - \mathbf{s}_\theta(x_t, t) \cdot \nabla \log p_t(x_t) \right) dx_t \\ & \stackrel{c}{=} \int p_{t|0}^*(x_t) p_{\text{data}}(x_0) \left( \|\mathbf{s}_\theta(x_t, t)\|^2 - \mathbf{s}_\theta(x_t, t) \cdot \nabla \log p_{t|0}^*(x_t | x_0) \right) dx_t dx_0 \\ & \stackrel{c}{=} \mathbb{E}_{p_{t|0}^* p_{\text{data}}} \left[ \|\mathbf{s}_\theta(X_t, t) - \nabla \log p_{t|0}^*(X_t | X_0)\|^2 \right], \end{aligned} \quad (52)$$

where  $\stackrel{c}{=}$  denotes the equality up to a constant.

# Conditional score matching

Since  $p_{t|0}^*$  is Gaussian, its score has a simple analytic form:

$$\nabla \log p_{t|0}^*(x_t | x_0) = -\frac{x_t - \alpha(t)x_0}{\sigma^2(t)}. \quad (53)$$

The conditional score matching loss is thus defined as,

$$\mathcal{L}_{\text{CSM}}(\theta; w) = \frac{1}{2} \mathbb{E}_{\mathcal{U}(0,T), p_{t|0}^*, p_{\text{data}}} \left[ w^2(t) \left\| \mathbf{s}_\theta(X_t, t) + \frac{X_t - \alpha(t)X_0}{\sigma^2(t)} \right\|^2 \right]. \quad (54)$$

# Conditional score matching

A step of training:

- ▶ Draw  $t \sim \mathcal{U}(0, T)$ .
- ▶ Draw  $X_0$  from training dataset and corrupt by  $X_t = \alpha(t)X_0 + \sigma(t)\varepsilon_t$ .
- ▶ Compute  $w^2(t)\nabla_{\theta} \left\| \mathbf{s}_{\theta}(X_t, t) + \frac{X_t - \alpha(t)X_0}{\sigma^2(t)} \right\|^2$  and update.

This training procedure is highly scalable because,

- ▶ Everything involved in unbiased gradient estimation is tractable.
- ▶ Simulation-free: we don't need to simulate any SDE during the training. No backpropagation through numerical SDE solvers.

# Parameterizations for score matching

We may also use alternative parameterizations for the score network  $\mathbf{s}_\theta$ .

- ▶ Noise prediction model: since  $\nabla \log p_{t|0}^*(x_t | x_0) = -\frac{\varepsilon_t}{\sigma(t)}$ , we can parameterize  $\mathbf{s}_\theta(x, t) = -\frac{\mathbf{e}_\theta(x, t)}{\sigma(t)}$  and minimize

$$\frac{1}{2} \mathbb{E}_{\mathcal{U}(0, T), p_{t|0}^* p_{\text{data}}} \left[ w^2(t) \|\mathbf{e}_\theta(X_t, t) - \varepsilon_t\|^2 \right]. \quad (55)$$

Ho et al. (2020) proposed to use this one with  $w(t) = 1$ .

- ▶ Denoiser: parameterize  $\mathbf{s}_\theta(x, t) = -\frac{x_t - \alpha(t)\mathbf{x}_\theta(x, t)}{\sigma^2(t)}$  and minimize

$$\frac{1}{2} \mathbb{E}_{\mathcal{U}(0, T), p_{t|0}^* p_{\text{data}}} \left[ w^2(t) \|\mathbf{x}_\theta(X_t, t) - X_0\|^2 \right]. \quad (56)$$

# Examples of score-based generative models

Existing score-based or diffusion-based generative models can be described as a special case of the framework we discussed so far, with varying choices of the functions  $\alpha(t)$  and  $\sigma(t)$  (Song, Sohl-Dickstein, et al., 2021).

- ▶ Variance Exploding SDE (VE-SDE): set  $\alpha(t) = 0$ , leading to

$$dX_t = \sqrt{\sigma^2(t)}dW_t. \quad (57)$$

Song and Ermon (2019) proposed a discrete time version of VE-SDE.

- ▶ Variance Preserving SDE (VP-SDE): set  $\alpha(t) = e^{-\frac{1}{2} \int_0^t \beta(s)ds}$  for some function  $\beta$ , and  $\sigma^2(t) = 1 - \alpha^2(t)$ . This lead to,

$$dX_t = -\frac{1}{2}\beta(t)X_tdt + \sqrt{\beta(t)}dW_t, \quad (58)$$

which is a continuous-time extension of DDPM (Ho et al., 2020).

# Sampling from score-based generative models

Once trained  $s_\theta$ , to generate a new sample, we have solve the SDE,

$$d\tilde{X}_t^\theta = -(\tilde{f} - \tilde{g}^2\tilde{s}_\theta)(\tilde{X}_t^\theta, t)dt + \tilde{g}(t)d\tilde{W}_t^\theta, \quad \tilde{X}_0^\theta \sim p_{\text{noise}}. \quad (59)$$

This requires a numerical solver, for instance, the linear approximation method (Euler-Maruyama solver):

$$\tilde{X}_{t_k}^\theta = \tilde{X}_{t_{k-1}}^\theta - \delta(\tilde{f} - \tilde{g}^2\tilde{s}_\theta)(\tilde{X}_{t_{k-1}}^\theta, t_{k-1}) + \sqrt{\delta}\tilde{g}(t_{k-1})\varepsilon_{k-1}. \quad (60)$$

Each step of the numerical solver requires a forward pass through a neural network  $s_\theta$ , so reducing the Number of Function Evaluations (NFE) is crucial for practical use.

# Sampling by solving ODEs

We can manipulate the same FPE a bit differently,

$$\begin{aligned}\partial_t p_t &= -\nabla_z \cdot (f p_t) + \frac{g^2}{2} \nabla \cdot (p_t \nabla \log p_t) \\ &= -\nabla \cdot \left( \left( f - \frac{g^2}{2} \nabla \log p_t \right) p_t \right).\end{aligned}\tag{61}$$

With a change of variable, we get another reverse process,

$$\partial_t \tilde{p}_t = \nabla \cdot \left( \left( \tilde{f} - \frac{\tilde{g}^2}{2} \nabla \log \tilde{p}_t \right) \tilde{p}_t \right),\tag{62}$$

which corresponds to an ODE,

$$d\tilde{X}_t = - \left( \tilde{f} - \frac{\tilde{g}^2}{2} \nabla \log \tilde{p}_t \right) (\tilde{X}_t, t) dt.\tag{63}$$

This is called a **Probability Flow ODE (PF-ODE)** (Song, Sohl-Dickstein, et al., 2021).

# Sampling by solving ODEs

Both the reverse SDE and PF-ODE share the same FPES, so the time marginals of them are the same.

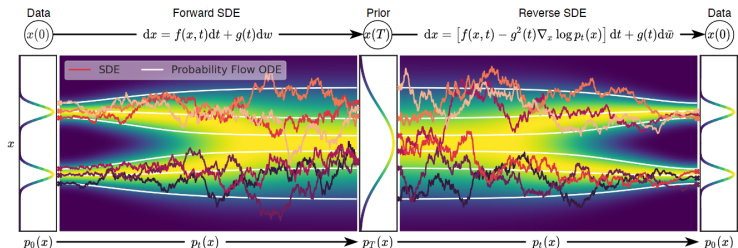


Figure 1: Figure from [Song, Sohl-Dickstein, et al. \(2021\)](#).

# Sampling by solving ODEs

PF-ODE vs. reverse SDE?

- ▶ A single step of PF-ODE: take a deterministic step with  $f$  and  $\mathbf{s}_\theta$ .

$$\tilde{X}_{t_k}^\theta = \tilde{X}_{t_{k-1}}^\theta - \delta \left( \tilde{f} - \frac{\tilde{g}^2}{2} \tilde{\mathbf{s}}_\theta \right) (\tilde{X}_{t_{k-1}}^\theta, t_{k-1}). \quad (64)$$

- ▶ A single step of reverse SDE: take a double-length step with  $\mathbf{s}_\theta$  in the reverse SDE, then correct it with Gaussian noise.

$$\tilde{X}_{t_k}^\theta = \tilde{X}_{t_{k-1}}^\theta - \delta \left( \tilde{f} - \tilde{g}^2 \tilde{\mathbf{s}}_\theta \right) (\tilde{X}_{t_{k-1}}^\theta, t_{k-1}) + \sqrt{\delta} \tilde{g}(t_{k-1}) \varepsilon_{k-1}. \quad (65)$$

In general, the reverse SDE takes longer to converge but produces higher-quality samples. On the other hand, while the PF-ODE may require fewer steps for convergence, its sample quality over time is typically inferior to that of the reverse SDE.

# Table of Contents

Motivation

Denoising Diffusion Probabilistic Models

From discrete to continuous time

Working in Latent Spaces

Diffusion Models for Learned Samplers

Summary & Conclusion

# Latent diffusion models

Vahdat et al. (2021) proposes Latent Score-based Generative Model (LSGM), where the SDE is defined on a latent space.

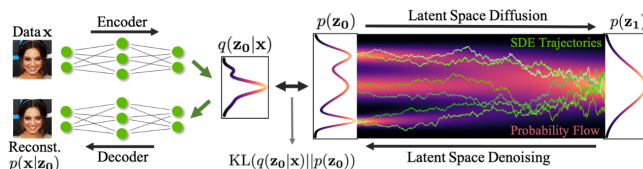


Figure 2: Vahdat et al. (2021).

# Latent diffusion models

Rombach et al. (2022) introduces Latent Diffusion Model (LDM), a discrete variant of LSGM enhanced with various engineering optimizations to enable scalability, forming the foundation of Stable Diffusion.

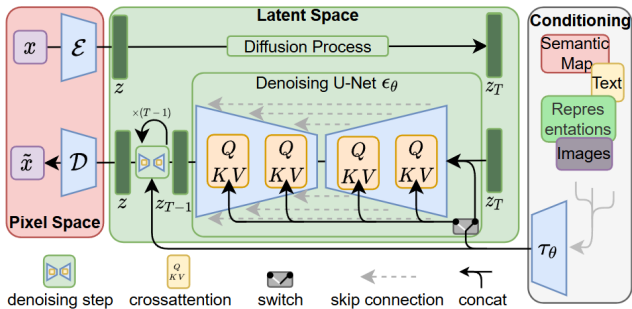


Figure 3: Rombach et al. (2022).

# Latent diffusion models

The forward process of LDM is defined as,

$$x \sim p_{\text{data}}, \quad z_0 = \mathcal{E}(x), \quad q(z_k | z_0) = \mathcal{N}(z_k | \alpha_k z_0, \sigma_k^2 I), \quad (66)$$

where  $\dim(z_k) \ll \dim(x)$ . The backward process to generate a new sample is then

$$\begin{aligned} \tilde{z}_0 &\sim p_{\text{noise}}, \\ \tilde{z}_k &= \tilde{z}_{k-1} - \delta(\bar{f} - \bar{g}^2 \bar{s}_\theta)(\tilde{z}_{k-1}, t_{k-1}) + \sqrt{\delta} \bar{g}(t_{k-1}) \varepsilon_{k-1}, \\ x &= \mathcal{D}(\tilde{z}_K). \end{aligned} \quad (67)$$

# Training LDM

- ▶ [Vahdat et al. \(2021\)](#) derive a variational inference algorithm that jointly learns the autoencoder ( $\mathcal{E}, \mathcal{D}$ ) and the score network. However, this method suffers from learning instability.
- ▶ [Rombach et al. \(2022\)](#) introduces an alternative strategy that pretrains the autoencoder and keeps it fixed during the training of the score network. This method has been shown to significantly improve stability, scalability, and generalization.

# Advantages of LDM

Working in pixel space wastes computation on imperceptible details. In the latent space, these details are already encoded, allowing the model to focus solely on semantic compression.

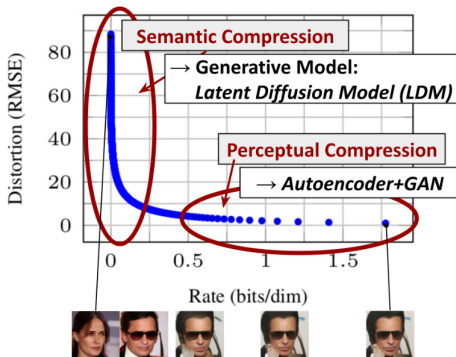


Figure 4: [Rombach et al. \(2022\)](#).

# Advantages of LDM

Many real-world applications require conditioning on some information  $y$  (e.g., class conditional, text-to-image, image-to-image translation).

- ▶ Unconditional Diffusion Model (DM):

$$\nabla_z \log p_{t_k}(z_k) \approx \mathbf{s}_\theta(z_k, t_k). \quad (68)$$

- ▶ Conditional DM:

$$\nabla_z \log p_{t_k}(z_k | \mathbf{y}) \approx \mathbf{s}_\theta(z_k, t_k, \mathbf{y}). \quad (69)$$

# Advantages of LDM

LDMs allow more flexible manipulation of conditioning by learning a shared latent space that seamlessly integrates both images and conditioning information.

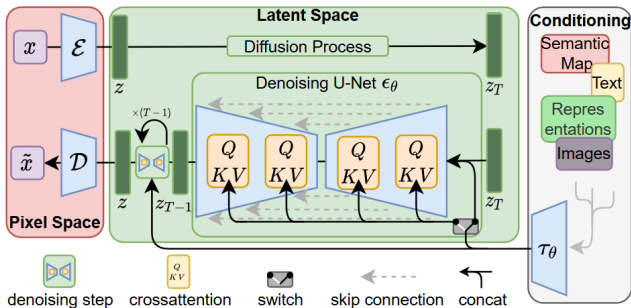


Figure 5: Rombach et al. (2022).

# Table of Contents

Motivation

Denoising Diffusion Probabilistic Models

From discrete to continuous time

Working in Latent Spaces

Diffusion Models for Learned Samplers

Summary & Conclusion

# The sampling problem

- ▶ So far, we have discussed generative modeling, where i.i.d. samples from a target distribution are given as a training set.
- ▶ In many scenarios (especially in Bayesian inference), we are often interested in sampling from a target distribution specified by an unnormalized density,

$$\pi(x) = \frac{\gamma(x)}{Z}, \quad (70)$$

where the normalization constant  $Z$  is typically unknown.

# The sampling problem

- ▶ Typical approaches to sample from a unnormalized density  $\gamma(x)$ ?
- ▶ MCMC (Hamiltonian Monte Carlo (HMC) (Neal, 2010) when  $\gamma$  is differentiable), without elaboration, may suffer mixing for high-dimensional and multi-modal distributions.
- ▶ Vanila MCMC is similar to “vanila” VAE in a sense that it is trying to “directly” sample from the target distribution.

# Annealed importance sampling

- ▶ Annealed Importance Sampling (AIS) (Neal, 2001) introduces a sequence of annealed distributions and sample from a path that gradually transitions from an easy one to the target.
- ▶ Neal (2001) proposed a geometric annealing path as an example,

$$\pi_k(x) \propto \pi_0(x)^{1-\beta_k} \gamma(x)^{\beta_k}, \quad (71)$$

where  $\pi_0$  is a easy to simulate density and  $0 = \beta_0 < \dots < \beta_K = 1$ .

- ▶ Similar to diffusion models, AIS breaks down the complex sampling problems into a sequence of simpler problems.

# Diffusions for learned samplers

- ▶ The idea behind AIS is related to diffusion or score-based models.
- ▶ Based on the connection, many works ([Doucet et al., 2022](#); [Vargas et al., 2023](#); [Richter et al., 2024](#); [Chen et al., 2024](#)) proposed learning diffusion models for sampling from an unnormalized density.

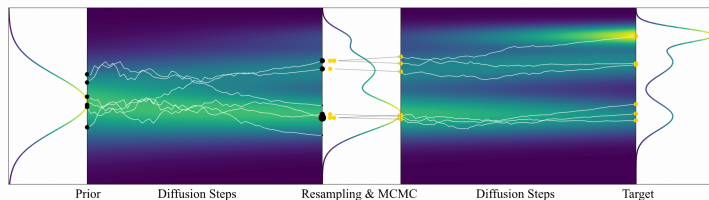


Figure 6: Figure from [Chen et al. \(2024\)](#).

# Denosing diffusion samplers

- ▶ Denoising Diffusion Sampler (DDS) (Vargas et al., 2023) adapts VP-SDE (continuous-time DDPM) for the sampling problem.
- ▶ Take the time-reversal of the forward process, starting from  $\pi$ , as a target process.

$$\begin{aligned}dX_t &= f(X_t, t)dt + g(t)dW_t, & X_0 &\sim \pi, \\d\tilde{X}_t &= -(\tilde{f} - \tilde{g}^2 \nabla \log \tilde{p}_t^*)(\tilde{X}_t, t)dt + \tilde{g}(t)d\tilde{W}_t, & \tilde{X}_0 &\sim p_T^*.\end{aligned}\tag{72}$$

- ▶ The sampling process is parameterized as,

$$d\tilde{X}_t^\theta = -(\tilde{f} - \tilde{g}^2 \tilde{s}_\theta)(\tilde{X}_t^\theta, t)dt + \tilde{g}(t)d\tilde{W}_t^\theta, \quad \tilde{X}_0^\theta \sim p_{\text{noise}}.\tag{73}$$

# The objective for DDS

Let  $\tilde{\mathbb{P}}^\star$  and  $\tilde{\mathbb{P}}^\theta$  be path measures of the target and sampling processes, respectively. In generative modeling, we could minimize

$$D_{\text{KL}}[\tilde{\mathbb{P}}^\star \|\tilde{\mathbb{P}}^\theta] = \mathbb{E}_{\tilde{\mathbb{P}}^\star} \left[ \log \frac{d\tilde{\mathbb{P}}^\star}{d\tilde{\mathbb{P}}^\theta} \right], \quad (74)$$

whose unbiased estimate could effectively be estimated as below,

$$\tilde{X}_T \sim p_{\text{data}}, \quad \tilde{X}_t | \tilde{X}_T \sim p_{T-t|0}^\star \text{ for } t \in [0, T]. \quad (75)$$

However, in the sampling problem, we cannot do the first step, because we cannot sample from  $\pi$ !

# The objective for DDS

To facilitate learning, we introduce a **reference process**,

$$dX_t^\circ = f(X_t^\circ, t)dt + g(t)dW_t^\circ, \quad X_0^\circ \sim p_0^\circ, \quad (76)$$

Here,  $p_0^\circ$  is chosen such that  $\nabla \log p_t^\circ$  is analytic for all  $t \in [0, T]$ . For instance, in DDS, we choose

$$f(x, t) = -\frac{\beta(t)}{2}x, \quad g(t) = \sigma\sqrt{\beta(t)}, \quad p_0^\circ(x) = \mathcal{N}(x|0, \sigma^2 I), \quad (77)$$

leading to

$$p_t^\circ(x) = \mathcal{N}(x|0, \sigma^2 I) \text{ for all } t \in [0, T], \quad \tilde{f} - \tilde{g}^2 \nabla \log \tilde{p}_t^\circ = -\tilde{f}, \quad (78)$$
$$d\tilde{X}_t^\circ = \tilde{f}(\tilde{X}_t^\circ, t)dt + \tilde{g}(t)d\tilde{W}_t^\circ, \quad \tilde{X}_0^\circ \sim p_T^\circ = p_0^\circ.$$

# The objective for DDS

With the help of the reference process, we can write

$$D_{\text{KL}}[\mathbb{P}^\theta \|\mathbb{P}^\star] = D_{\text{KL}}[\mathbb{P}^\theta \|\mathbb{P}^\circ] - \mathbb{E}_{\mathbb{P}^\theta} \left[ \log \frac{d\mathbb{P}^\star}{d\mathbb{P}^\circ} \right]. \quad (79)$$

Here, we have

$$\log \frac{d\mathbb{P}^\star}{d\mathbb{P}^\circ}(X^\theta) = \log \frac{\pi(X_0^\theta)}{p_0^\circ(X_0^\theta)} = \log \frac{\pi(\tilde{X}_T^\theta)}{p_0^\circ(\tilde{X}_T^\theta)}. \quad (80)$$

The first KL term can be computed with [Theorem 2](#),

$$\begin{aligned} D_{\text{KL}}[\mathbb{P}^\theta \|\mathbb{P}^\circ] &= D_{\text{KL}}[\tilde{\mathbb{P}}^\theta \|\tilde{\mathbb{P}}^\circ] \\ &= \frac{1}{2} \mathbb{E}_{\tilde{\mathbb{P}}^\theta} \left[ \int_0^T \tilde{g}(t)^2 \|\tilde{\mathbf{s}}_\theta(\tilde{X}^\theta, t) - \nabla \log \tilde{p}_t^\circ(\tilde{X}^\theta)\|^2 \right]. \end{aligned} \quad (81)$$

# The objective for DDS

We can consider a parameterization,

$$\mathbf{u}_\theta(x, t) := g(t)(\mathbf{s}_\theta(x, t) - \nabla \log p_t^\circ(x)) = \sigma \sqrt{\beta(t)}(\mathbf{s}_\theta(x, t) + x/\sigma^2), \quad (82)$$

leading to the overall objective,

$$D_{\text{KL}}[\mathbb{P}^\theta \|\mathbb{P}^\star] = \mathbb{E}_{\tilde{\mathbb{P}}^\theta} \left[ \frac{1}{2} \int_0^T \|\tilde{\mathbf{u}}_\theta(\tilde{X}_t^\theta, t)\|^2 dt + \log \frac{\pi(\tilde{X}_T^\theta)}{p_0^\circ(\tilde{X}_T)} \right]. \quad (83)$$

The corresponding objective is actually a special instance of a **Schrödinger bridge problem** (Schrödinger, 1931; Schrödinger, 1932; Léonard, 2013), an entropy-regularized optimal transport problem.

# Table of Contents

Motivation

Denoising Diffusion Probabilistic Models

From discrete to continuous time

Working in Latent Spaces

Diffusion Models for Learned Samplers

Summary & Conclusion

# Diffusion models are everywhere!

- ▶ Generative models for various domains:
  - ▶ Image synthesis.
  - ▶ Molecule, protein, and drug discovery.
  - ▶ Improving on language models.
- ▶ Learned samplers.
  - ▶ Sampling from Boltzmann distributions.
  - ▶ Approximate Bayesian inference.
  - ▶ Solving optimal transport problems (e.g., distribution to distribution translation problems).

# Key takeaways

Two things to learn from the methods discussed today:

- ▶ Break down a complex problem into easier problems.
- ▶ Tricks to keep learning tractable and scalable, for instance, conditional matching.

# References I

- Anderson, B. D. O. (1982). "Reverse-time diffusion equation models". In: *Stochastic Processes and their Applications* 12.3, pp. 313–326 (cit. on p. 32).
- Chen, J. et al. (2024). "Sequential controlled Langevin diffusions". In: *Proceedings of The 13th International Conference on Learning Representations (ICLR 2025)* (cit. on p. 60).
- Doucet, A. et al. (2022). "Score-based diffusion meets annealed importance sampling". In: *Advances in Neural Information Processing Systems 35 (NeurIPS 2022)* (cit. on p. 60).
- Goodfellow, I. et al. (2014). "Generative Adversarial Nets". In: *Advances in Neural Information Processing Systems 27 (NIPS 2014)* (cit. on p. 4).
- Ho, J., A. Jain, and P. Abbeel (2020). "Denoising diffusion probabilistic models". In: *Advances in Neural Information Processing Systems 33 (NeurIPS 2020)* (cit. on pp. 12, 15, 38, 42, 43).
- Kingma, D. P. and M. Welling (2014). "Auto-encoding variational Bayes". In: *Proceedings of The 2nd International Conference on Learning Representations (ICLR 2014)* (cit. on p. 4).
- Léonard, C. (2013). "A survey of the Schrödinger problem and some of its connections with optimal transport". In: *Discrete and Continuous Dynamical Systems* 34.4, pp. 1533–1574 (cit. on p. 65).
- – (2014). "Some properties of path measures". In: *Séminaire de Probabilités XLVI*. Vol. 2123. Lecture Notes in Mathematics. Springer, pp. 207–230 (cit. on p. 35).
- Neal, R. M. (2001). "Annealed importance sampling". In: *Statistics and Computing* 11.2, pp. 125–139 (cit. on p. 59).
- – (2010). "MCMC using Hamiltonian dynamics". In: *Handbook of Markov Chain Monte Carlo* 54, pp. 113–162 (cit. on p. 58).

# References II

- Rezende, D. and S. Mohamed (2015). "Variational inference with normalizing flows". In: *Proceedings of The 32th International Conference on Machine Learning (ICML 2015)* (cit. on p. 4).
- Richter, L. and J. Berner (2024). "Improved sampling via learned diffusions". In: *Proceedings of The 12th International Conference on Learning Representations (ICLR 2024)* (cit. on p. 60).
- Rombach, R. et al. (2022). "High-resolution image synthesis with latent diffusion models". In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (cit. on pp. 50, 52, 53, 55).
- Schrödinger, Erwin (1931). "Über die Umkehrung der Naturgesetze". In: *Sitzungsberichte der Preussischen Akademie der Wissenschaften, Physikalisch-mathematische Klasse*, pp. 144–153 (cit. on p. 65).
- – (1932). "Sur la théorie relativiste de l'électron et l'interprétation de la mécanique quantique". In: *Annales de l'Institut Henri Poincaré* 2.4, pp. 269–310 (cit. on p. 65).
- Sohl-Dickstein, J. et al. (2015). "Deep unsupervised learning using nonequilibrium thermodynamics". In: *Proceedings of The 32th International Conference on Machine Learning (ICML 2015)* (cit. on p. 12).
- Song, Y. and S. Ermon (2019). "Generative modeling by estimating gradients of the data distribution". In: *Advances in Neural Information Processing Systems 32 (NeurIPS 2019)* (cit. on p. 43).
- Song, Y., J. Sohl-Dickstein, et al. (2021). "Score-based generative modeling through stochastic differential equations". In: *Proceedings of The 9th International Conference on Learning Representations (ICLR 2021)* (cit. on pp. 32, 43, 45, 46).
- Uria, B. et al. (2016). "Neural autoregressive distribution estimation". In: *Journal of Machine Learning Research* 17.205, pp. 1–37 (cit. on p. 4).

# References III

- Vahdat, A., K. Kreis, and J. Kautz (2021). "Score-based generative modeling in latent space". In: *Advances in Neural Information Processing Systems 34 (NeurIPS 2021)* (cit. on pp. 49, 52).
- Vargas, F., W. Grathwohl, and A. Doucet (2023). "Denoising diffusion samplers". In: *Proceedings of The 11th International Conference on Learning Representations (ICLR 2023)* (cit. on pp. 60, 61).