

Denoising Diffusion Probabilistic Models

김현수



Denoising Diffusion Models

Emerging as powerful generative models. outperforming GANs

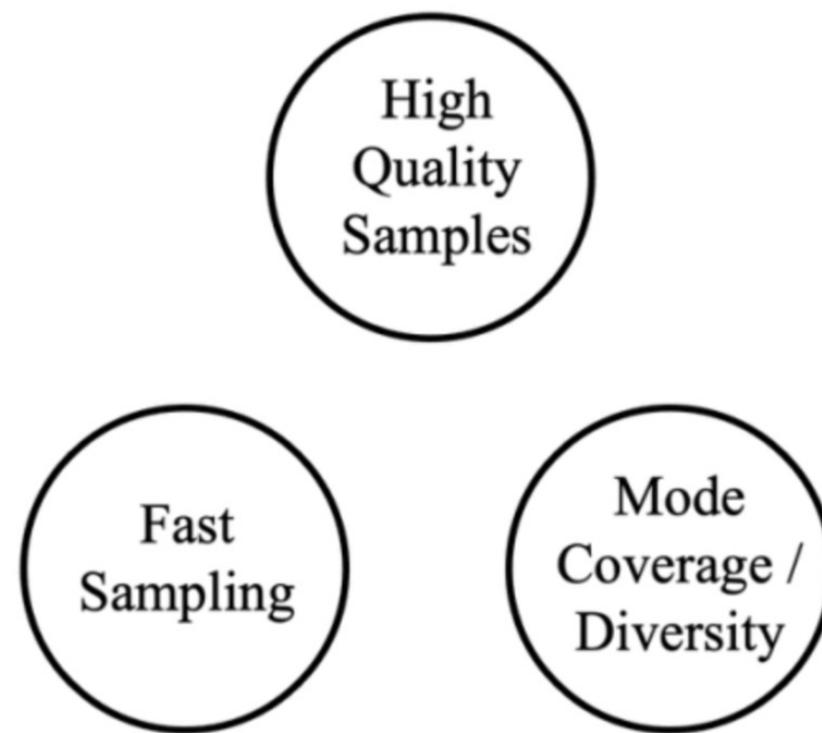


Diffusion Models Beat GANs on Image Synthesis, 2021

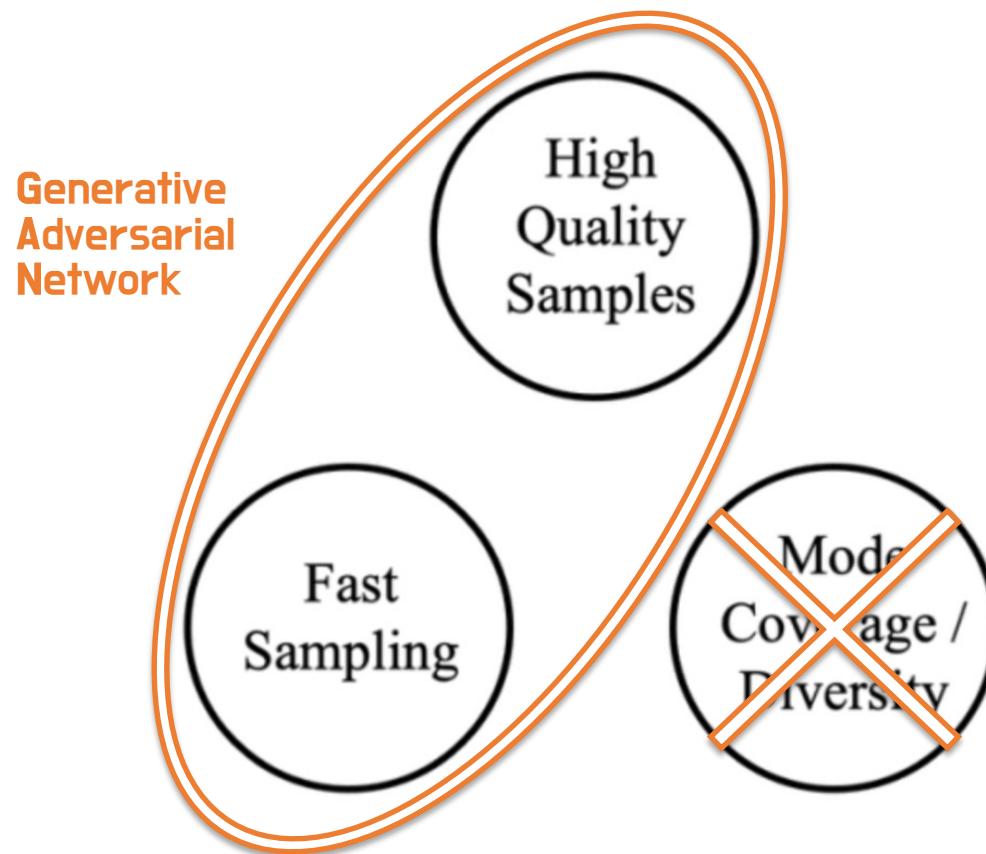


Cascaded Diffusion Models for High Fidelity Image Generation, 2021

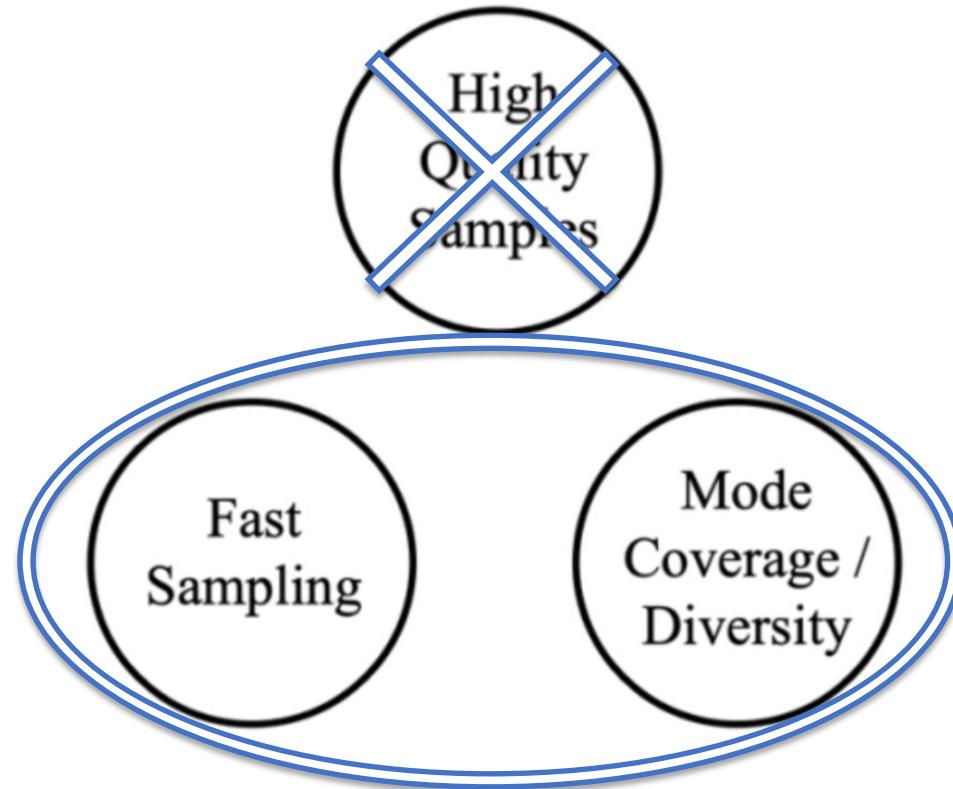
The Generative Learning Trilemma



The Generative Learning Trilemma

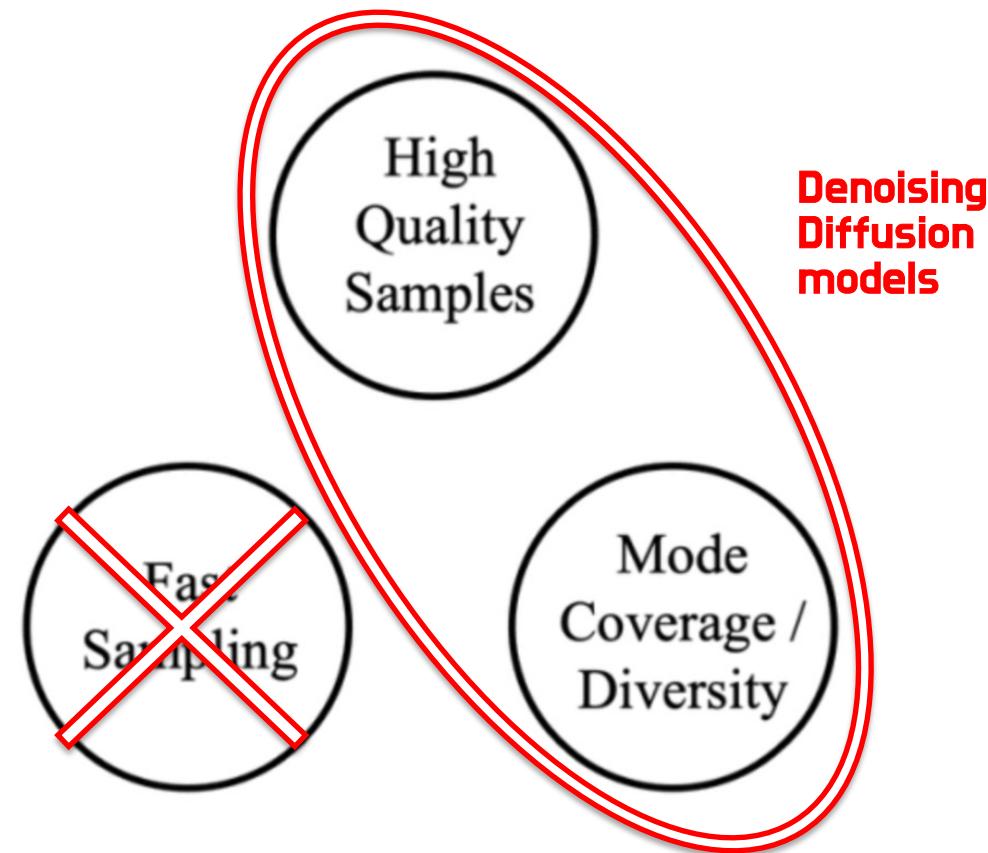


The Generative Learning Trilemma

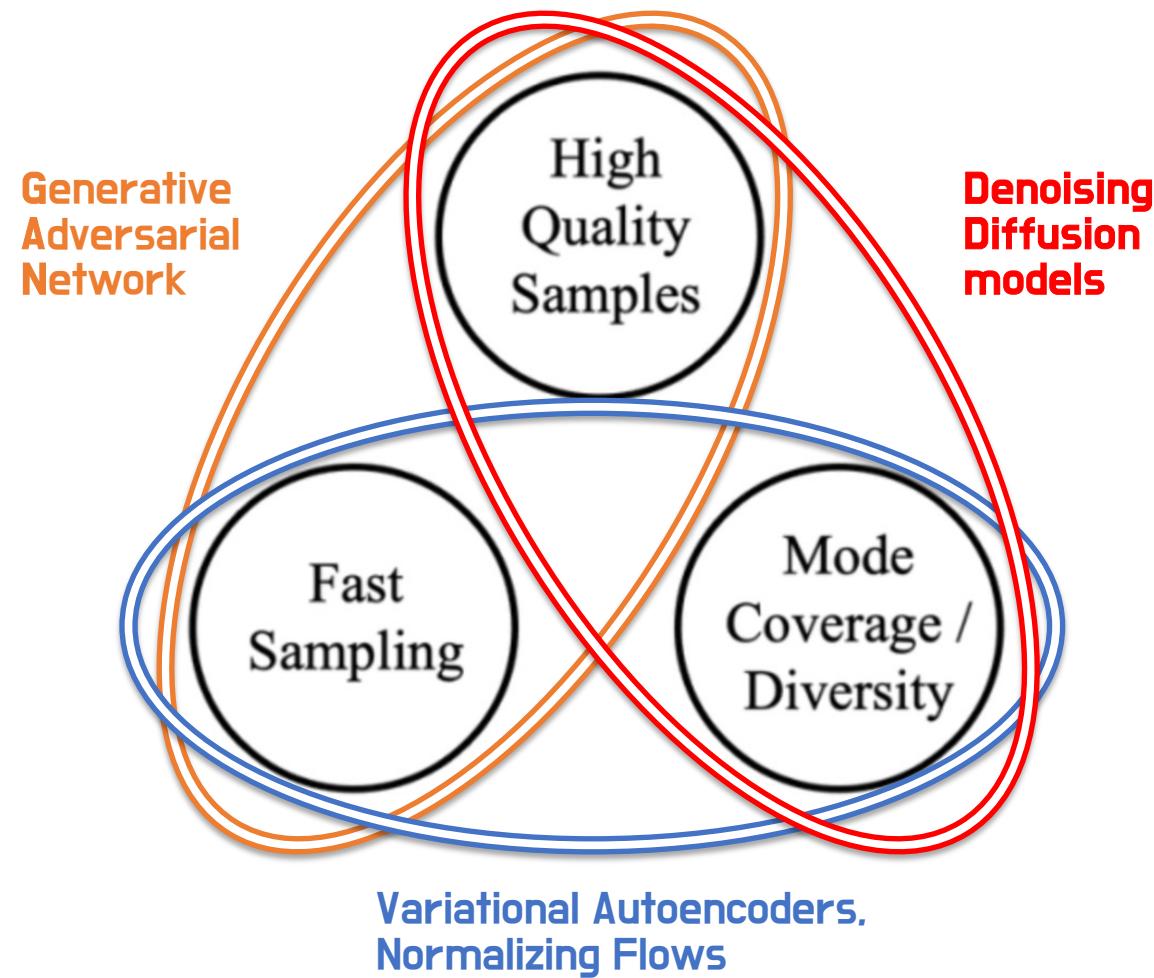


Variational Autoencoders.
Normalizing Flows

The Generative Learning Trilemma



The Generative Learning Trilemma



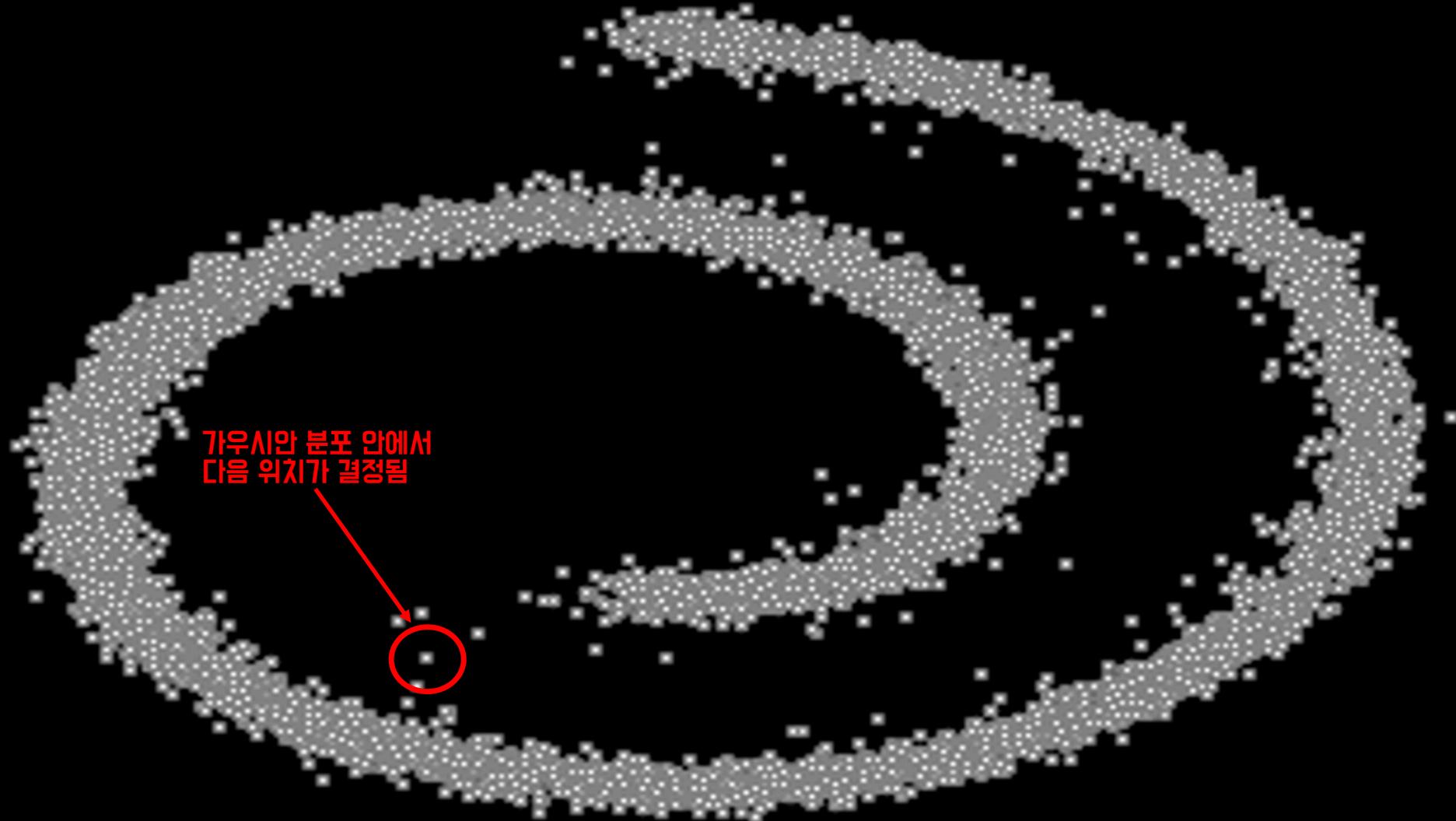
Denoising Diffusion Probabilistic Models

확산을 제거하는 확률적 모델

Denoising Diffusion Probabilistic Models

확산을 제거하는 확률적 모델



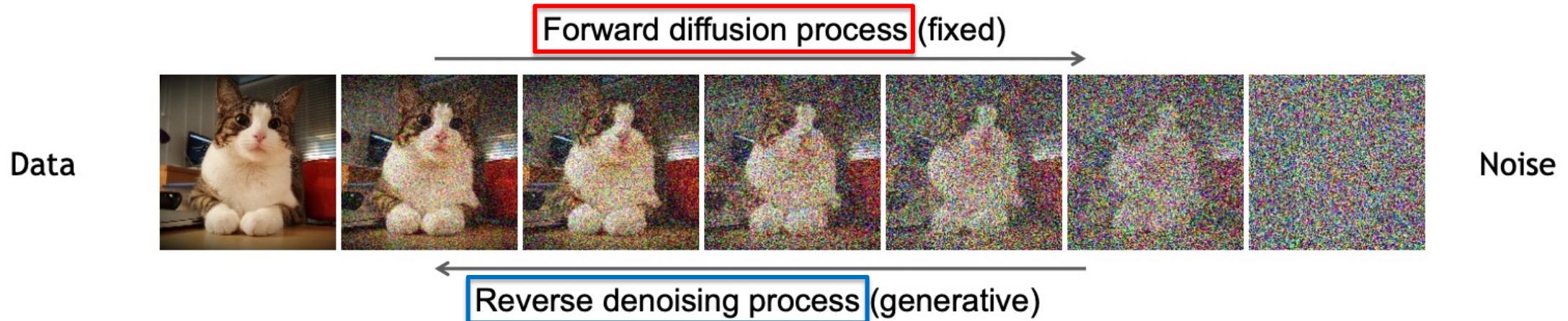


Denoising Diffusion Models

Learning to generate by denoising

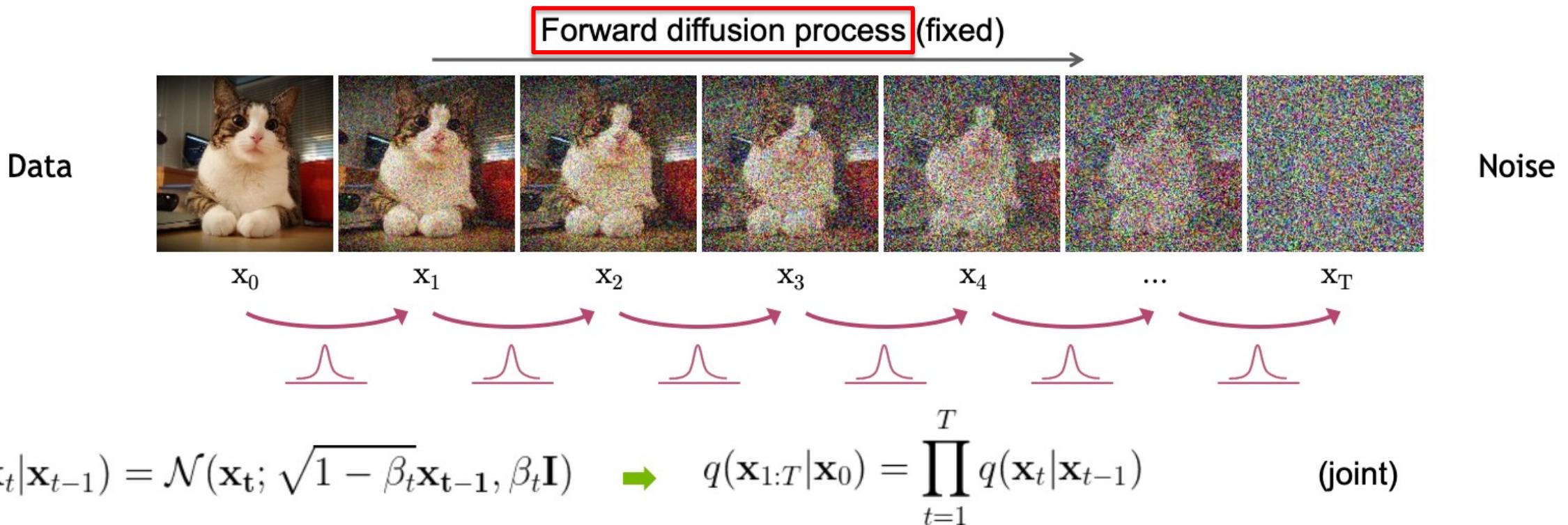
Denoising diffusion models consist of two processes:

- Forward diffusion process that gradually adds noise to input → 수학
- Reverse denoising process that learns to generate data by denoising → 학습(딥러닝)



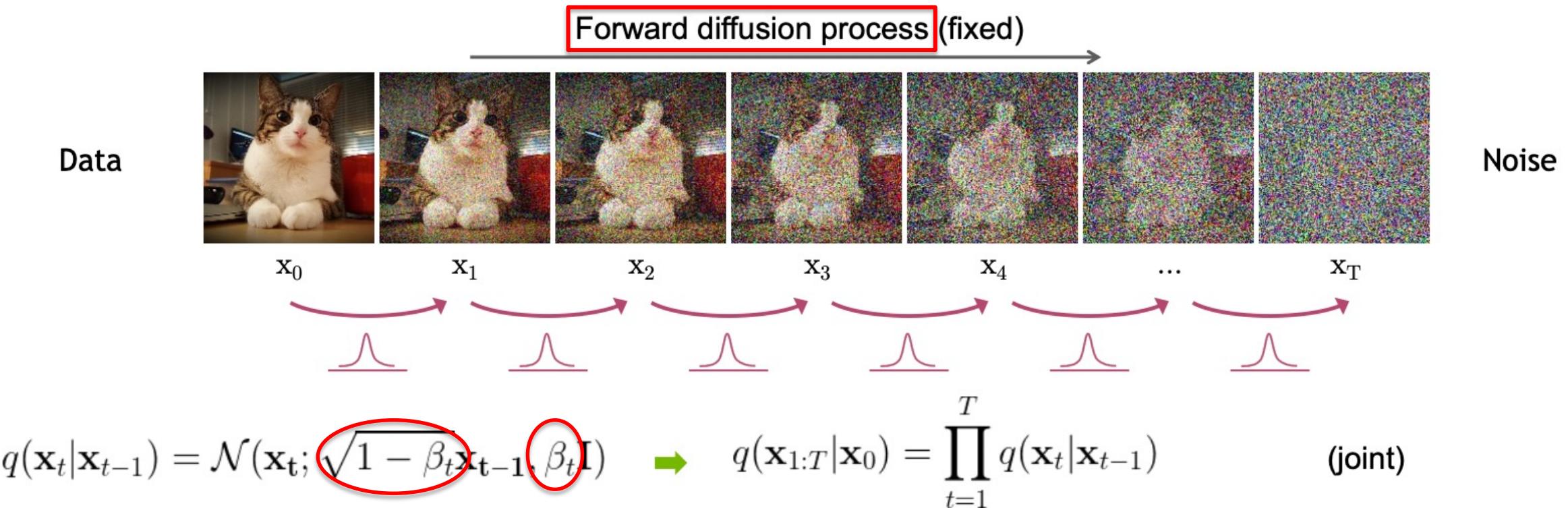
Forward (Diffusion) Process

The formal definition of the forward process in T steps:



Forward (Diffusion) Process

The formal definition of the forward process in T steps:



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

β_t 고정시키고
 ε 만 정규분포에서 sampling

Why?

stochastic node는 stochastic + deterministic한 부분으로
분해시켜 deterministic한 부분으로 backpropagation을 편하게 하자.

$$= \sqrt{1-\beta_t} x_{t-1} + \sqrt{\beta_t} \varepsilon_{t-1} \quad (\text{by reparametrization trick}) \\ (\varepsilon \sim N(0, I))$$

$$\text{Var}(q(x_t | x_{t-1})) = \text{Var}(\sqrt{1-\beta_t} x_{t-1} + \sqrt{\beta_t} \varepsilon_{t-1})$$

$$= (1-\beta_t) \text{Var}(x_{t-1}) + \beta_t \text{Var}(\varepsilon_{t-1})$$

$$\text{Var}(ax+by)$$

$$= a^2 \text{Var}(x) + b^2 \text{Var}(y) + 2ab \text{Cov}(x, y)$$

x, y 서로 독립이므로
 $\text{Cov}(x, y) = 0$

$$= (1-\beta_t) + \beta_t = \underline{1}$$

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t I) \quad \rightarrow \quad q(\mathbf{x}_{1:T} | \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1}) \quad (\text{joint})$$

* Markov Property

$$q(x_{t+1} | x_t) = q(x_{t+1} | x_1, x_2, \dots, x_t)$$

The formal definition

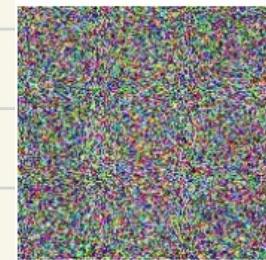
$$\prod_{t=1}^T q(x_t | x_{t-1}) = \frac{q(x_1, x_0)}{q(x_0)} \cdot \frac{q(x_2, x_1)}{q(x_1)} \cdot \dots \cdot \frac{q(x_T, x_{T-1})}{q(x_{T-1})}$$

$$= \frac{q(x_1, x_0)}{q(x_0)} \cdot \frac{q(x_2, x_1, x_0)}{q(x_1, x_0)} \cdot \dots \cdot \frac{q(x_T, x_{T-1}, \dots, x_1, x_0)}{q(x_{T-1}, \dots, x_0)}$$

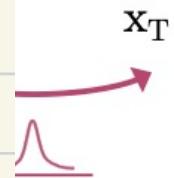
Data

$$= \frac{q(x_0, x_1, \dots, x_{T-1}, x_T)}{q(x_0)}$$

$$= q(x_{1:T} | x_0)$$



Noise

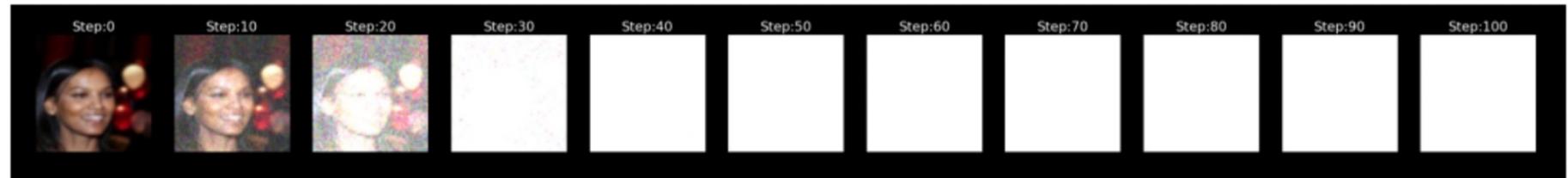
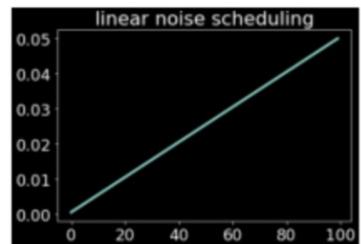


$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \quad \rightarrow \quad q(\mathbf{x}_{1:T} | \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1}) \quad (\text{joint})$$

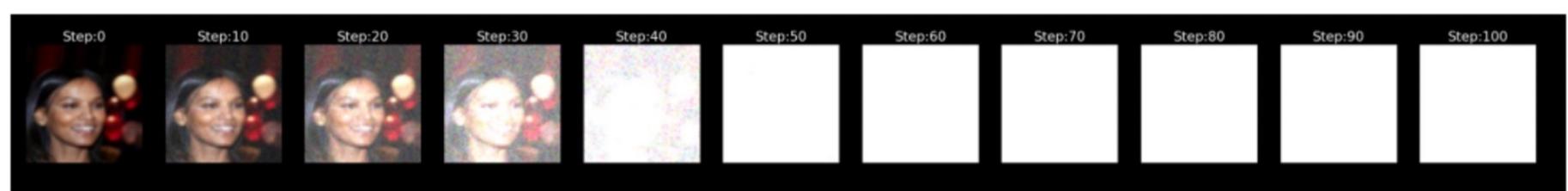
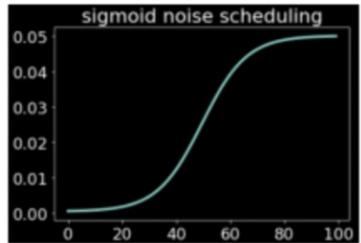
$$q(x_{1:T} | x_0) = \frac{q(x_0, x_1, \dots, x_T)}{q(x_0)}$$

$$\prod_{t=1}^T q(x_t | x_{t-1}) = \frac{q(x_1, x_0)}{q(x_0)} \cdot \frac{q(x_2, x_1)}{q(x_1)} \cdot \dots \cdot \frac{q(x_T, x_{T-1})}{q(x_{T-1})}$$

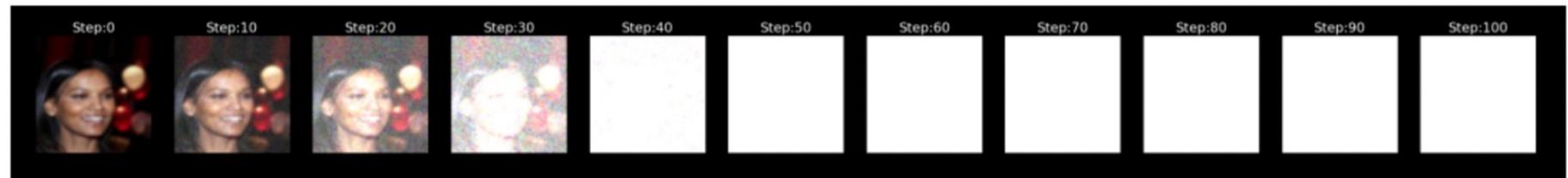
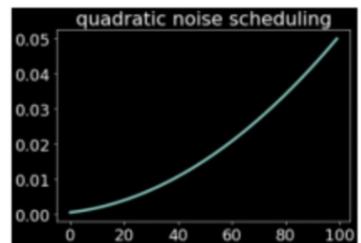
✓ Linear scheduling



✓ Sigmoid scheduling



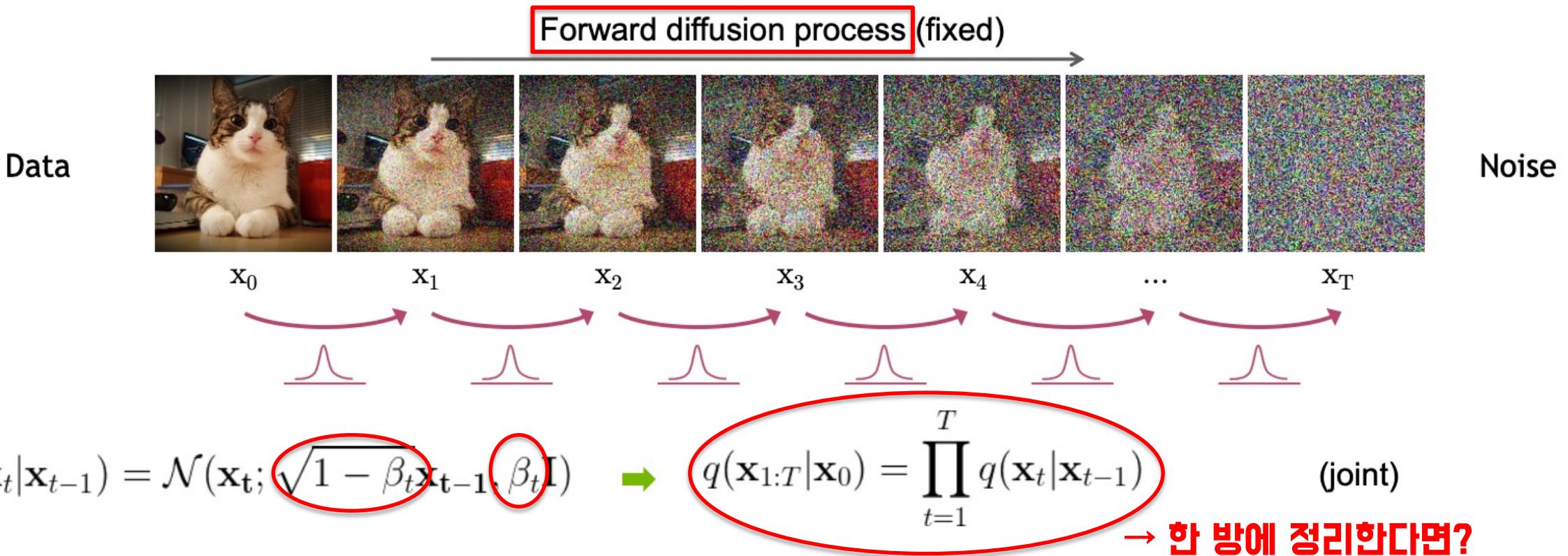
✓ Quadratic scheduling



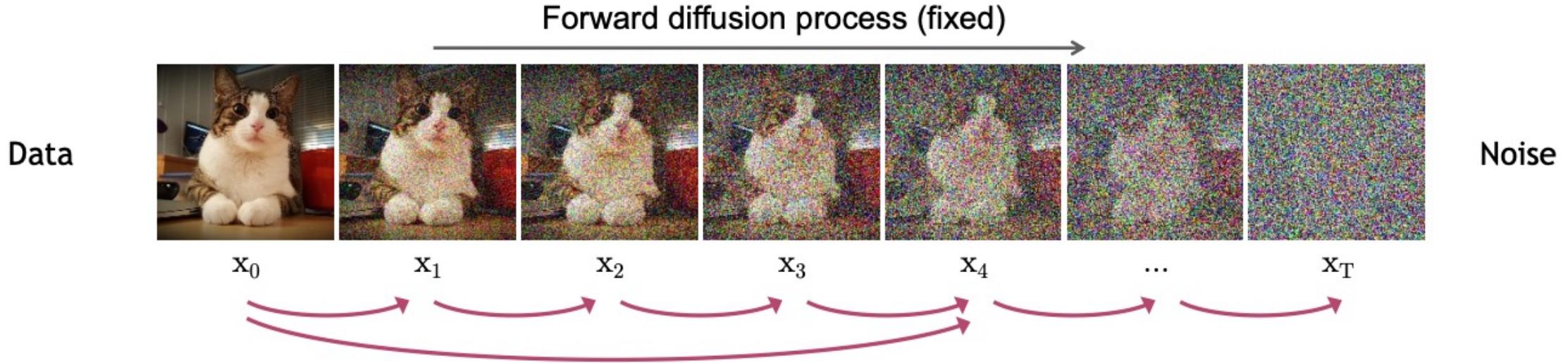
$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \quad \rightarrow \quad q(\mathbf{x}_{1:T} | \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1}) \quad (\text{joint})$$

Forward (Diffusion) Process

The formal definition of the forward process in T steps:



Forward (Diffusion) Process



Define $\bar{\alpha}_t = \prod_{s=1}^t (1 - \beta_s)$  $q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$ **(Diffusion Kernel)**

For sampling: $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{(1 - \bar{\alpha}_t)} \boldsymbol{\epsilon}$ where $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

β_t values schedule (i.e., the noise schedule) is designed such that $\bar{\alpha}_T \rightarrow 0$ and $q(\mathbf{x}_T | \mathbf{x}_0) \approx \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$

Let, $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$

$$X_t = \sqrt{1-\beta_t} X_{t-1} + \sqrt{\beta_t} \varepsilon_{t-1}$$

$$= \sqrt{\alpha_t} \underbrace{X_{t-1}}_{\text{blue}} + \sqrt{1-\alpha_t} \varepsilon_{t-1}$$

$$= \sqrt{\alpha_t} (\underbrace{\sqrt{\alpha_{t-1}} X_{t-2} + \sqrt{1-\alpha_{t-1}} \varepsilon_{t-2}}_{\text{blue}}) + \sqrt{1-\alpha_t} \varepsilon_{t-1}$$

$$= \sqrt{\alpha_t \alpha_{t-1}} X_{t-2} + \underbrace{\sqrt{1-\alpha_t} \cdot \sqrt{\alpha_t} \cdot \varepsilon_{t-2} + \sqrt{1-\alpha_t} \varepsilon_{t-1}}_{\text{red}}$$

$$= \sqrt{\alpha_t \alpha_{t-1}} X_{t-2} + \underbrace{\sqrt{1-\alpha_t \alpha_{t-1}} \varepsilon}_{\text{red}}$$

...

$$= \sqrt{\bar{\alpha}_t} X_0 + \sqrt{1-\bar{\alpha}_t} \varepsilon$$

* $\varepsilon_{t-1}, \varepsilon_{t-2}, \dots \sim N(0, I)$

* Joint Distribution (two Gaussians)

$$A \sim N(0, \sigma_1^2 I), B \sim N(0, \sigma_2^2 I)$$

$$\longrightarrow A+B \sim N(0, (\sigma_1^2 + \sigma_2^2) I)$$

$$\therefore f(X_t | X_0) = N(X_t; \sqrt{\bar{\alpha}_t} X_0, (1-\bar{\alpha}_t)I)$$

Define $\bar{\alpha}_t = \prod_{s=1}^t (1 - \beta_s) \quad \rightarrow \quad q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t)I) \quad (\text{Diffusion Kernel})$

For sampling: $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{(1 - \bar{\alpha}_t)} \epsilon \quad \text{where } \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

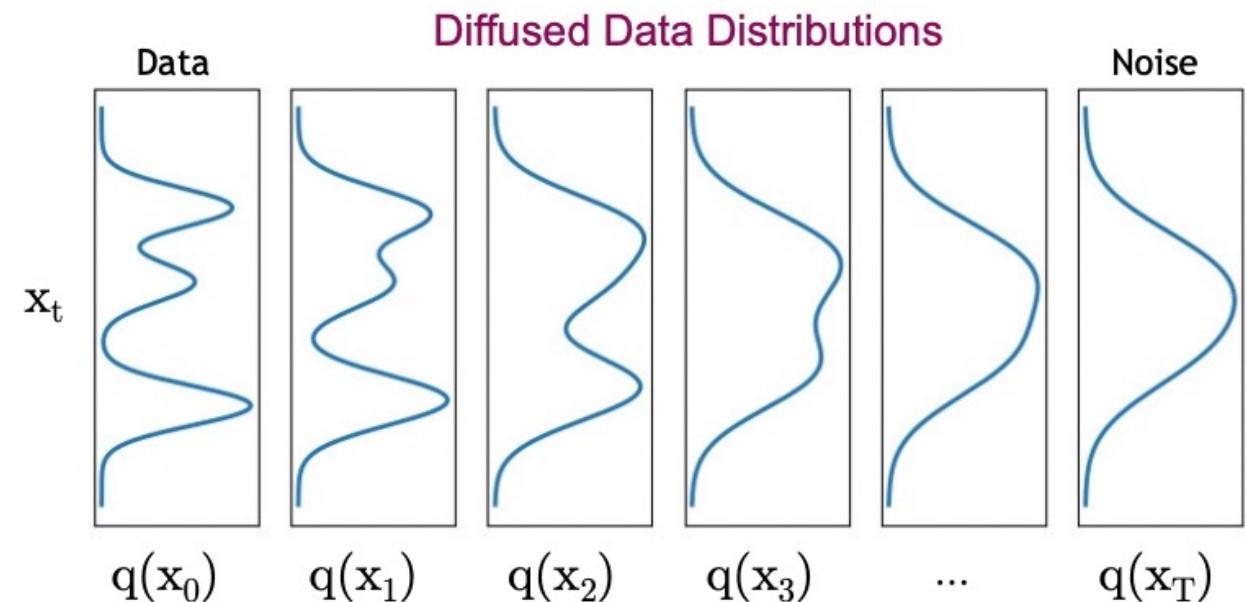
β_t values schedule (i.e., the noise schedule) is designed such that $\bar{\alpha}_T \rightarrow 0$ and $q(\mathbf{x}_T | \mathbf{x}_0) \approx \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$

What happens to a distribution in the forward diffusion?

So far, we discussed the diffusion kernel $q(\mathbf{x}_t|\mathbf{x}_0)$ but what about $q(\mathbf{x}_t)$?

$$q(\mathbf{x}_t) = \underbrace{\int q(\mathbf{x}_0, \mathbf{x}_t) d\mathbf{x}_0}_{\text{Diffused data dist.}} = \underbrace{\int q(\mathbf{x}_0) q(\mathbf{x}_t|\mathbf{x}_0) d\mathbf{x}_0}_{\text{Joint dist.} \quad \text{Input data dist.} \quad \text{Diffusion kernel}}$$

The diffusion kernel is Gaussian convolution.



We can sample $\mathbf{x}_t \sim q(\mathbf{x}_t)$ by first sampling $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ and then sampling $\mathbf{x}_t \sim q(\mathbf{x}_t|\mathbf{x}_0)$ (i.e., ancestral sampling).

What happens to a distribution in the forward diffusion?

DEFINITION

The **marginal probability density functions** of X and Y , denoted by $f_X(x)$ and $f_Y(y)$, respectively, are given by

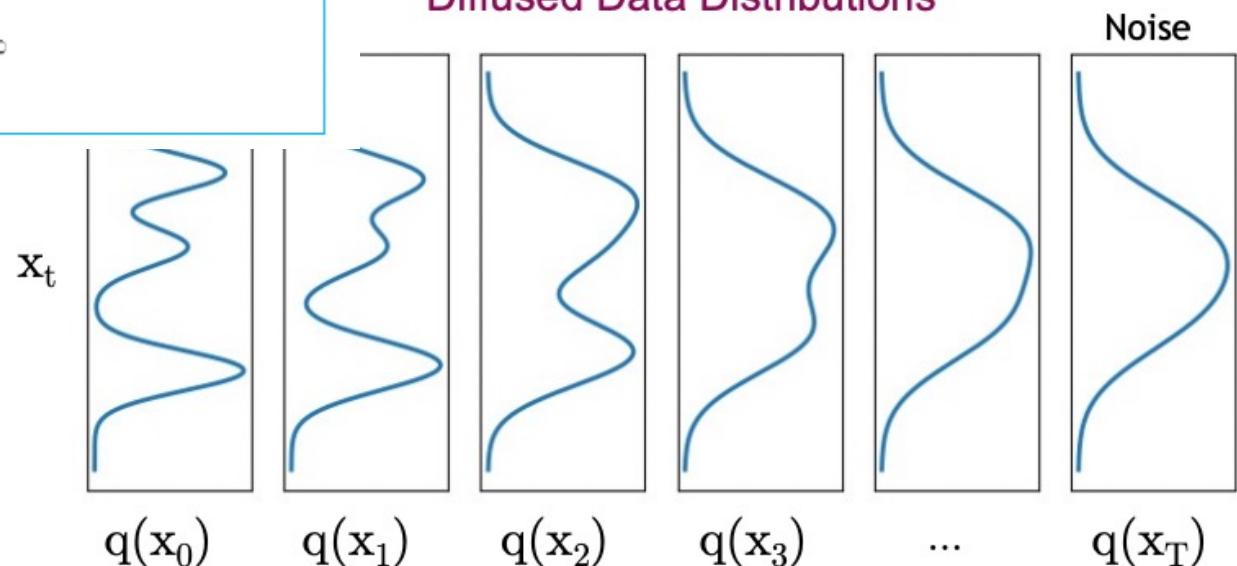
$$f_X(x) = \int_{-\infty}^{\infty} f(x, y) dy \quad \text{for } -\infty < x < \infty$$

$$f_Y(y) = \int_{-\infty}^{\infty} f(x, y) dx \quad \text{for } -\infty < y < \infty$$

$$q(\mathbf{x}_t) = \underbrace{\int q(\mathbf{x}_0, \mathbf{x}_t) d\mathbf{x}_0}_{\text{Diffused data dist.}} = \underbrace{\int q(\mathbf{x}_0) q(\mathbf{x}_t | \mathbf{x}_0) d\mathbf{x}_0}_{\text{Joint dist.}} = \underbrace{q(\mathbf{x}_0)}_{\text{Input data dist.}} \underbrace{q(\mathbf{x}_t | \mathbf{x}_0)}_{\text{Diffusion kernel}}$$

The diffusion kernel is Gaussian convolution.

Diffused Data Distributions



We can sample $\mathbf{x}_t \sim q(\mathbf{x}_t)$ by first sampling $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ and then sampling $\mathbf{x}_t \sim q(\mathbf{x}_t | \mathbf{x}_0)$ (i.e., ancestral sampling).

Generative Learning by Denoising

Recall, that the diffusion parameters are designed such that $q(\mathbf{x}_T) \approx \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$

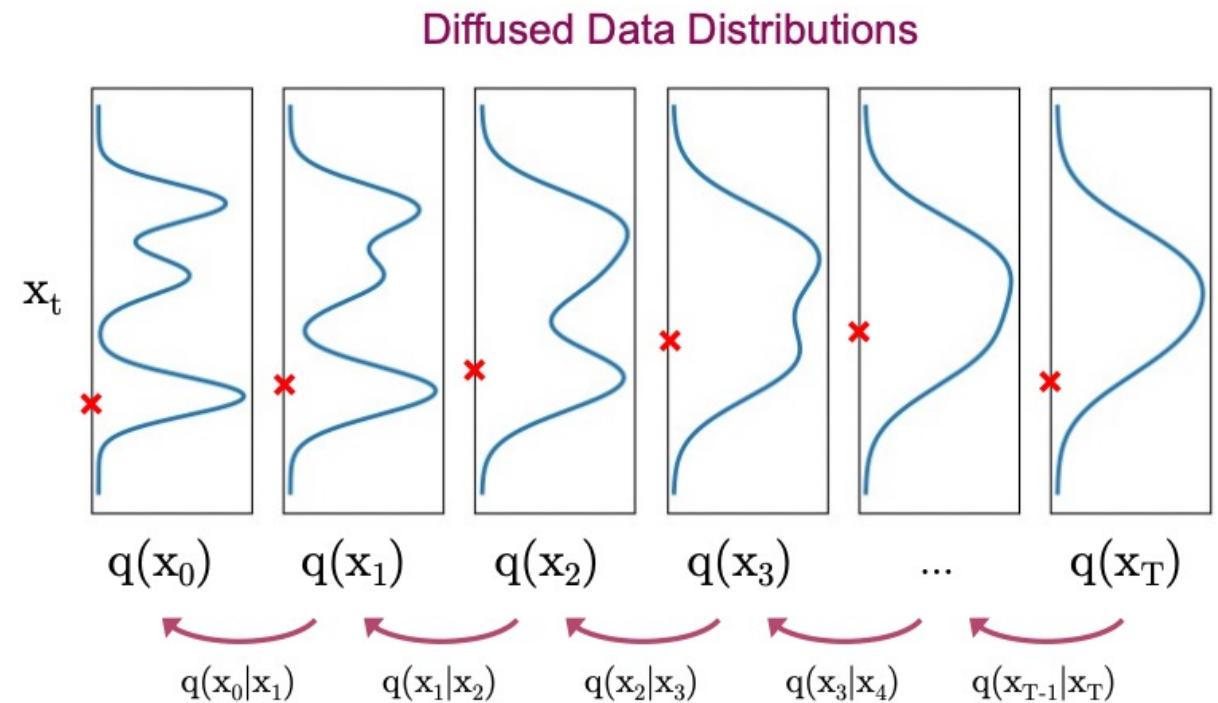
Generation:

Sample $\mathbf{x}_T \sim \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$

Iteratively sample $\mathbf{x}_{t-1} \sim \underbrace{q(\mathbf{x}_{t-1} | \mathbf{x}_t)}_{\text{True Denoising Dist.}}$

In general, $q(\mathbf{x}_{t-1} | \mathbf{x}_t) \propto q(\mathbf{x}_{t-1})q(\mathbf{x}_t | \mathbf{x}_{t-1})$ is intractable.

Can we approximate $q(\mathbf{x}_{t-1} | \mathbf{x}_t)$? Yes, we can use a **Normal distribution** if β_t is small in each forward diffusion step.



Generative Learning by Denoising

Recall, that the diffusion parameters are designed such that $q(\mathbf{x}_T) \approx \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$

Generation:

Sample $\mathbf{x}_T \sim \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$

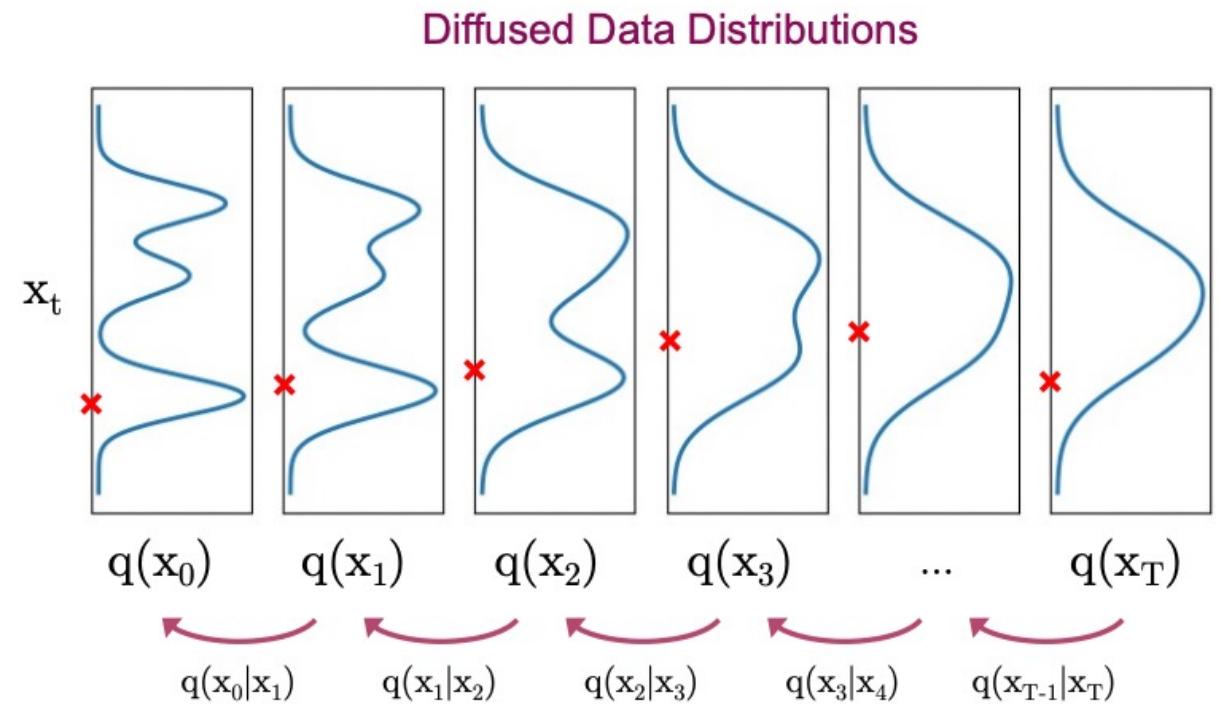
Iteratively sample $\mathbf{x}_{t-1} \sim \underbrace{q(\mathbf{x}_{t-1} | \mathbf{x}_t)}_{\text{True Denoising Dist.}}$

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

(Bayes rule)

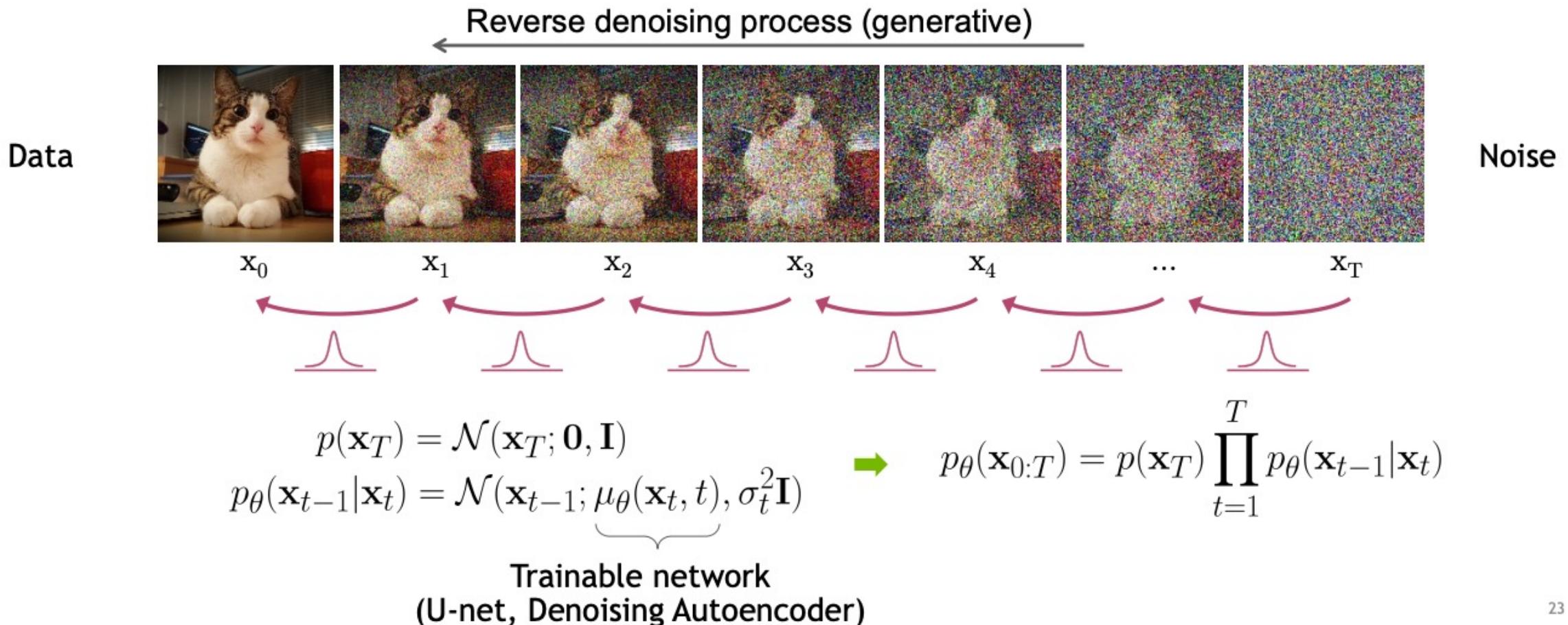
In general, $q(\mathbf{x}_{t-1} | \mathbf{x}_t) \propto q(\mathbf{x}_{t-1})q(\mathbf{x}_t | \mathbf{x}_{t-1})$ is intractable.

Can we approximate $q(\mathbf{x}_{t-1} | \mathbf{x}_t)$? Yes, we can use a **Normal distribution** if β_t is small in each forward diffusion step.



Reverse Process

Formal definition of forward and reverse processes in T steps:



Learning Denoising Model

Variational upper bound

For training, we can form variational upper bound that is commonly used for training variational autoencoders:

$$\mathbb{E}_{q(\mathbf{x}_0)} [-\log p_\theta(\mathbf{x}_0)] \leq \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] =: L$$

Sohl-Dickstein et al. ICML 2015 and Ho et al. NeurIPS 2020 show that:

$$L = \mathbb{E}_q \left[\underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) || p(\mathbf{x}_T))}_{L_T} + \sum_{t>1} \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) || p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))}_{L_{t-1}} \underbrace{- \log p_\theta(\mathbf{x}_0|\mathbf{x}_1)}_{L_0} \right]$$

where $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ is the tractable posterior distribution:

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

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

Learning Denoising Model

Variational upper bound

* Goal: VAE처럼 data likelihood $P_\theta(x) = P_{\text{model}}(x; \theta)$ 을 maximize

⇒ Negative Log Likelihood (NLL)을 사용: $-\log P_\theta(x)$ 을 minimize

$$-\log P_\theta(x_0) \leq -\log P_\theta(x_0) + D_{\text{KL}}(q(X_{1:\tau}|x_0) \| P_\theta(X_{1:\tau}|x_0))$$

$$= -\log P_\theta(x_0) + \mathbb{E}_{X_{1:\tau} \sim q(X_{1:\tau}|x_0)} \left[\log \frac{q(X_{1:\tau}|x_0)}{P_\theta(X_{1:\tau}|x_0)} \right]$$

$$= -\log P_\theta(x_0) + \mathbb{E}_{X_{1:\tau} \sim q(X_{1:\tau}|x_0)} \left[\log \frac{q(X_{1:\tau}|x_0)}{P_\theta(X_{0:\tau})/P_\theta(x_0)} \right]$$

~~$$= -\log P_\theta(x_0) + \mathbb{E}_q \left[\log \frac{q(X_{1:\tau}|x_0)}{P_\theta(X_{0:\tau})} + \log P_\theta(x_0) \right]$$~~

$$= \mathbb{E}_q \left[\log \frac{q(X_{1:\tau}|x_0)}{P_\theta(X_{0:\tau})} \right]$$

* KL Divergence

$$D_{\text{KL}}(P||Q) = \mathbb{E}_{x \sim P} \left[\log \frac{P(x)}{Q(x)} \right]$$

$$= \mathbb{E}_{x \sim P} \left[-\log \frac{Q(x)}{P(x)} \right] \geq 0$$

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

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

$$= \mathbb{E}_{q(x_{0:T})} \left[\log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{0:T})} \right]$$

* Reverse Process : $p_\theta(x_{0:T}) := p(x_T) \prod_{t=1}^T p_\theta(x_{t+1}|x_t)$

Forward Process : $q(x_{1:T}|x_0) := \prod_{t=1}^T q(x_t|x_{t+1})$

$$= \mathbb{E}_q \left[\log \frac{\prod_{t=1}^T q(x_t|x_{t+1})}{p_\theta(x_T) \prod_{t=1}^T p_\theta(x_{t+1}|x_t)} \right]$$

$$= \mathbb{E}_q \left[-\log p_\theta(x_T) + \sum_{t=1}^T \log \frac{q(x_t|x_{t+1})}{p_\theta(x_{t+1}|x_t)} \right]$$

$$= \mathbb{E}_q \left[-\log p_\theta(x_T) + \sum_{t=2}^T \log \frac{q(x_t|x_{t+1})}{p_\theta(x_{t+1}|x_t)} + \log \frac{q(x_1|x_0)}{p_\theta(x_0|x_1)} \right]$$

For training, we

do

$$= \mathbb{E}_q \left[-\log p_\theta(x_T) + \sum_{t=2}^T \log \frac{q(x_t|x_{t+1}, x_0)}{p_\theta(x_{t+1}|x_t)} + \log \frac{q(x_1|x_0)}{p_\theta(x_0|x_1)} \right] \quad (\text{by Markov Property})$$

$$= \mathbb{E}_q \left[-\log p_\theta(x_T) + \sum_{t=2}^T \log \frac{1}{p_\theta(x_{t+1}|x_t)} \cdot \frac{q(x_{t+1}|x_t, x_0)}{q(x_{t+1}|x_0)} + \log \frac{q(x_1|x_0)}{p_\theta(x_0|x_1)} \right]$$

$$= \mathbb{E}_q \left[-\log p_\theta(x_T) + \sum_{t=2}^T \log \frac{q(x_{t+1}|x_t, x_0)}{p_\theta(x_{t+1}|x_t)} + \sum_{t=2}^T \log \frac{q(x_t|x_0)}{q(x_{t+1}|x_0)} + \log \frac{q(x_1|x_0)}{p_\theta(x_0|x_1)} \right]$$

Sohl-Dickstein et al.

$$= \mathbb{E}_q \left[-\log p_\theta(x_T) + \sum_{t=2}^T \log \frac{q(x_{t+1}|x_t, x_0)}{p_\theta(x_{t+1}|x_t)} + \log \left(\frac{q(x_2|x_0)}{q(x_1|x_0)} \times \frac{q(x_3|x_0)}{q(x_2|x_0)} \times \dots \times \frac{q(x_T|x_0)}{q(x_{T-1}|x_0)} \right) + \log \frac{q(x_1|x_0)}{p_\theta(x_0|x_1)} \right]$$

$$= \mathbb{E}_q \left[-\log p_\theta(x_T) + \sum_{t=2}^T \log \frac{q(x_{t+1}|x_t, x_0)}{p_\theta(x_{t+1}|x_t)} + \log \frac{q(x_T|x_0)}{q(x_1|x_0)} + \log \frac{q(x_1|x_0)}{p_\theta(x_0|x_1)} \right]$$

$$= \mathbb{E}_q \left[\log \frac{q(x_T|x_0)}{p_\theta(x_0|x_1)p_\theta(x_T)} + \sum_{t=2}^T \log \frac{q(x_{t+1}|x_t, x_0)}{p_\theta(x_{t+1}|x_t)} \right]$$

$$= \mathbb{E}_q \left[\log \frac{q(x_T|x_0)}{p_\theta(x_T)} + \sum_{t=2}^T \log \frac{q(x_{t+1}|x_t, x_0)}{p_\theta(x_{t+1}|x_t)} - \log p_\theta(x_0|x_1) \right]$$

$$= \mathbb{E}_q \left[D_{KL}(q(x_T|x_0) || p_\theta(x_T)) + \sum_{t=2}^T D_{KL}(q(x_{t+1}|x_t, x_0) || p_\theta(x_{t+1}|x_t)) - \log p_\theta(x_0|x_1) \right]$$

where $q(x_{t-1}|x_t)$

L_T

L_{t-1}

L_0

Learning Denoising Model

Variational upper bound

For training, we can form variational upper bound that is commonly used for training variational autoencoders:

$$\mathbb{E}_{q(\mathbf{x}_0)} [-\log p_\theta(\mathbf{x}_0)] \leq \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] =: L$$

Sohl-Dickstein et al. ICML 2015 and Ho et al. NeurIPS 2020 show that:

$$L = \mathbb{E}_q \left[\underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) || p(\mathbf{x}_T))}_{L_T} + \sum_{t>1} \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) || p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))}_{L_{t-1}} \underbrace{- \log p_\theta(\mathbf{x}_0|\mathbf{x}_1)}_{L_0} \right]$$

where $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ is the tractable posterior distribution:

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

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

- 정규분포 pdf $f(x) = \frac{1}{\sqrt{2\pi} \cdot \sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$
 - $q(X_t | X_0) = N(X_t; \sqrt{\alpha_t} X_0, (1 - \bar{\alpha}_t) I)$
 - $q(X_{t-1} | X_t) = N(X_{t-1}; \sqrt{1 - \bar{\alpha}_t} X_t, \beta_t I)$
- $= q(X_{t-1} | X_t, X_0)$ (by Markov chain Rule)

$$q(X_{t-1} | X_t, X_0) = q(X_t | X_{t-1}, X_0) \cdot \frac{q(X_{t-1} | X_0)}{q(X_t | X_0)} \quad (\text{by Bayes rule})$$

(Reverse Process)

$$\propto \exp \left(-\frac{1}{2} \left(\frac{(X_t - \sqrt{\alpha_t} X_{t-1})^2}{\beta_t} + \frac{(X_{t-1} - \sqrt{\alpha_{t-1}} X_0)^2}{1 - \bar{\alpha}_{t-1}} - \frac{(X_t - \sqrt{\alpha_t} X_0)^2}{1 - \bar{\alpha}_t} \right) \right)$$

$$= \exp \left(-\frac{1}{2} \left(\frac{X_t^2 - 2\sqrt{\alpha_t} X_t X_{t-1} + \alpha_t X_{t-1}^2}{\beta_t} + \frac{X_{t-1}^2 - 2\sqrt{\alpha_{t-1}} X_0 X_{t-1} + \alpha_{t-1} X_0^2}{1 - \bar{\alpha}_{t-1}} - \frac{(X_t - \sqrt{\alpha_t} X_0)^2}{1 - \bar{\alpha}_t} \right) \right)$$

$$= \exp \left(-\frac{1}{2} \left(\left(\frac{\alpha_t}{\beta_t} + \frac{1}{1 - \bar{\alpha}_{t-1}} \right) X_{t-1}^2 - \left(\frac{2\sqrt{\alpha_t}}{\beta_t} X_t + \frac{2\sqrt{\alpha_{t-1}}}{1 - \bar{\alpha}_{t-1}} X_0 \right) X_{t-1} + C(X_t, X_0) \right) \right)$$

$$\left(N \sim (\mu, \sigma^2) \text{의 pdf: } \frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{1}{2} \left(\frac{1}{\sigma^2} X^2 - \frac{2\mu}{\sigma^2} X + \frac{\mu^2}{\sigma^2} \right) \right) \right)$$

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

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

$$q_t(X_{t-1} | X_t, x_0) = N(X_{t-1}; \tilde{\mu}(x_t, x_0), \tilde{\Sigma}_t I)$$

$$\textcircled{1} \quad \tilde{\Sigma}_t = \frac{1}{\left(\frac{d_t}{\varepsilon_t} + \frac{1}{1 - d_{t-1}} \right)} = \frac{1}{\left(\frac{d_t - d_t \bar{d}_{t-1} + \varepsilon_t}{\varepsilon_t (1 - \bar{d}_{t-1})} \right)} = \frac{1 - \bar{d}_{t-1}}{1 - \bar{d}_t} \cdot \varepsilon_t \quad \left(\begin{array}{l} \because d_t = 1 - \varepsilon_t, \\ \bar{d}_t = \pi_{i=1}^T d_i \end{array} \right)$$

$$\textcircled{2} \quad \tilde{\mu}(X_t, x_0) = \left(\frac{\sqrt{d_t}}{\varepsilon_t} X_t + \frac{\sqrt{d_{t-1}}}{1 - \bar{d}_{t-1}} x_0 \right) / \left(\frac{d_t}{\varepsilon_t} + \frac{1}{1 - \bar{d}_{t-1}} \right)$$

$$= \left(\frac{\sqrt{d_t}}{\varepsilon_t} X_t + \frac{\sqrt{d_{t-1}}}{1 - \bar{d}_{t-1}} x_0 \right) \cdot \left(\frac{1 - \bar{d}_{t-1}}{1 - \bar{d}_t} \cdot \varepsilon_t \right) \quad (\text{by } \textcircled{1})$$

$$X_t = \sqrt{d_t} X_t + \sqrt{1 - \bar{d}_t} \varepsilon_t \leftarrow$$

$$X_0 = \frac{1}{\sqrt{d_0}} (X_t - \sqrt{1 - \bar{d}_t} \varepsilon_t)$$

$$= \frac{\sqrt{d_t} (1 - \bar{d}_{t-1})}{1 - \bar{d}_t} X_t + \frac{\sqrt{d_{t-1}} \varepsilon_t}{1 - \bar{d}_t} X_0$$

$$= \frac{\sqrt{d_t} (1 - \bar{d}_{t-1})}{1 - \bar{d}_t} X_t + \frac{\sqrt{d_{t-1}} \varepsilon_t}{1 - \bar{d}_t} \left(\frac{1}{\sqrt{d_0}} (X_t - \sqrt{1 - \bar{d}_t} \varepsilon_t) \right)$$

$$= \left(\frac{\sqrt{d_t} (1 - \bar{d}_{t-1})}{1 - \bar{d}_t} + \frac{\sqrt{d_{t-1}} \varepsilon_t}{1 - \bar{d}_t \sqrt{d_0}} \right) X_t - \frac{\sqrt{d_{t-1}} \varepsilon_t}{1 - \bar{d}_t \sqrt{d_0}} \sqrt{1 - \bar{d}_t} \varepsilon_t$$

$$= \left(\frac{\sqrt{d_t} (1 - \bar{d}_{t-1})}{1 - \bar{d}_t} + \frac{\sqrt{d_{t-1}} (1 - d_t)}{1 - \bar{d}_t \sqrt{d_0}} \right) X_t - \frac{\sqrt{d_{t-1}} (1 - d_t)}{\sqrt{1 - \bar{d}_t} \sqrt{d_0}} \varepsilon_t$$

$$= \left(\frac{\sqrt{d_t} (1 - \bar{d}_{t-1})}{1 - \bar{d}_t} + \frac{\sqrt{d_{t-1}} (1 - d_t)}{1 - \bar{d}_t \sqrt{d_0}} \right) X_t - \frac{\sqrt{d_{t-1}} (1 - d_t)}{\sqrt{d_t} \sqrt{1 - \bar{d}_t}} \varepsilon_t$$

$$= \frac{\sqrt{d_t} \sqrt{d_0} (1 - \bar{d}_{t-1}) + \sqrt{d_{t-1}} - d_t (\bar{d}_{t-1})}{(1 - \bar{d}_t) \sqrt{d_{t-1}} \sqrt{d_t}} X_t - \frac{1}{\sqrt{d_t} \sqrt{1 - \bar{d}_t}} \sqrt{d_{t-1}} (1 - d_t) \varepsilon_t$$

$$= \frac{d_t (\cancel{\sqrt{d_{t-1}}}) (1 - \bar{d}_{t-1}) + \cancel{\sqrt{d_{t-1}}} (1 - d_t)}{(1 - \bar{d}_t) \cancel{\sqrt{d_{t-1}}} \sqrt{d_t}} X_t - \frac{1}{\sqrt{d_t} \cancel{\sqrt{d_{t-1}}}} \cancel{\sqrt{d_{t-1}}} \frac{(1 - d_t)}{\sqrt{1 - \bar{d}_t}} \varepsilon_t$$

$$= \frac{\cancel{d_t - \bar{d}_t} + \cancel{1 - d_t}}{(1 - \bar{d}_t) \sqrt{d_t}} X_t - \frac{1}{\sqrt{d_t}} \frac{(1 - d_t)}{\sqrt{1 - \bar{d}_t}} \varepsilon_t$$

$$= \frac{1}{\sqrt{d_t}} \left(X_t - \frac{1 - d_t}{\sqrt{1 - \bar{d}_t}} \varepsilon_t \right)$$

Parameterizing the Denoising Model

Since both $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ and $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ are Normal distributions, the KL divergence has a simple form:

$$L_{t-1} = D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) || p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)) = \mathbb{E}_q \left[\frac{1}{2\sigma_t^2} \|\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) - \mu_\theta(\mathbf{x}_t, t)\|^2 \right] + C$$

Recall that $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{(1 - \bar{\alpha}_t)} \epsilon$. [Ho et al. NeurIPS 2020](#) observe that:

$$\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{1 - \beta_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon \right)$$

They propose to represent the mean of the denoising model using a *noise-prediction* network:

$$\mu_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{1 - \beta_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right)$$

With this parameterization

$$L_{t-1} = \mathbb{E}_{\mathbf{x}_0 \sim q(\mathbf{x}_0), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\frac{\beta_t^2}{2\sigma_t^2(1 - \beta_t)(1 - \bar{\alpha}_t)} \|\epsilon - \underbrace{\epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)}_{\mathbf{x}_t}\|^2 \right] + C$$

Training Objective Weighting

Trading likelihood for perceptual quality

$$L_{t-1} = \mathbb{E}_{\mathbf{x}_0 \sim q(\mathbf{x}_0), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\underbrace{\frac{\beta_t^2}{2\sigma_t^2(1 - \beta_t)(1 - \bar{\alpha}_t)} ||\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)||^2}_{\lambda_t} \right]$$

The time dependent λ_t ensures that the training objective is weighted properly for the maximum data likelihood training.

However, this weight is often very large for small t's.

[Ho et al. NeurIPS 2020](#) observe that simply setting $\lambda_t = 1$ improves sample quality. So, they propose to use:

$$L_{\text{simple}} = \mathbb{E}_{\mathbf{x}_0 \sim q(\mathbf{x}_0), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), t \sim \mathcal{U}(1, T)} \left[||\epsilon - \underbrace{\epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)}_{\mathbf{x}_t}||^2 \right]$$

For more advanced weighting see [Choi et al., Perception Prioritized Training of Diffusion Models, CVPR 2022](#).

Summary

Training and Sample Generation

Algorithm 1 Training

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

6: until converged
```

Algorithm 2 Sampling

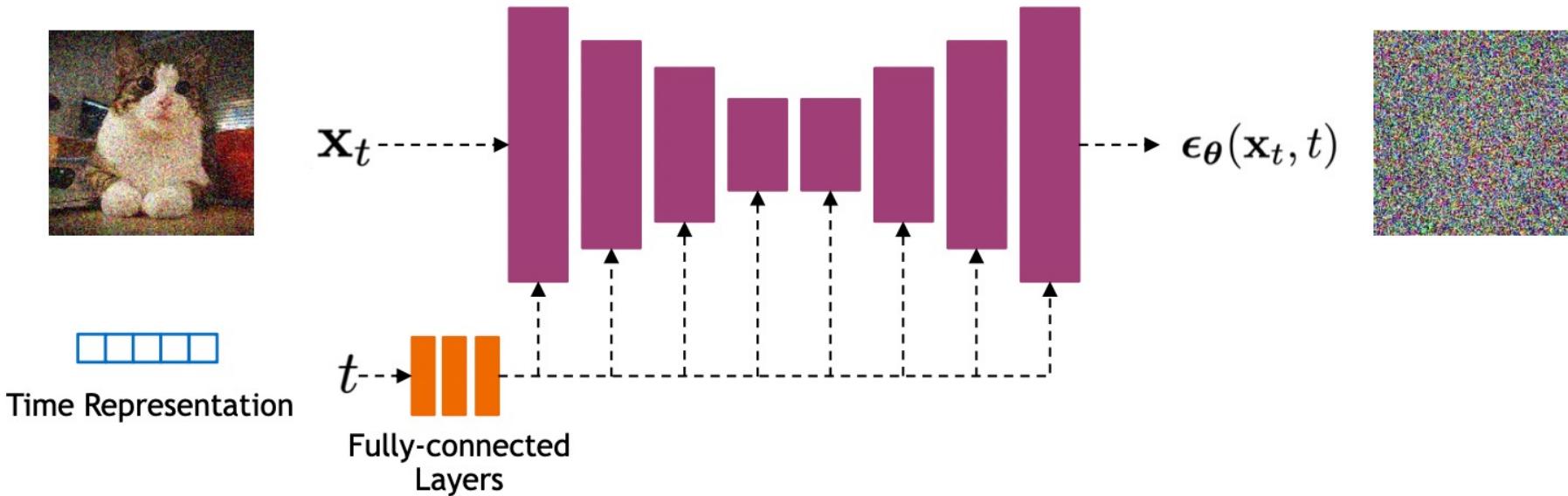
```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
4:   
$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$$

5: end for
6: return  $\mathbf{x}_0$ 
```

Implementation Considerations

Network Architectures

Diffusion models often use U-Net architectures with ResNet blocks and self-attention layers to represent $\epsilon_\theta(\mathbf{x}_t, t)$

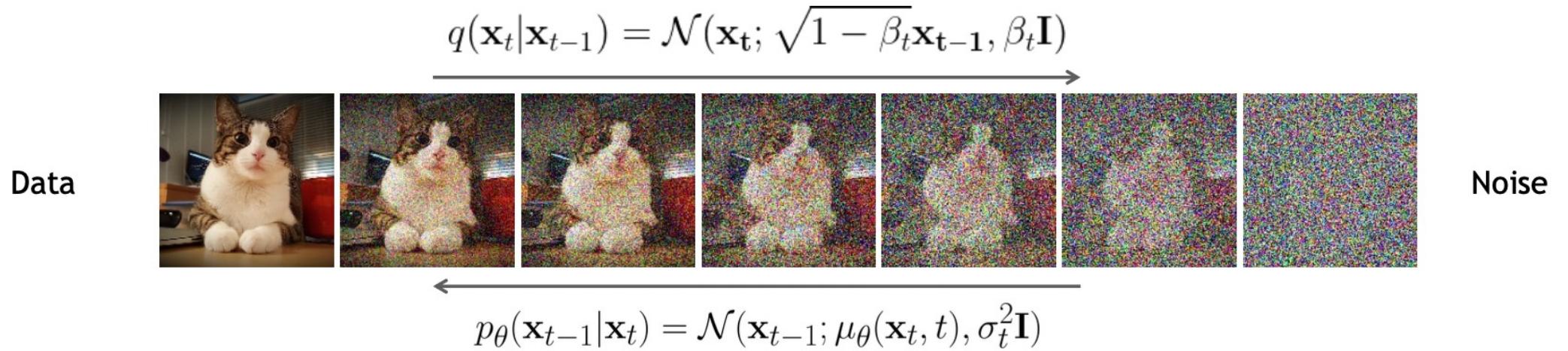


Time representation: sinusoidal positional embeddings or random Fourier features.

Time features are fed to the residual blocks using either simple spatial addition or using adaptive group normalization layers. (see [Dhariwal and Nichol NeurIPS 2021](#))

Diffusion Parameters

Noise Schedule



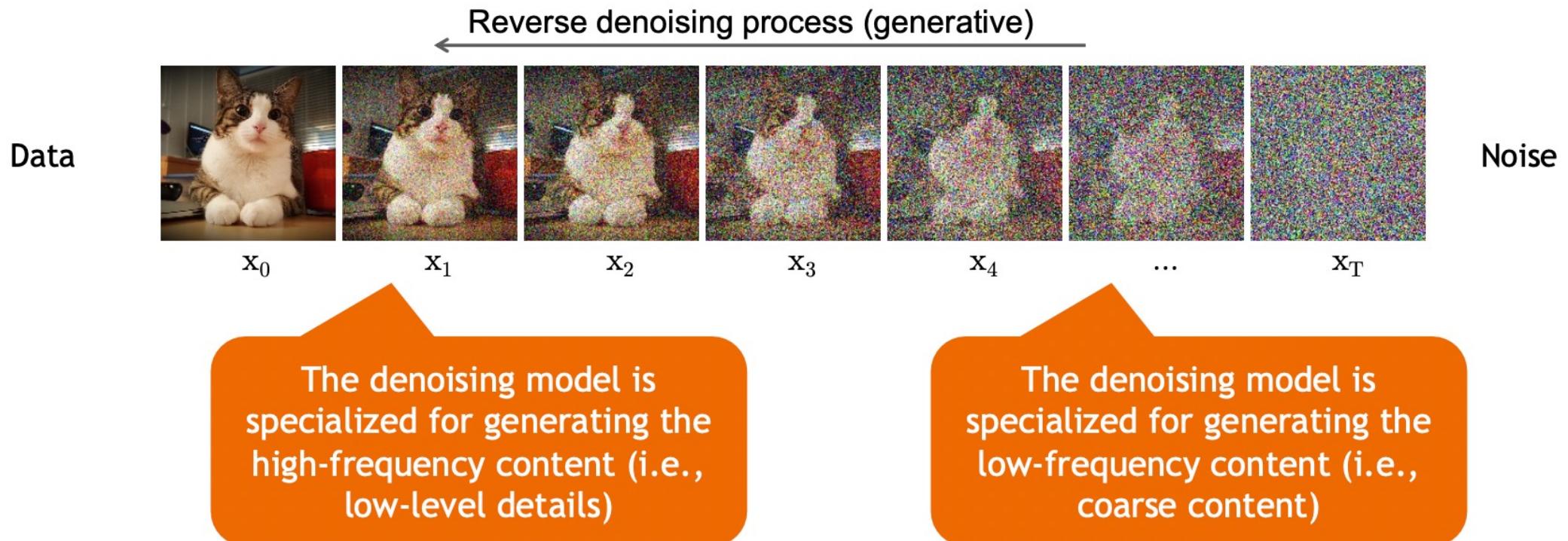
Above, β_t and σ_t^2 control the variance of the forward diffusion and reverse denoising processes respectively.

Often a linear schedule is used for β_t , and σ_t^2 is set equal to β_t .

[Kingma et al. NeurIPS 2022](#) introduce a new parameterization of diffusion models using signal-to-noise ratio (SNR), and show how to learn the noise schedule by minimizing the variance of the training objective.

We can also train σ_t^2 while training the diffusion model by minimizing the variational bound ([Improved DPM by Nichol and Dhariwal ICML 2021](#)) or after training the diffusion model ([Analytic-DPM by Bao et al. ICLR 2022](#)).

Content-Detail Tradeoff



The weighting of the training objective for different timesteps is important!

Summary

Denoising Diffusion Probabilistic Models

In this part, we reviewed denoising diffusion probabilistic models.

The model is trained by sampling from the forward diffusion process and training a denoising model to predict the noise.

We discussed how the forward process perturbs the data distribution or data samples.

The devil is in the details:

- Network architectures
- Objective weighting
- Diffusion parameters (i.e., noise schedule)

See [“Elucidating the Design Space of Diffusion-Based Generative Models” by Karras et al.](#) for important design decisions.

Reference

- <https://cvpr2022-tutorial-diffusion-models.github.io/>
- <https://minibatchai.com/diffusion/generative-models/text2image/2022/06/17/Diffusion.html>
- <https://lilianweng.github.io/posts/2021-07-11-diffusion-models/#nice>
- <https://happy-jihye.github.io/diffusion/diffusion-1/>

감사합니다