# Diffusion Model

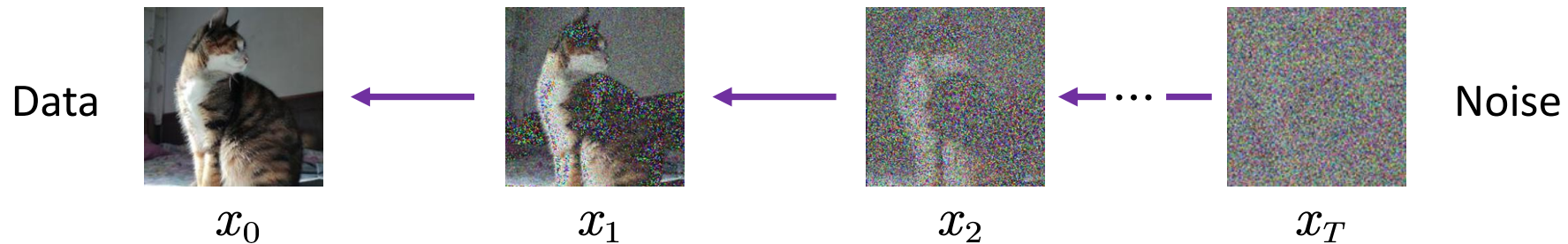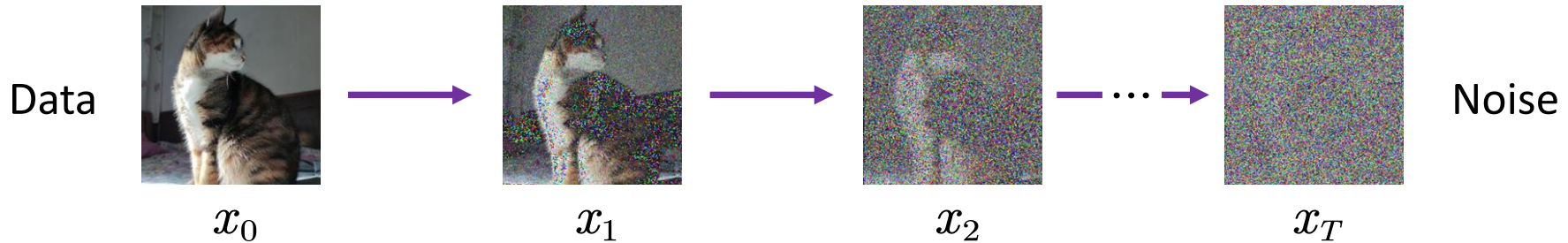## Denoising Diffusion Probabilistic Models

- What is the diffusion model?
- How to visually understand the diffusion model?
- How to derive the diffusion model mathematically?
- How to train a diffusion model and infer it?

Xin Zhang

# What is Diffusion Model ?

## **D**enoising **D**iffusion **P**robabilistic **M**odels[1]
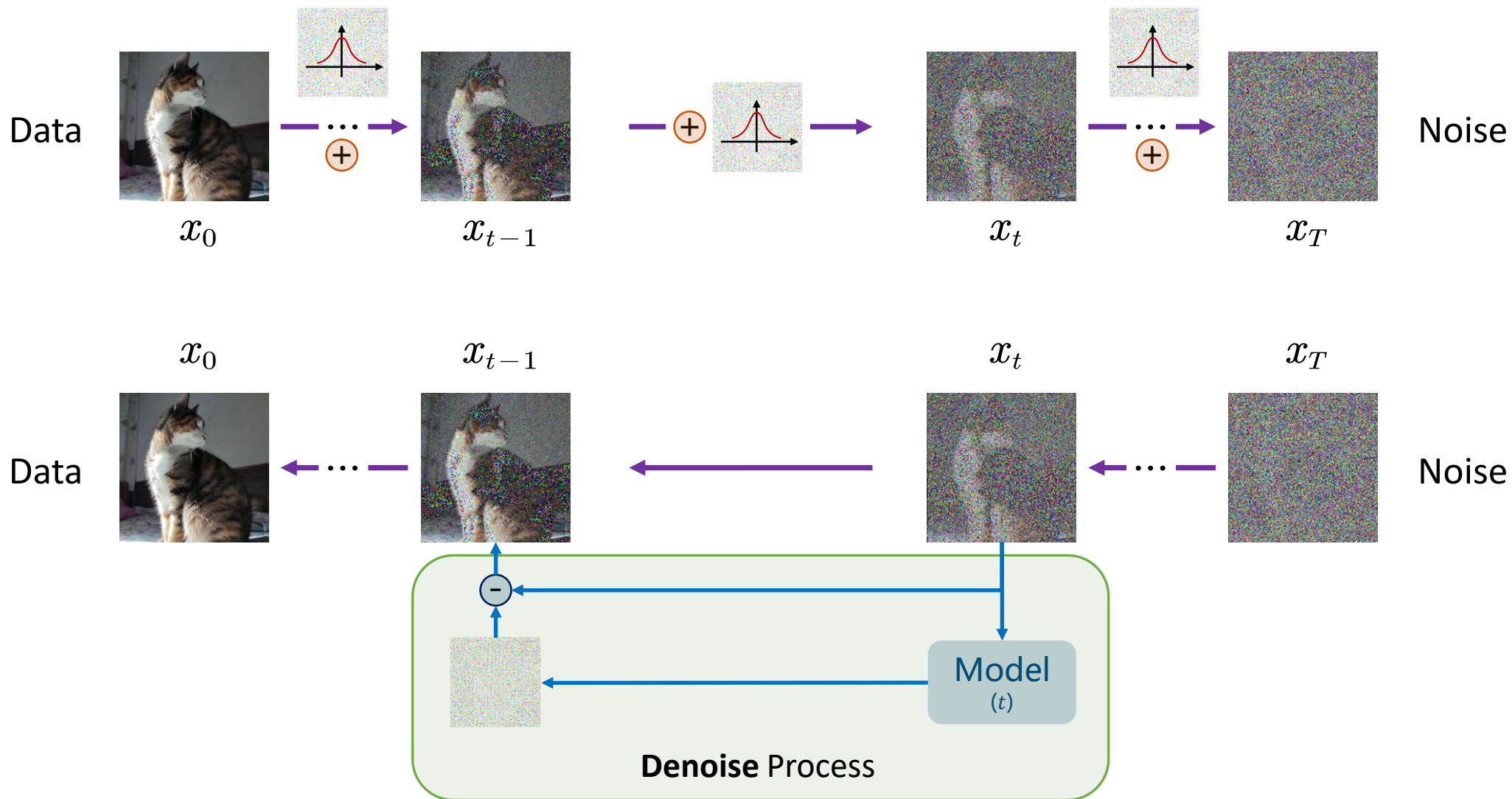
[1] Ho J, Jain A, Abbeel P. Denoising diffusion probabilistic models. NeuraIPS, 2020.

# What is Diffusion Model ?



Data    $x_0$    $x_{t-1}$    $x_t$    $x_T$    Noise

Data    $x_0$    $x_{t-1}$    $x_t$    $x_T$    Noise
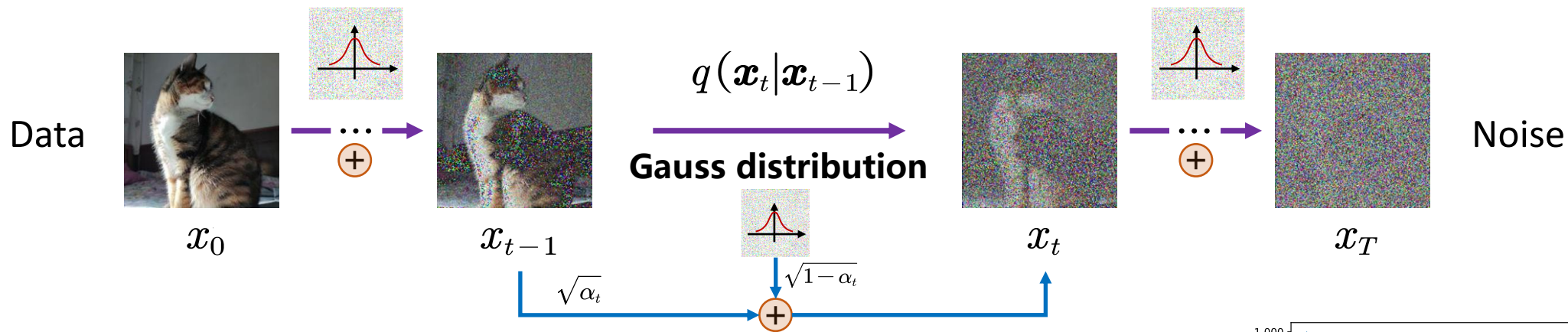
Model $(t)$

**Denoise** Process

" 

**The sculpture is already complete within the marble block before I start my work. It is already there, I just have to chisel away the superfluous material.**

"

**——Michelangelo**

# **Forward** Diffusion Process

Data

$$q(\boldsymbol{x}_t|\boldsymbol{x}_{t-1})$$

**Gauss distribution**

Noise

$$x_0 \qquad\qquad x_{t-1} \qquad\qquad x_t \qquad\qquad x_T$$

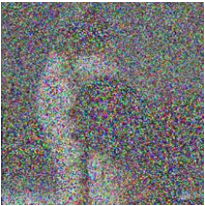$$\sqrt{\alpha_t} \qquad\qquad \sqrt{1-\alpha_t}$$

$$z \sim \mathcal{N}(\mu,\ \sigma^2)$$

$$\frac{z-\mu}{\sigma} \sim \mathcal{N}(0,\ I)$$

$$z = \mu + \sigma \cdot \varepsilon$$

$$\varepsilon \sim \mathcal{N}(0,I)$$

gaussian  $=\ \sqrt{\alpha_t}$  signal  $+\ \sqrt{1-\alpha_t}$  noise

$$\boldsymbol{x}_t = \sqrt{\alpha_t}\,\boldsymbol{x}_{t-1} + \sqrt{1-\alpha_t}\,\boldsymbol{\varepsilon}_{t-1}$$

$$q(\boldsymbol{x}_t|\boldsymbol{x}_{t-1}) = \mathcal{N}\big(\boldsymbol{x}_t; \underset{\text{mean}}{\sqrt{\alpha_t}\,\boldsymbol{x}_{t-1}}, \underset{\text{variance}}{(1-\alpha_t)\mathbf{I}}\big)$$

$$\sqrt{\alpha_t}$$

$$\sqrt{1-\alpha_t} = \sqrt{\beta_t}$$

# **Forward** Diffusion Process

$$\boldsymbol{x}_t = \sqrt{\alpha_t}\,\boldsymbol{x}_{t-1} + \sqrt{1-\alpha_t}\,\boldsymbol{\varepsilon}_{t-1}$$

$$= \sqrt{\alpha_t}\left(\sqrt{\alpha_{t-1}}\,x_{t-2} + \sqrt{1-\alpha_{t-1}}\,\boldsymbol{\varepsilon}_{t-2}\right) + \sqrt{1-\alpha_t}\,\boldsymbol{\varepsilon}_{t-1}$$

$$= \sqrt{\alpha_t\alpha_{t-1}}\,x_{t-2} + \left(\sqrt{\alpha_t - \alpha_t\alpha_{t-1}}\,\boldsymbol{\varepsilon}_{t-2} + \sqrt{1-\alpha_t}\,\boldsymbol{\varepsilon}_{t-1}\right)$$

...

$$\boxed{\begin{array}{c} \text{If } X \sim \mathcal{N}(\mu_X, \sigma_X^2) \quad Y \sim \mathcal{N}(\mu_Y, \sigma_Y^2) \\ Z = X + Y \\ \text{Then } Z \sim \mathcal{N}(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2) \end{array}}$$

$$= \sqrt{\alpha_t\alpha_{t-1}\cdots\alpha_1}\,\boldsymbol{x}_0 + \sqrt{1-\alpha_t\alpha_{t-1}\cdots\alpha_1}\,\boldsymbol{\varepsilon}$$

$$= \sqrt{\overline{\alpha}_t}\,\boldsymbol{x}_0 + \sqrt{1-\overline{\alpha}_t}\,\boldsymbol{\varepsilon}$$



$$x_{t-1} \qquad x_t$$

$$\boldsymbol{x}_t = \sqrt{\alpha_t}\,\boldsymbol{x}_{t-1} + \sqrt{1-\alpha_t}\,\boldsymbol{\varepsilon_t}$$



$$x_0 \qquad x_t$$

$$\boldsymbol{x}_t = \sqrt{\overline{\alpha}_t}\,\boldsymbol{x}_0 + \sqrt{1-\overline{\alpha}_t}\,\boldsymbol{\varepsilon_t}$$

$$\overline{\alpha}_t = \alpha_1\alpha_2...\alpha_t$$

# **Reverse** Diffusion Process



Markov

$$q(\boldsymbol{x}_t|\boldsymbol{x}_{t-1})$$

Data

Noise

$$x_0 \qquad x_{t-1}$$

$$q(\boldsymbol{x}_{t-1}|\boldsymbol{x}_t)$$

impossible!

$$x_t \qquad x_T$$

$$q(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_t, \boldsymbol{x}_0) \ \textbf{?}$$

$$p_\theta(\boldsymbol{x}_{t-1}|\boldsymbol{x}_t) \quad \text{Neural Network}$$

Assume: the output is gaussian

**Target Distribution**
$$q(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_t(x_t), \Sigma_t(x_t))$$

**Approximated Distribution**
$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$$

# Maximum Likelihood Estimation



$$\arg\max_{\theta} \prod_{i=1}^{t} p_\theta\left(\boldsymbol{x}_i\right)$$

$$\arg\max_{\theta} \sum_{i=1}^{t} \log p_\theta\left(\boldsymbol{x}_i\right)$$

**② optimization (view 1)**

$$\min \ -\log p_\theta(x_0) \leqslant -\log p_\theta(x_0) + D_{KL}\left(q(x_{1:T}|x_0)||p_\theta(x_{1:T}|x_0)\right)$$

$$\min \ -\log p_\theta(x_0) \leqslant \mathbb{E}_{q(x_{1:T}|x_0)}\left[\log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{0:T})}\right] \quad \textbf{ELBO}$$

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

[1] Luo C. Understanding diffusion models: A unified perspective. arXiv, 2022.

# What is $q(x_{t-1}|x_t, x_0)$



$x_t$ $\quad\quad$ $x_{t-1}$ $\quad\quad$ $x_0$

If we know $x_0$ and $x_t$

$q(x_{t-1}|x_t, x_0)$ is deterministic

**Assume: Markov**

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

$$q(x_t \mid x_{t-1}) \sim \mathcal{N}\big(x_t;\ \sqrt{\alpha_t}\, x_{t-1},\ 1 - \alpha_t\big)$$

$$q(x_{t-1} \mid x_0) \sim \mathcal{N}\big(x_{t-1};\ \sqrt{\overline{\alpha}_{t-1}}\, x_0,\ 1 - \overline{\alpha}_{t-1}\big)$$

$$q(x_t \mid x_0) \sim \mathcal{N}\big(x_t;\ \sqrt{\overline{\alpha}_t}\, x_0,\ 1 - \overline{\alpha}_t\big)$$



[1] Luo C. Understanding diffusion models: A unified perspective. arXiv, 2022.

# What is $q(x_{t-1}|x_t, x_0)$

$$q(x_t \mid x_{t-1}) \sim \mathcal{N}(x_t; \sqrt{\alpha_t}\, x_{t-1}, 1 - \alpha_t)$$

$$q(x_{t-1} \mid x_0) \sim \mathcal{N}(x_{t-1}; \sqrt{\bar{\alpha}_{t-1}}\, x_0, 1 - \bar{\alpha}_{t-1})$$

$$q(x_t \mid x_0) \sim \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}\, x_0, 1 - \bar{\alpha}_t)$$

③ **reverse**    If we know $x_0$

$$q(\boldsymbol{x}_{t-1}|\boldsymbol{x}_t, \boldsymbol{x}_0) = \mathcal{N}\left(\boldsymbol{x}_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})\boldsymbol{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1-\alpha_t)\boldsymbol{x}_0}{1-\bar{\alpha}_t}}_{\boldsymbol{\mu}_q(\boldsymbol{x}_t, t)}, \underbrace{\frac{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}\boldsymbol{I}}_{\Sigma_q(t)}\right)$$

Assume: fixed

**Minimize the distance between two Gaussian distributions (refer to PRML)**

$$D_{KL}\left(\mathcal{N}(\boldsymbol{x}; \boldsymbol{\mu}_x, \Sigma_x) \| \mathcal{N}(\boldsymbol{y}; \boldsymbol{\mu}_y, \Sigma_y)\right) = \frac{1}{2}\left[\log\frac{|\Sigma_y|}{|\Sigma_x|} - d + \mathrm{tr}(\Sigma_y^{-1}\Sigma_x) + (\boldsymbol{\mu}_y - \boldsymbol{\mu}_x)^T \Sigma_y^{-1}(\boldsymbol{\mu}_y - \boldsymbol{\mu}_x)\right]$$

$$\underset{\boldsymbol{\theta}}{\mathrm{argmin}}\ D_{KL}\left(q(\boldsymbol{x}_{t-1}|\boldsymbol{x}_t, \boldsymbol{x}_0) \| p_\theta(\boldsymbol{x}_{t-1}|\boldsymbol{x}_t)\right)$$

$$= \underset{\boldsymbol{\theta}}{\mathrm{argmin}}\ \frac{1}{2\sigma_q^2(t)}\left[\|\boldsymbol{\mu}_\theta - \boldsymbol{\mu}_q\|_2^2\right]$$

[1] Luo C. Understanding diffusion models: A unified perspective. arXiv, 2022.

# Remove $x_0$

③ **reverse**    If we know $x_0$

$$q\left(\boldsymbol{x}_{t-1}|\boldsymbol{x}_t,\boldsymbol{x}_0\right)=\mathcal{N}\left(\boldsymbol{x}_{t-1};\underbrace{\frac{\sqrt{\alpha_t}\left(1-\overline{\alpha}_{t-1}\right)\boldsymbol{x}_t+\sqrt{\overline{\alpha}_{t-1}}\left(1-\alpha_t\right)\boxed{\boldsymbol{x}_0}}{1-\overline{\alpha}_t}}_{\boldsymbol{\mu}_q\left(\boldsymbol{x}_t,t\right)},\underbrace{\frac{\left(1-\alpha_t\right)\left(1-\overline{\alpha}_{t-1}\right)}{1-\overline{\alpha}_t}\boldsymbol{I}}_{\Sigma_q\left(t\right)}\right)$$

① **forward (close-form)**    $\boldsymbol{x}_t=\sqrt{\overline{\alpha}_t}\,\boldsymbol{x}_0+\sqrt{1-\overline{\alpha}_t}\,\boldsymbol{\varepsilon}_t$    $\Longrightarrow$    $\boldsymbol{x}_0=\dfrac{\boldsymbol{x}_t-\sqrt{1-\overline{\alpha}_t}\,\boldsymbol{\varepsilon}_t}{\sqrt{\overline{\alpha}_t}}$

$$\frac{\sqrt{\alpha_t}\left(1-\overline{\alpha}_{t-1}\right)\boldsymbol{x}_t+\sqrt{\overline{\alpha}_{t-1}}\left(1-\alpha_t\right)\boldsymbol{x}_0}{1-\overline{\alpha}_t}\Longrightarrow\frac{1}{\sqrt{\alpha_t}}\left(\boldsymbol{x}_t-\frac{1-\alpha_t}{\sqrt{1-\overline{\alpha}_t}}\boldsymbol{\varepsilon}_t\right)$$

noise predictor

$$\boldsymbol{\varepsilon}_t\left(\boldsymbol{x}_0\to\boldsymbol{x}_t\right)$$

Why not predict $\boldsymbol{x}_0$ directly?

[1] Luo C. Understanding diffusion models: A unified perspective. arXiv, 2022.

# Training and Sampling

**③ reverse**

$$q(\boldsymbol{x}_{t-1}|\boldsymbol{x}_t,\boldsymbol{x}_0) = \mathcal{N}\left(\boldsymbol{x}_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1-\overline{\alpha}_{t-1})\boldsymbol{x}_t + \sqrt{\overline{\alpha}_{t-1}}(1-\alpha_t)\boxed{\boldsymbol{x}_0}}{1-\overline{\alpha}_t}}_{\boldsymbol{\mu}_q(\boldsymbol{x}_t,t)}, \underbrace{\frac{(1-\alpha_t)(1-\overline{\alpha}_{t-1})}{1-\overline{\alpha}_t}\boldsymbol{I}}_{\Sigma_q(t)}\right)$$

$$\frac{\sqrt{\alpha_t}(1-\overline{\alpha}_{t-1})\boldsymbol{x}_t + \sqrt{\overline{\alpha}_{t-1}}(1-\alpha_t)\boldsymbol{x}_0}{1-\overline{\alpha}_t} \implies \frac{1}{\sqrt{\alpha_t}}\left(\boldsymbol{x}_t - \frac{1-\alpha_t}{\sqrt{1-\overline{\alpha}_t}}\boldsymbol{\varepsilon}_t\right)$$

**④ training**

$$\underset{\boldsymbol{\theta}}{\arg\min} \frac{1}{2\sigma_q^2(t)}[\|\boldsymbol{\mu}_\theta - \boldsymbol{\mu}_q\|_2^2]$$

$$\underset{\boldsymbol{\theta}}{\arg\min} \frac{1}{2\sigma_q^2(t)}\frac{(1-\alpha_t)^2}{(1-\overline{\alpha}_t)\alpha_t}[\|\boldsymbol{\varepsilon}_t - \hat{\boldsymbol{\varepsilon}}_\theta(\boldsymbol{x}_t,t)\|_2^2]$$

**⑤ sampling**

$$\boldsymbol{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}}\left(\boldsymbol{x}_t - \frac{1-\alpha_t}{\sqrt{1-\overline{\alpha}_t}}\boldsymbol{\varepsilon}_\theta(\boldsymbol{x}_t,t)\right) + \sigma_t\boldsymbol{z}$$

[1] Luo C. Understanding diffusion models: A unified perspective. arXiv, 2022.

# Training and Sampling

**Algorithm 1** Training

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

**Algorithm 2** Sampling

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

④ **training**

$$\underset{\boldsymbol{\theta}}{\text{argmin}} \, \frac{1}{2\sigma_q^2(t)} [\|\boldsymbol{\mu}_\theta - \boldsymbol{\mu}_q\|_2^2]$$

$$\underset{\boldsymbol{\theta}}{\text{argmin}} \, \boxed{\frac{1}{2\sigma_q^2(t)} \frac{(1-\alpha_t)^2}{(1-\bar{\alpha}_t)\alpha_t}} [\|\boldsymbol{\varepsilon}_t - \hat{\boldsymbol{\varepsilon}}_\theta(\boldsymbol{x}_t, t)\|_2^2]$$

simplify

⑤ **sampling**

$$\boldsymbol{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \boldsymbol{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\varepsilon}_\theta(\boldsymbol{x}_t, t) \right) + \sigma_t \boldsymbol{z}$$

[1] Luo C. Understanding diffusion models: A unified perspective. arXiv, 2022.

# The variance

$$q(\boldsymbol{x}_{t-1}|\boldsymbol{x}_t,\boldsymbol{x}_0) = \mathcal{N}\left(\boldsymbol{x}_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1-\overline{\alpha}_{t-1})\boldsymbol{x}_t + \sqrt{\overline{\alpha}_{t-1}}(1-\alpha_t)\boxed{\boldsymbol{x}_0}}{1-\overline{\alpha}_t}}_{\boldsymbol{\mu}_q(\boldsymbol{x}_t,t)}, \underbrace{\frac{(1-\alpha_t)(1-\overline{\alpha}_{t-1})}{1-\overline{\alpha}_t}\boldsymbol{I}}_{\Sigma_q(t)}\right)$$

$$\frac{(1-\alpha_t)(1-\overline{\alpha}_{t-1})}{1-\overline{\alpha}_t} = \frac{\overline{\beta}_{t-1}}{\overline{\beta}_t}\beta_t$$

**Fixed-small**   $\sigma_t^2 = \dfrac{\overline{\beta}_{t-1}}{\overline{\beta}_t}\beta_t = \tilde{\beta}_t$   [GLIGEN]

**Fixed-large**   $\sigma_t^2 = \beta_t$   [DDPM]

**hybrid**   $\sigma_t^2 = \exp\left(v\log\beta_t + (1-v)\log\tilde{\beta}_t\right)$   [IDDPM]
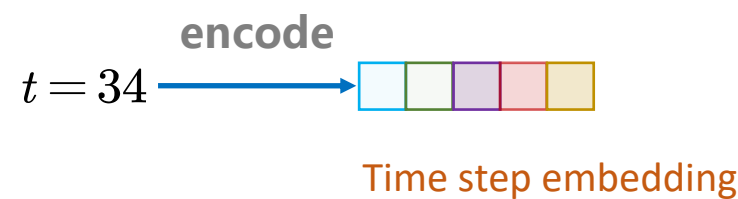
**optimal**   [Analytic DPM]

[1] Nichol A Q, Dhariwal P. Improved denoising diffusion probabilistic models. ICML, 2021.
[2] Li Y, Liu H, Wu Q, et al. Gligen: Open-set grounded text-to-image generation. CVPR, 2023.
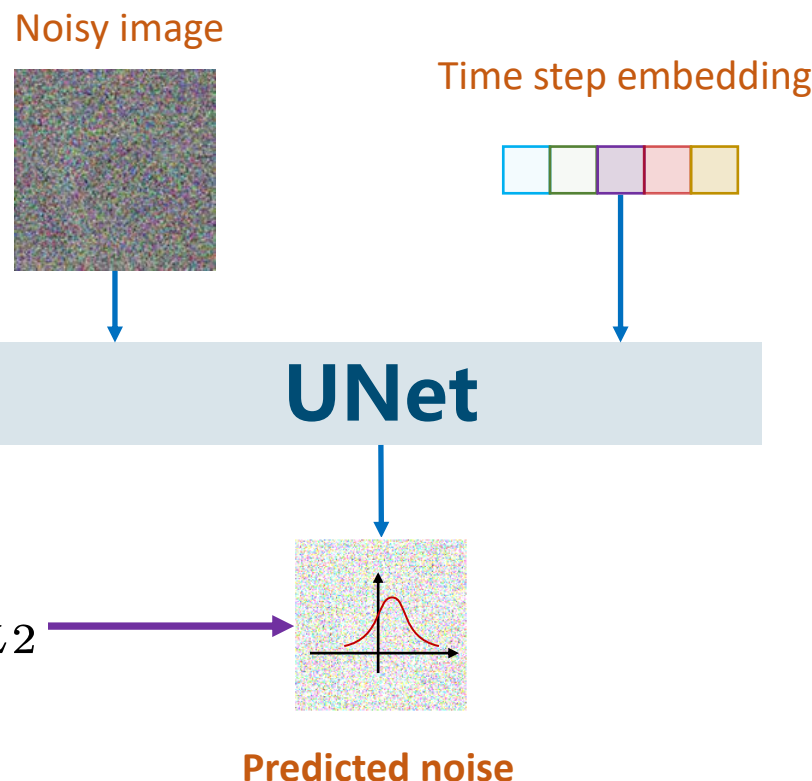
# Illustration of training

**1. Randomly select a time step and encode it.**

**encode**

$t = 34$ ⟶

Time step embedding

**2. Add noise to image.**

Noisy image $= \sqrt{\overline{\alpha}_t}$ [Original image] $+ \sqrt{1 - \overline{\alpha}_t}$ [Gaussian noise] ⟵ $\mathcal{L}_{L2}$ ⟶ [Predicted noise]

Noisy image          Original image          Gaussian noise          **Predicted noise**

**3. Train the UNet.**

Noisy image

Time step embedding

**UNet**

# Illustration of sampling

**1. Iteratively denoise the image ($T = 1000$)**



$\boldsymbol{x}_{1000}$

**UNet**

$\boldsymbol{\varepsilon}_\theta(x_t, t)$

$t = 1000$

$\boldsymbol{x}_{999}$     $\boldsymbol{x}_{1000}$     $\boldsymbol{\varepsilon}_\theta(\boldsymbol{x}_t, t)$

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

**2. Iteratively denoise the image ($T = 999..2$)**

**3. Iteratively denoise the image ($T = 1$)**

$\boldsymbol{x}_1$

**UNet**

$\boldsymbol{\varepsilon}_\theta(x_t, t)$

$t = 1$

$\boldsymbol{x}_0$     $\boldsymbol{x}_1$     $\boldsymbol{\varepsilon}_\theta(\boldsymbol{x}_t, t)$

$$\boldsymbol{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \boldsymbol{x}_t - \frac{1-\alpha_t}{\sqrt{1-\overline{\alpha}_t}} \boldsymbol{\varepsilon}_\theta(\boldsymbol{x}_t, t) \right)$$

# Reference

1. Ho J, Jain A, Abbeel P. Denoising diffusion probabilistic models. NeuraIPS, 2020.

2. Luo C. Understanding diffusion models: A unified perspective. arXiv, 2022.

3. Nichol A Q, Dhariwal P. Improved denoising diffusion probabilistic models. ICML, 2021.

4. https://jalammar.github.io/illustrated-stable-diffusion/

5. https://lilianweng.github.io/posts/2021-07-11-diffusion-models/ (What are diffusion models? )

6. https://www.youtube.com/watch?v=ifCDXFdeaaM&t=210s (李宏毅, 【生成式AI】Diffusion Model 原理剖析 (1/4) (optional))

7. https://kexue.fm/archives/9119. (苏剑林, 生成扩散模型漫谈（一）：DDPM = 拆楼 + 建楼)

8. https://www.bilibili.com/video/BV19H4y1G73r (SY_007, 【较真系列】讲人话-Diffusion Model全解(原理+代码+公式))

9. https://www.bilibili.com/video/BV1b541197HX (deep_thoughts, 54、Probabilistic Diffusion Model概率扩散模型理论与完整PyTorch代码详细解读)

10. https://www.bilibili.com/video/BV1p24y1K7Pf (Nik_Li, 一个视频看懂扩散模型DDPM原理推导|AI绘画底层模型)

# The end

**非常感谢你能看到这，希望该课件对你有帮助，视频讲解版在B站。**

**课件中出现的是我家的猫咪的照片，她已经陪伴了我很多年了，感谢她的友情出镜。o(*￣▽￣*)ブ**

Thank you so much for seeing this, I hope the slide is helpful, the video explanation version is on Bilibili.

The picture on the slide is of my cat, who has been with me for many years now, thanks for her friendly appearance! o(*￣▽￣*)ブ

" **What I cannot create, I do not understand.** "

**——Richard Feynman**