

Deep Generative Models

14. Denoising Diffusion Probabilistic Models



• 국가수리과학연구소 산업수학혁신센터 김민중

Recap. of score-based model

- Fisher divergence between $p(\mathbf{x})$ and $q(\mathbf{x})$:

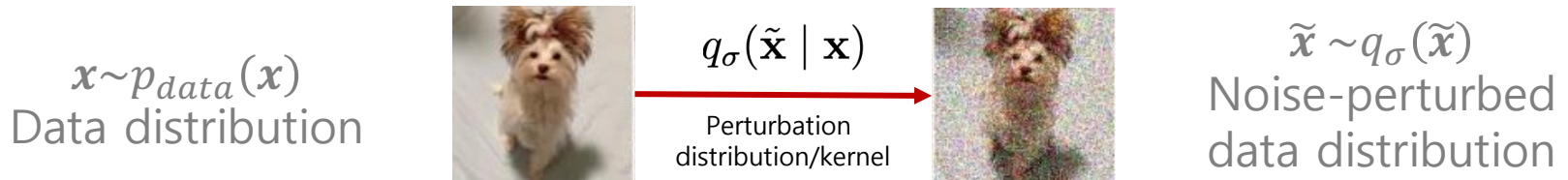
$$D_F(p, q) := \frac{1}{2} E_{\mathbf{x} \sim p} [\|\nabla_{\mathbf{x}} \log p(\mathbf{x}) - \nabla_{\mathbf{x}} \log q(\mathbf{x})\|_2^2]$$

- Score matching (Hyvärinen, 2005)

$$\begin{aligned} & \frac{1}{2} E_{\mathbf{x} \sim p_{data}} [\|\mathbf{s}_{\theta}(\mathbf{x}) - \nabla_{\mathbf{x}} \log p_{data}(\mathbf{x})\|_2^2] \\ &= E_{\mathbf{x} \sim p_{data}} \left[\frac{1}{2} \|\mathbf{s}_{\theta}(\mathbf{x})\|_2^2 + \text{tr}(\nabla_{\mathbf{x}} \mathbf{s}_{\theta}(\mathbf{x})) \right] + \text{const.} \end{aligned}$$

- Not scalable for deep score-based models and high dimensional data

Denoising score matching with Langevin dynamics



$$\begin{aligned} & E_{\tilde{\mathbf{x}} \sim q_{\sigma}} [\|\nabla_{\tilde{\mathbf{x}}} \log q_{\sigma}(\tilde{\mathbf{x}}) - \mathbf{s}_{\theta}(\tilde{\mathbf{x}})\|_2^2] \\ &= E_{\mathbf{x} \sim p_{data}(\mathbf{x})} E_{\tilde{\mathbf{x}} \sim q_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x})} [\|\nabla_{\tilde{\mathbf{x}}} \log q_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x}) - \mathbf{s}_{\theta}(\tilde{\mathbf{x}})\|_2^2] + \text{const.} \\ &= E_{\mathbf{x} \sim p_{data}(\mathbf{x})} E_{\mathbf{z} \sim N(\mathbf{0}, I)} \left[\left\| \frac{1}{\sigma} \mathbf{z} + \mathbf{s}_{\theta}(\mathbf{x} + \sigma \mathbf{z}) \right\|_2^2 \right] + \text{const.} \end{aligned}$$

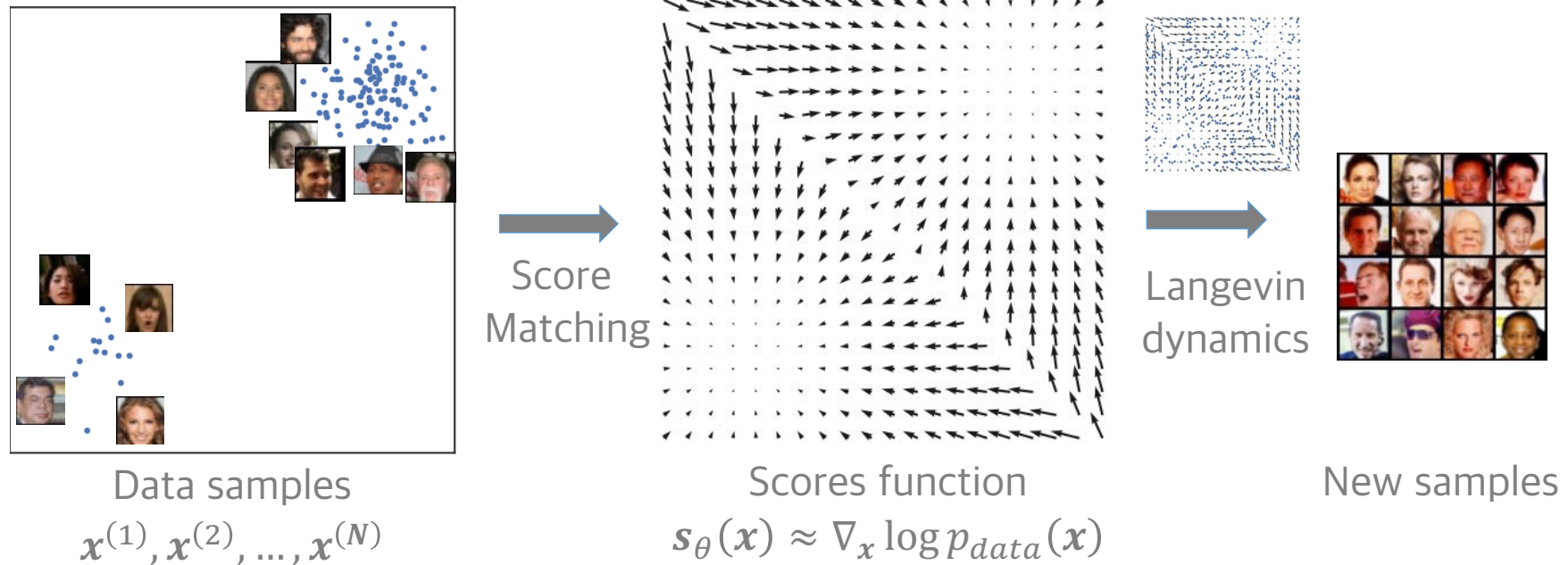
- **Pros**
 - more scalable than score matching
 - reduces score estimation to a denoising task
- **Con:** cannot estimate the score of clean data (noise-free)

$$\mathbf{s}_{\theta}(\mathbf{x}) \approx \nabla_{\mathbf{x}} \log q_{\sigma}(\mathbf{x}) \neq \nabla_{\mathbf{x}} \log p_{data}(\mathbf{x})$$

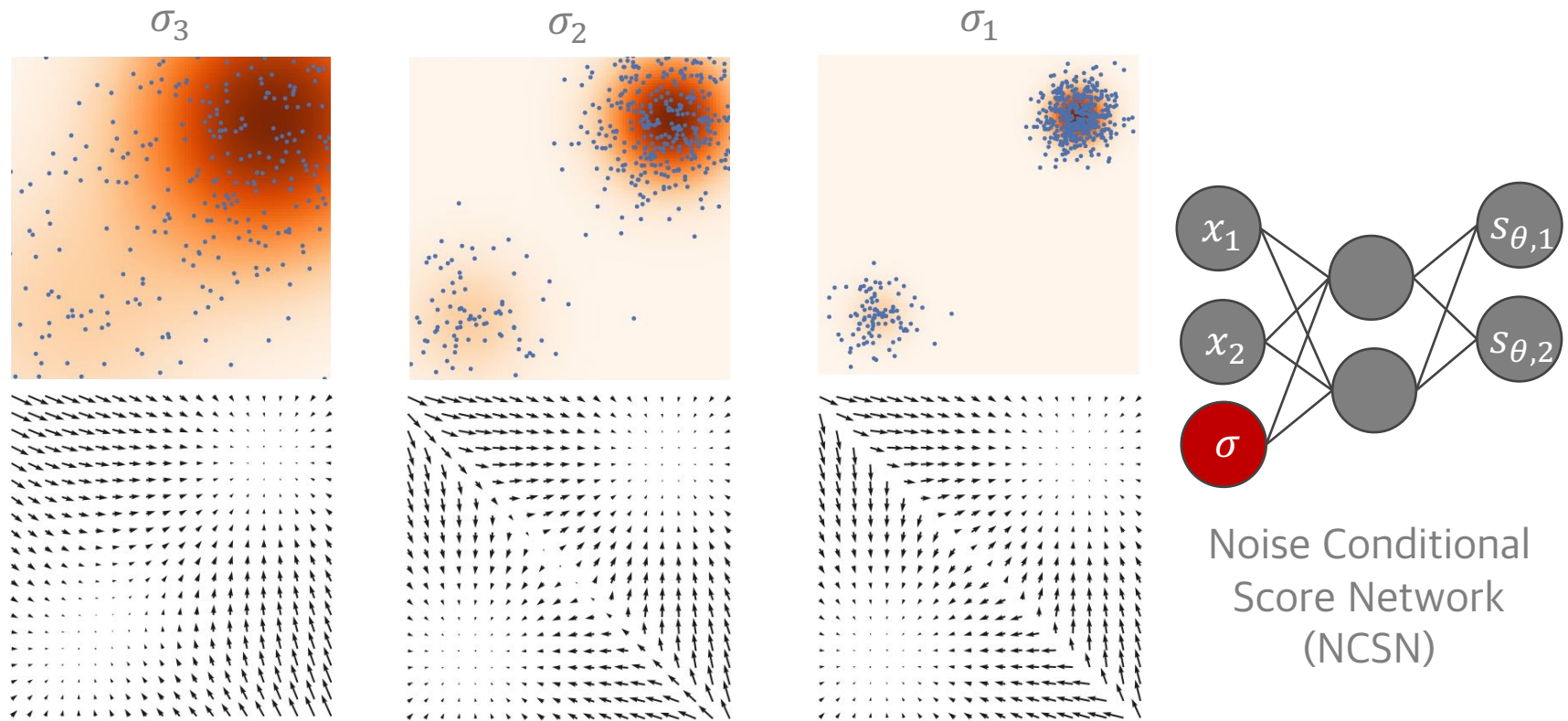
Denoising score matching with Langevin dynamics

- Let $q_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) := N(\tilde{\mathbf{x}}|\mathbf{x}, \sigma^2 I)$, $q_\sigma(\tilde{\mathbf{x}}) := \int p_{data}(\mathbf{x}) q_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) d\mathbf{x}$
- The objective
$$E_{\tilde{\mathbf{x}} \sim q_\sigma(\tilde{\mathbf{x}})} [\|\mathbf{s}_\theta(\tilde{\mathbf{x}}) - \nabla_{\tilde{\mathbf{x}}} \log q_\sigma(\tilde{\mathbf{x}})\|_2^2]$$
$$= E_{\mathbf{x} \sim p_{data}(\mathbf{x})} E_{\tilde{\mathbf{x}} \sim q_\sigma(\tilde{\mathbf{x}}|\mathbf{x})} [\|\mathbf{s}_\theta(\tilde{\mathbf{x}}) - \nabla_{\tilde{\mathbf{x}}} \log q_\sigma(\tilde{\mathbf{x}}|\mathbf{x})\|_2^2] + const.$$
- Consider a sequence of positive noise scales
$$\sigma_1 < \sigma_2 < \dots < \sigma_L$$
- σ_1 is small enough $q_{\sigma_1}(\mathbf{x}) \approx p_{data}(\mathbf{x})$
- σ_L is large enough $q_{\sigma_L}(\mathbf{x}) \approx N(\mathbf{x}|\mathbf{0}, \sigma_L^2 I)$

Score-based generative modeling



Joint score estimation via noise conditional score networks



Denoising score matching with Langevin dynamics

- For each $q_{\sigma_i}(\mathbf{x})$ with $\sigma_1 < \sigma_2 < \dots < \sigma_L$, Song & Ermon run T steps of Langevin MCMC to get a sample sequentially

$$\mathbf{x}_i^t := \mathbf{x}_i^{t-1} + \frac{\alpha_i}{2} \mathbf{s}_{\theta^*}(\mathbf{x}_i^{t-1}, \sigma_i) + \sqrt{\alpha_i} \mathbf{z}, \quad t = 1, 2, \dots, T$$

- where $\alpha_i > 0$ is the step size and $\mathbf{z} \sim N(\mathbf{0}, \mathbf{I})$

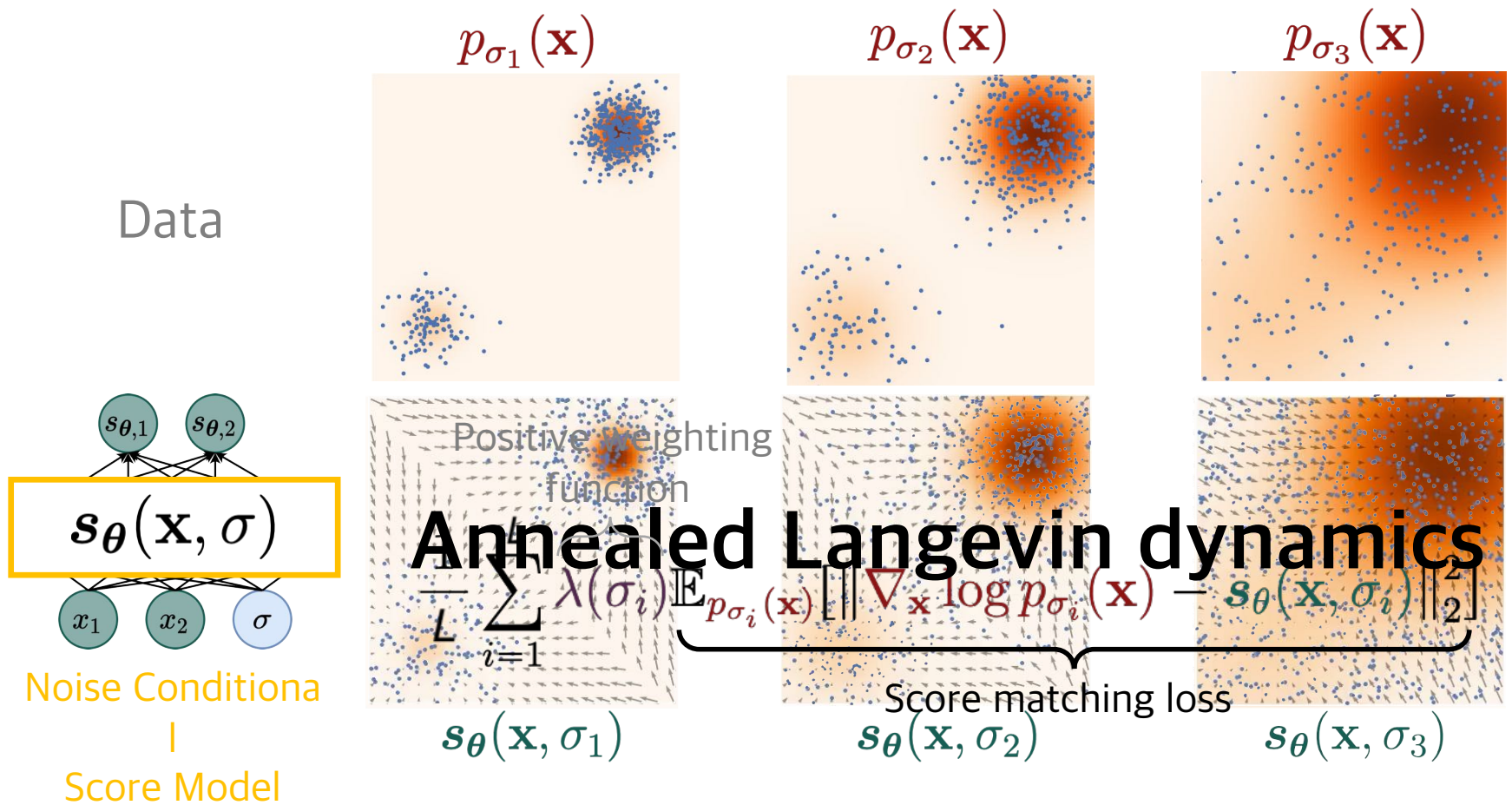
$$\alpha_i := \epsilon \frac{\sigma_i^2}{\sigma_1^2}$$

- $\epsilon > 0$

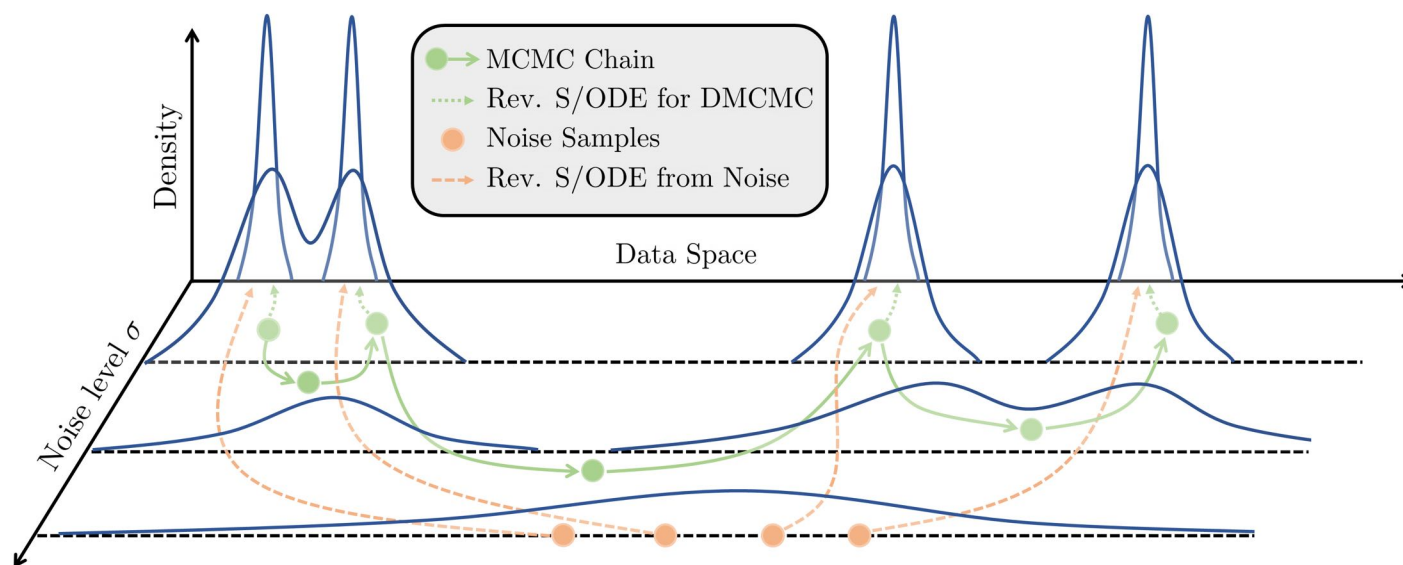
Generative Modeling by Estimating Gradients of the Data Distribution

Song Yang, and Stefano Ermon. NeurIPS 2019

Using multiple noise levels



Denoising score matching with Langevin dynamics



Conceptual illustration of a multiple noise score matching with Langevin sampling process

DENOISING MCMC FOR ACCELERATING DIFFUSION-BASED GENERATIVE MODELS

Beomsu Kim, Jong Chul Ye. ICML 2023

Denoising score matching with Langevin dynamics

- Let $q_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) := N(\tilde{\mathbf{x}}|\mathbf{x}, \sigma^2 I)$, $q_\sigma(\tilde{\mathbf{x}}) := \int p_{data}(\mathbf{x}) q_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) d\mathbf{x}$
- Consider a sequence of positive noise scales

$$\sigma_1 < \sigma_2 < \dots < \sigma_L$$

- σ_1 is small enough $q_{\sigma_1}(\mathbf{x}) \approx p_{data}(\mathbf{x})$
- σ_L is large enough $q_{\sigma_L}(\mathbf{x}) \approx N(\mathbf{x}|\mathbf{0}, \sigma_L^2 I)$

Data space

Noise space



Denoising score matching with Langevin dynamics

- Let $q_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) := N(\tilde{\mathbf{x}}|\mathbf{x}, \sigma^2 I)$, $q_\sigma(\tilde{\mathbf{x}}) := \int p_{data}(\mathbf{x}) q_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) d\mathbf{x}$
- Consider a sequence of positive noise scales

$$\sigma_1 < \sigma_2 < \dots < \sigma_L$$

- **Noise conditional score network**

$$\sum_{i=1}^L \sigma_i^2 E_{\mathbf{x} \sim p_{data}(\mathbf{x})} E_{\tilde{\mathbf{x}} \sim q_{\sigma_i}(\tilde{\mathbf{x}}|\mathbf{x})} \left[\left\| \mathbf{s}_\theta(\tilde{\mathbf{x}}, \sigma_i) - \nabla_{\tilde{\mathbf{x}}} \log q_{\sigma_i}(\tilde{\mathbf{x}}|\mathbf{x}) \right\|_2^2 \right]$$

- Given sufficient data and model capacity, the optimal score-based model

$$\mathbf{s}_{\theta^*}(\mathbf{x}, \sigma_i) \approx \nabla_{\mathbf{x}} \log q_{\sigma_i}(\mathbf{x}) \quad \text{for } \sigma \in \{\sigma_1, \dots, \sigma_L\}$$

- The weights σ_i^2 are related to $\sigma_i^2 \propto 1/E \left[\left\| \nabla_{\tilde{\mathbf{x}}} \log p_{\sigma_i}(\tilde{\mathbf{x}}|\mathbf{x}) \right\|_2^2 \right]$

Generation with annealed Langevin dynamics

- For each $q_{\sigma_i}(\mathbf{x})$ with $\sigma_1 < \sigma_2 < \dots < \sigma_L$, Song & Ermon run T steps of Langevin MCMC to get a sample sequentially

$$\mathbf{x}_i^t := \mathbf{x}_i^{t-1} + \frac{\alpha_i}{2} \mathbf{s}_{\theta^*}(\mathbf{x}_i^{t-1}, \sigma_i) + \sqrt{\alpha_i} \mathbf{z}, \quad t = 1, 2, \dots, T$$

- where $\alpha_i > 0$ is the step size and $\mathbf{z} \sim N(\mathbf{0}, \mathbf{I})$

$$\alpha_i := \epsilon \frac{\sigma_i^2}{\sigma_1^2}$$

- $\epsilon > 0$

Generative Modeling by Estimating Gradients of the Data Distribution

Song Yang, and Stefano Ermon. NeurIPS 2019

Denoising diffusion probabilistic models(DDPM)

- Consider a seq. of positive noise scales $0 < \beta_1 < \beta_2 \cdots < \beta_T < 1$
- $\mathbf{x}_0 \sim p_{data}(\mathbf{x})$, construct latent variables $\{\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T\}$ s.t.

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) := N(\mathbf{x}_t | \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

- i.e., $q(\mathbf{x}_t | \mathbf{x}_0) = N(\mathbf{x}_t | \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$ where $\alpha_t := 1 - \beta_t$,
 $\bar{\alpha}_t := \prod_{s=1}^t \alpha_s$
- Similar to SMLD, we can denote the perturbed data distribution

$$q(\mathbf{x}_t) := \int q(\mathbf{x}_t | \mathbf{x}) p_{data}(\mathbf{x}) d\mathbf{x}$$

- The noise scales are prescribed s.t. $\mathbf{x}_T \sim q(\mathbf{x}_T) \approx N(\mathbf{0}, \mathbf{I})$



Denoising diffusion probabilistic models(DDPM)

- A variational Markov chain in the reverse direction is parametrized with

$$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) = N(\mathbf{x}_{t-1}|\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \beta_t \mathbf{I})$$

- where $\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}}(\mathbf{x}_t + \beta_t \mathbf{s}_{\theta}(\mathbf{x}_t, t))$
- Re-weighted variant of the evidence lower bound

$$\sum_{t=1}^T (1 - \bar{\alpha}_t) E_{\mathbf{x} \sim p_{data}(\mathbf{x})} E_{\mathbf{x}_t \sim q(\mathbf{x}_t|\mathbf{x})} \left[\|\mathbf{s}_{\theta}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t|\mathbf{x})\|_2^2 \right]$$

- which is a weighted sum of denoising score matching

$$\mathbf{s}_{\theta^*}(\mathbf{x}_t, t) \approx \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t)$$

- The weights $(1 - \bar{\alpha}_t)$ are related to

$$(1 - \bar{\alpha}_t) \propto 1/E \left[\|\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t|\mathbf{x})\|_2^2 \right]$$

Denoising diffusion probabilistic models(DDPM)

- Generate samples by starting from $\mathbf{x}_T \sim N(\mathbf{0}, I)$
- $\mathbf{x}_{t-1} := \underbrace{\frac{1}{\sqrt{\alpha_t}}(\mathbf{x}_t + \beta_t \mathbf{s}_{\theta^*}(\mathbf{x}_t, t))}_{=\mu_{\theta^*}(\mathbf{x}_t, t)} + \sqrt{\beta_t} \mathbf{z}, \quad t = T, T-1, \dots, 2$
- We call this method **ancestral sampling** ($\prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$)

Denoising Diffusion Probabilistic Models

Jonathan Ho, Ajay Jain, Pieter Abbeel. NeurIPS 2020

Summary of score-based models

- **SMLD** and **DDPM** involve sequentially corrupting training data with slowly increasing noise, and then learning to reverse this corruption to form a generative model of the data
- **SMLD** estimates the score at each noise scale and then use Langevin dynamics to sample from a sequence of decreasing noise scales during generation
- **DDPM** trains a sequence of probabilistic models to reverse each step of the noise corruption, using knowledge of the functional form of the reverse distributions to make training tractable

Recap: Latent Variable Models

- Observable variables $\mathbf{x} \in \mathbb{R}^d$
- Latent variables $\mathbf{z} \in \mathbb{R}^h$ (unobservable)

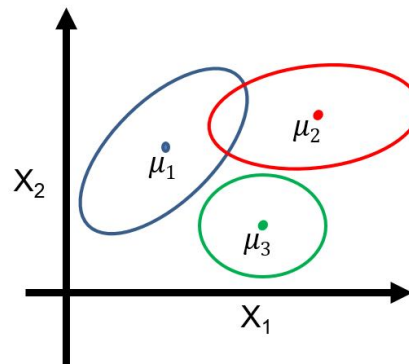
$$p_{data}(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z})$$

or

$$= \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z}$$

Recap: Mixture of Gaussians

- Mixture of Gaussians. Bayes net: $z \rightarrow \mathbf{x}$
 - $z = \text{Categorical}(z|\gamma_1, \dots, \gamma_K)$
 - $p(\mathbf{x}|z = k) = N(\mathbf{x}|\mu_k, \Sigma_k)$



- Generative Process
 - Pick a mixture component k by sampling z
 - Generate a data point by sampling from that Gaussian

Denoising diffusion probabilistic models(DDPM)

- DDPM is a latent variable model

$$p_{\theta}(\mathbf{x}_0) := \int p_{\theta}(\mathbf{x}_0, \mathbf{x}_1 \dots, \mathbf{x}_T) d\mathbf{x}_{1:T}$$

- $\mathbf{x}_0 = q(\mathbf{x}_0) = p_{data}$
- The joint distribution $p_{\theta}(\mathbf{x}_{0:T})$ is called the **reverse process** starting at $p_{\theta}(\mathbf{x}_T) = N(\mathbf{x}_T | \mathbf{0}, I)$

$$p_{\theta}(\mathbf{x}_{0:T}) = p_{\theta}(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t),$$
$$p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t) = N(\mathbf{x}_{t-1} | \boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \Sigma_{\theta}(\mathbf{x}_t, t))$$

Denoising diffusion probabilistic models(DDPM)

- **Forward process** or diffusion process $q(\mathbf{x}_{1:T}|\mathbf{x}_0)$ is fixed to a Markov chain that gradually adds Gaussian noise to the data according to a variance schedule $0 < \beta_1 < \dots < \beta_T < 1$

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

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) := N(\mathbf{x}_t | \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

- β_t can be learned by reparameterization or held constants as hyperparameters
- Let $\alpha_t := 1 - \beta_t$ and $\bar{\alpha}_t := \prod_{s=1}^t \alpha_s$, then

$$q(\mathbf{x}_t|\mathbf{x}_0) = N(\mathbf{x}_t | \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$$

Denoising diffusion probabilistic models(DDPM)

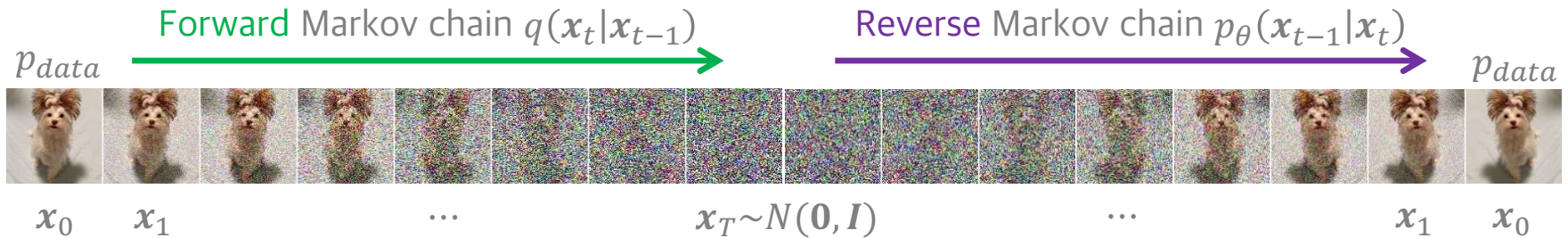


Image to noise(prescribed)

Noise to image(learnable)

$$q(x_t|x_{t-1}) = N(x_t|\sqrt{\alpha_t}x_{t-1}, \beta_t I)$$

$$q(x_t|x_0) = N(x_t|\sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t) I)$$

$$p_\theta(x_{t-1}|x_t) = N(x_{t-1}|\mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$$

where $\Sigma_\theta(x_t, t) = \sigma_t^2 I = \beta_t I$

- What is **target** of $p_\theta(x_{t-1}|x_t) = N(x_{t-1}|\mu_\theta(x_t, t), \beta_t I)$?
 - $p_\theta(x_{t-1}|x_t) \approx q(x_{t-1}|x_t)$?



Denoising diffusion probabilistic models(DDPM)

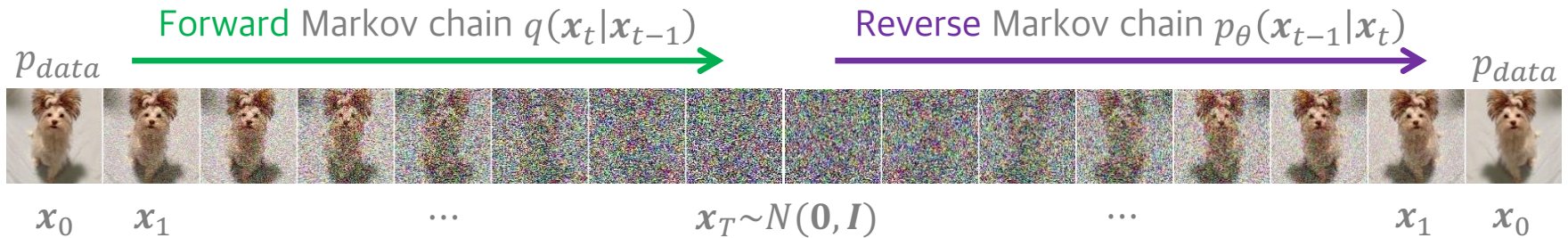


Image to noise(prescribed)

Noise to image(learnable)

$$q(x_t|x_{t-1}) = N(x_t|\sqrt{\alpha_t}x_{t-1}, \beta_t I)$$

$$q(x_t|x_0) = N(x_t|\sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t) I)$$

$$p_\theta(x_{t-1}|x_t) = N(x_{t-1}|\mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$$

where $\Sigma_\theta(x_t, t) = \sigma_t^2 I = \beta_t I$

- What is **target** of $p_\theta(x_{t-1}|x_t) = N(x_{t-1}|\mu_\theta(x_t, t), \beta_t I)$?
 - $p_\theta(x_{t-1}|x_t) \approx q(x_{t-1}|x_t)$?
 - $q(x_{t-1}|x_t)$ is not tractable

$$q(x_{t-1}|x_t) = \frac{q(x_t|x_{t-1})q(x_{t-1})}{q(x_t)}, q(x_t) = \int q(x_t|x_0)q(x_0)dx_0$$

Denoising diffusion probabilistic models(DDPM)

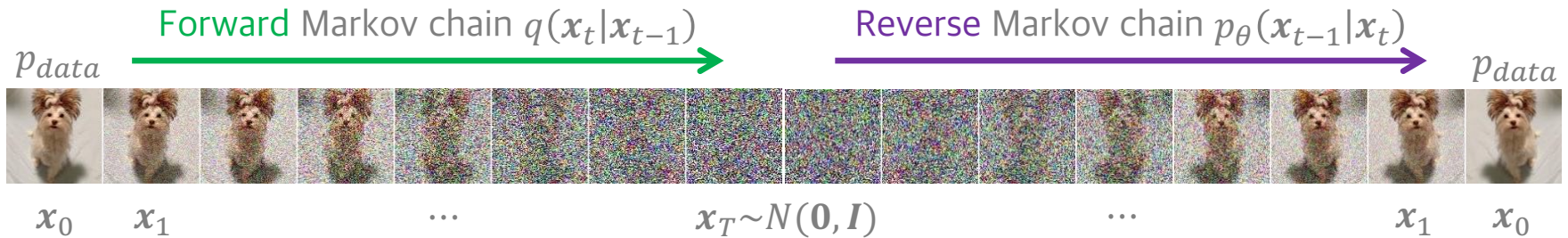


Image to noise(prescribed)

Noise to image(learnable)

$$q(x_t|x_{t-1}) = N(x_t|\sqrt{\alpha_t}x_{t-1}, \beta_t I)$$

$$q(x_t|x_0) = N(x_t|\sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t) I)$$

$$p_\theta(x_{t-1}|x_t) = N(x_{t-1}|\mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$$

where $\Sigma_\theta(x_t, t) = \sigma_t^2 I = \beta_t I$

- What is **target** of $p_\theta(x_{t-1}|x_t) = N(x_{t-1}|\mu_\theta(x_t, t), \beta_t I)$?
 - $q(x_{t-1}|x_t, x_0)$ is tractable! Why?



Denoising diffusion probabilistic models(DDPM)

Image to noise(prescribed)

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = N(\mathbf{x}_t|\sqrt{\alpha_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$$
$$q(\mathbf{x}_t|\mathbf{x}_0) = N(\mathbf{x}_t|\sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})$$

Noise to image(learnable)

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = N(\mathbf{x}_{t-1}|\boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t))$$

where $\Sigma_\theta(\mathbf{x}_t, t) = \sigma_t^2\mathbf{I} = \beta_t\mathbf{I}$

- What is **target** of $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = N(\mathbf{x}_{t-1}|\boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \beta_t\mathbf{I})$?
 - $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ is tractable. Why?

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)}$$

$$= N(\mathbf{x}_{t-1}|\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, t), \tilde{\beta}_t\mathbf{I})$$

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

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

Denoising diffusion probabilistic models(DDPM)

Image to noise(prescribed)

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = N(\mathbf{x}_t|\sqrt{\alpha_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$$
$$q(\mathbf{x}_t|\mathbf{x}_0) = N(\mathbf{x}_t|\sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})$$

Noise to image(learnable)

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = N(\mathbf{x}_{t-1}|\boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t))$$

where $\Sigma_\theta(\mathbf{x}_t, t) = \sigma_t^2\mathbf{I} = \beta_t\mathbf{I}$

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = N(\mathbf{x}_{t-1}|\boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \beta_t\mathbf{I}) \approx q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$$
$$= N(\mathbf{x}_{t-1}|\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, t), \tilde{\beta}_t\mathbf{I})$$

- I.e.,

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

- If $\boldsymbol{\mu}_\theta(\mathbf{x}_t, t) := \frac{1}{\sqrt{\alpha_t}}(\mathbf{x}_t + \beta_t\mathbf{s}_\theta(\mathbf{x}_t, t))$, then

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

Foundation of DDPM

$$\begin{aligned} & \operatorname{argmin}_{\theta} D(q(\mathbf{x}_{t-1}|\mathbf{x}_t) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)) \\ &= \operatorname{argmin}_{\theta} E_{\mathbf{x}_0 \sim p_{data}} [D(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))] \end{aligned}$$

- Proof: By Fubini theorem,

$$\begin{aligned} & - \int q(\mathbf{z}) \log p_{\theta}(\mathbf{z}) d\mathbf{z} = - \int \left[\int q(\mathbf{z}|\mathbf{x}_0) p_{\theta}(\mathbf{z}) d\mathbf{z} \right] p_{data}(\mathbf{x}_0) d\mathbf{x}_0 \\ &= D(q(\mathbf{z}) \parallel p_{\theta}(\mathbf{z})) + \text{const. w.r.t. } \theta \\ &= E_{\mathbf{x}_0 \sim p_{data}} [D(q(\mathbf{z}|\mathbf{x}_0) \parallel p_{\theta}(\mathbf{z}))] + \text{const. w.r.t. } \theta \end{aligned}$$

Learning objective of DDPM

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

Learning objective of DDPM

- Optimizing the variational bound on log-likelihood

$$\begin{aligned}\log p_{\theta}(\mathbf{x}_0) &= \log \int p_{\theta}(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T} = \log \int p_{\theta}(\mathbf{x}_{0:T}) \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} d\mathbf{x}_{1:T} \\ &\geq \int q(\mathbf{x}_{1:T}|\mathbf{x}_0) \log \frac{p_{\theta}(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} d\mathbf{x}_{1:T}\end{aligned}$$

Learning objective of DDPM

- Optimizing the variational bound on log-likelihood

$$\log p_{\theta}(\mathbf{x}_0) = \log \int p_{\theta}(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T} = \log \int p_{\theta}(\mathbf{x}_{0:T}) \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} d\mathbf{x}_{1:T}$$

$$\geq \int q(\mathbf{x}_{1:T}|\mathbf{x}_0) \log \frac{p_{\theta}(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} d\mathbf{x}_{1:T}$$

- The Markov property of $q(\mathbf{x}_t|\mathbf{x}_{t-1})$ and $p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$ implies

$$p_{\theta}(\mathbf{x}_{0:T}) = p_{\theta}(\mathbf{x}_T) \prod_{t=2}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t),$$

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) = q(\mathbf{x}_T|\mathbf{x}_0) \prod_{t=2}^T q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$$

Learning objective of DDPM

$$\begin{aligned} & \log \frac{p_{\theta}(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \\ &= \log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1) + \sum_{t=2}^T \log \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} + \log \frac{p_{\theta}(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_0)} \end{aligned}$$

$$\begin{aligned} & \int q(\mathbf{x}_{1:T}|\mathbf{x}_0) \log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1) d\mathbf{x}_{1:T} = \int q(\mathbf{x}_1|\mathbf{x}_0) \log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1) d\mathbf{x}_1 \\ &= E_{q(\mathbf{x}_1|\mathbf{x}_0)}[\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)] \\ & \int q(\mathbf{x}_{1:T}|\mathbf{x}_0) \log \frac{p_{\theta}(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_0)} d\mathbf{x}_{1:T} = \int q(\mathbf{x}_T|\mathbf{x}_0) \log \frac{p_{\theta}(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_0)} d\mathbf{x}_T \\ &= -D[q(\mathbf{x}_T|\mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_T)] \end{aligned}$$

Learning objective of DDPM

$$\begin{aligned} & \int q(\mathbf{x}_{1:T}|\mathbf{x}_0) \log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} d\mathbf{x}_{1:T} \\ &= \int \int q(\mathbf{x}_{t-1}, \mathbf{x}_t|\mathbf{x}_0) \log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} d\mathbf{x}_{t-1} d\mathbf{x}_t \\ &= \int \int \frac{q(\mathbf{x}_{t-1}, \mathbf{x}_t, \mathbf{x}_0)}{q(\mathbf{x}_0)} \frac{q(\mathbf{x}_t, \mathbf{x}_0)}{q(\mathbf{x}_t, \mathbf{x}_0)} \log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} d\mathbf{x}_{t-1} d\mathbf{x}_t \\ &= - \int q(\mathbf{x}_t|\mathbf{x}_0) D[q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)] d\mathbf{x}_t \\ &= -E_{q(\mathbf{x}_t|\mathbf{x}_0)} [D[q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)]] \end{aligned}$$

Learning objective of DDPM

$$\begin{aligned} & E_{\mathbf{x}_0 \sim q(\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0)] \\ & \geq E_{q(\mathbf{x}_0)} \left[E_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)] \right] - E_{q(\mathbf{x}_0)} [D[q(\mathbf{x}_T|\mathbf{x}_0) \parallel p_\theta(\mathbf{x}_T)]] \\ & - \sum_{t=2}^T E_{q(\mathbf{x}_0)} \left[E_{q(\mathbf{x}_t|\mathbf{x}_0)} [D[q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)]] \right] \end{aligned}$$

Learning objective of DDPM

- $p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) = N(\mathbf{x}_{t-1}|\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \beta_t \mathbf{I})$
- $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = N(\mathbf{x}_{t-1}|\tilde{\boldsymbol{\mu}}_t, \tilde{\beta}_t \mathbf{I})$
 - $\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, t) = \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1-\bar{\alpha}_t} \mathbf{x}_0 + \frac{\sqrt{\bar{\alpha}_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} \mathbf{x}_t$
 - $\tilde{\beta}_t = \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t} \beta_t$
- $q(\mathbf{x}_t|\mathbf{x}_0) = N(\mathbf{x}_t|\sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}), \quad \mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon},$
 $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{I})$

Learning objective of DDPM

$$\begin{aligned} L_{t-1}(\theta) &= E_{\mathbf{x}_t \sim q(\mathbf{x}_t)} \left[D(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)) \right] \\ &= E_{\mathbf{x}_0 \sim p_{data}} E_{\mathbf{x}_t \sim q(\mathbf{x}_t | \mathbf{x}_0)} \left[D(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)) \right] \\ &= E_{\mathbf{x}_0 \sim p_{data}} E_{\mathbf{x}_t \sim q(\mathbf{x}_t | \mathbf{x}_0)} \|\boldsymbol{\mu}_\theta(\mathbf{x}_t, t) - \tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, t)\|_2^2 \\ &= E_{\mathbf{x}_0 \sim p_{data}} E_{\mathbf{x}_t \sim q(\mathbf{x}_t | \mathbf{x}_0)} \left\| \boldsymbol{\mu}_\theta(\mathbf{x}_t, t) - \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} \mathbf{x}_0 \right. \\ &\quad \left. - \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}_t \right\|_2^2 \\ &= \frac{\beta_t^2}{\alpha_t(1 - \bar{\alpha}_t)} E_{\mathbf{x}_0 \sim p_{data}} E_{\boldsymbol{\epsilon} \sim N(\mathbf{0}, I)} \|\boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t) - \boldsymbol{\epsilon}\|_2^2 \\ \bullet \text{ if } \boldsymbol{\mu}_\theta(\mathbf{x}_t, t) &= \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t + \frac{\beta_t}{\sqrt{1 - \alpha_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) \end{aligned}$$

Training and noise predictor

$$\begin{aligned} & \text{(mean predictor)} \quad \mu_{\theta}(\mathbf{x}_t, t) \\ &= \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t + \frac{\beta_t}{\sqrt{1 - \alpha_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) \quad \text{(noise predictor)} \end{aligned}$$

- i.e., $\frac{1}{\sqrt{1 - \alpha_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \approx \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t | \mathbf{x}_0)$

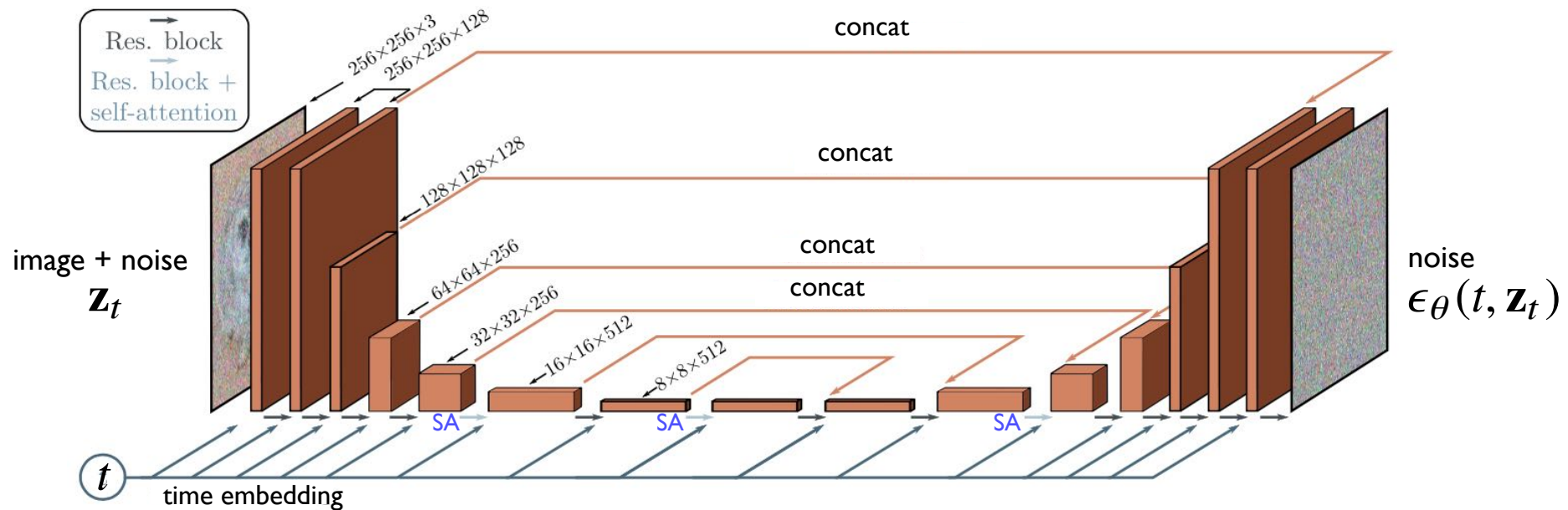
Algorithm 1 Training

- 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
 $\nabla_{\theta} \left\| \epsilon - \epsilon_{\theta}(\sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t} \epsilon, t) \right\|^2$
 - 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: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
 - 5: **end for**
 - 6: **return** \mathbf{x}_0
-

Architecture: U-net + self-attention + time Embedding



Experiments

- $T = 1000$, linear variance schedule $\beta_1 = 10^{-4}$ to $\beta_T = 0.02$
- U-Net backbone similar to an unmasked PixelCNN++ with group normalization



Figure 3: LSUN Church samples. FID=7.89

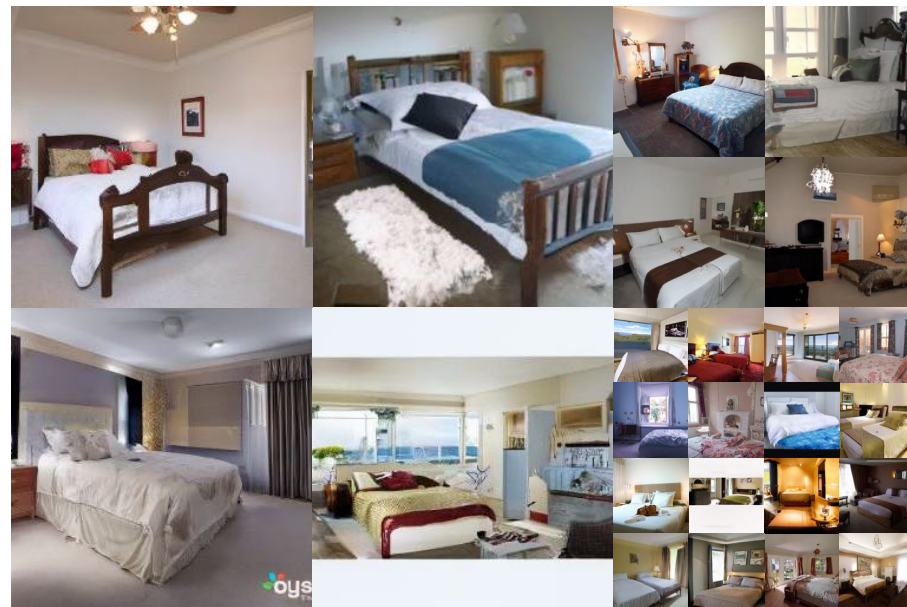


Figure 4: LSUN Bedroom samples. FID=4.90

Denoising Diffusion Probabilistic Models

Jonathan Ho, Ajay Jain, Pieter Abbeel. NeurIPS 2020

Thanks
