

# Diffusion Model

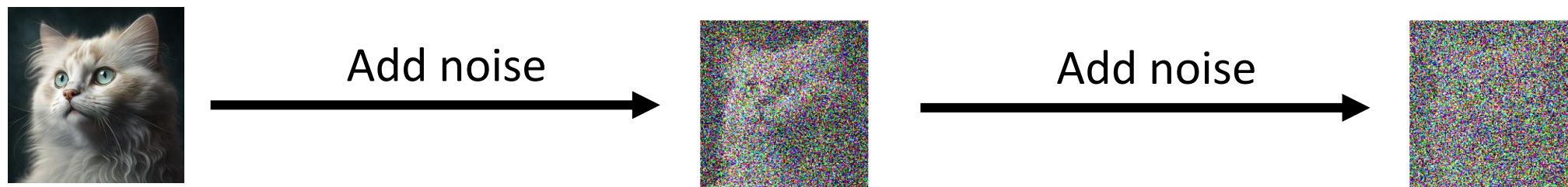
## 背後的數學原理

感謝姜成翰同學大力協助

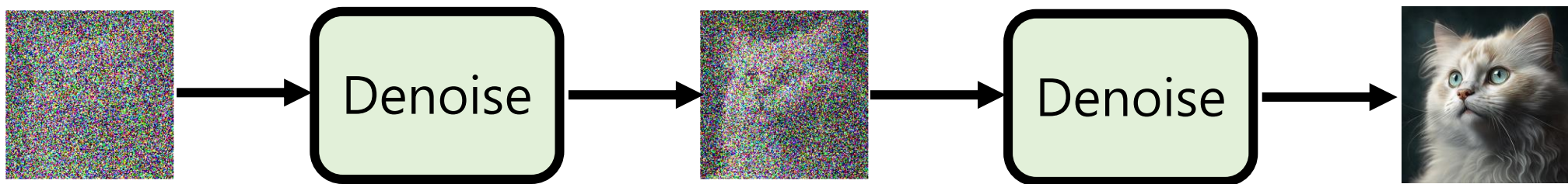
---

# 基本概念

## Forward Process

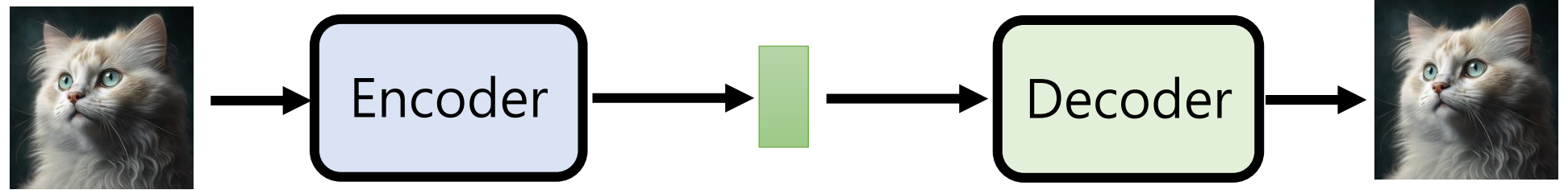


## Reverse Process

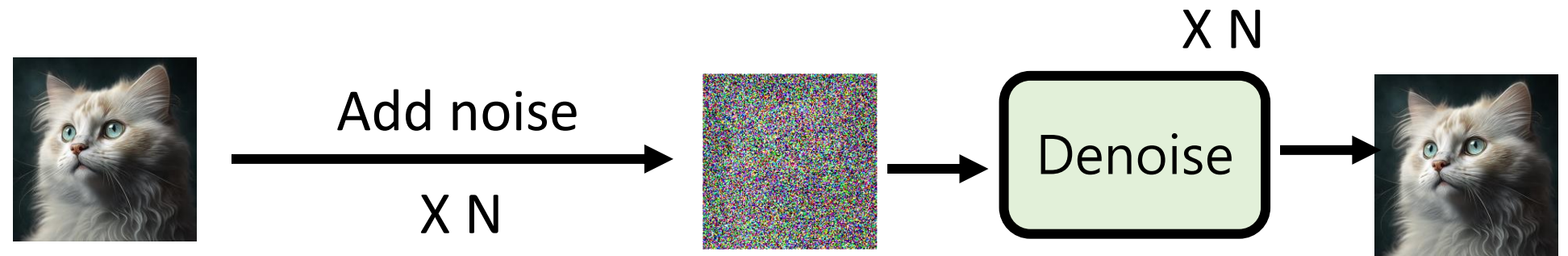


# VAE vs. Diffusion Model

VAE



Diffusion



# Denoising Diffusion Probabilistic Models

## 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

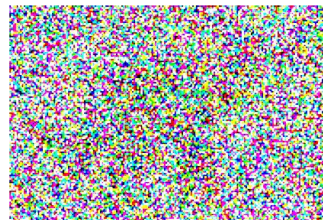
- 1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 2: **for**  $t = T, \dots, 1$  **do**
- 3:  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$
- 4:  $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1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

暗藏玄機!

# Training



$x_0$ : clean image



$\epsilon$ : noise

---

## Algorithm 1 Training

---

1: **repeat**

2:  $x_0 \sim q(x_0)$   $\leftarrow$  sample clean image

3:  $t \sim \text{Uniform}(\{1, \dots, T\})$

4:  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$   $\leftarrow$  sample a noise

5: Take gradient descent step on

$$\nabla_{\theta} \left\| \epsilon - \epsilon_{\theta}(\underbrace{\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon}_{\text{Noisy image}}, t) \right\|^2$$

6: **until** converged

Target  
Noise

Noisy image

Noise  
predictor

$\bar{\alpha}_1, \bar{\alpha}_2, \dots, \bar{\alpha}_T$

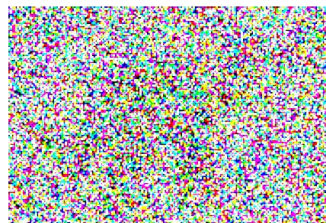
smaller

# Training

$\bar{\alpha}_1, \bar{\alpha}_2, \dots, \bar{\alpha}_T$



$x_0$



$\epsilon$

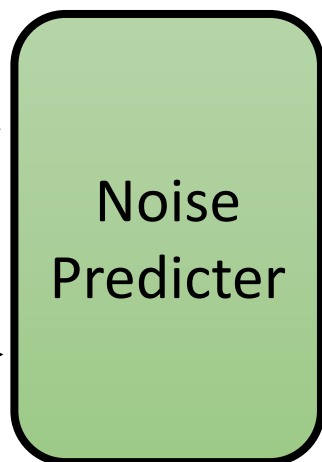
Sample  $t$

$$\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon = \text{Noisy Cat Image}$$

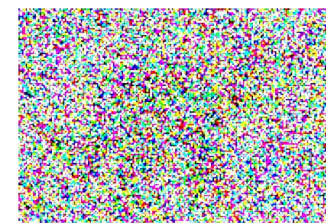
The equation shows the combination of a scaled clear image of the cat and a scaled noise image to produce a noisy version of the cat image.



$t$

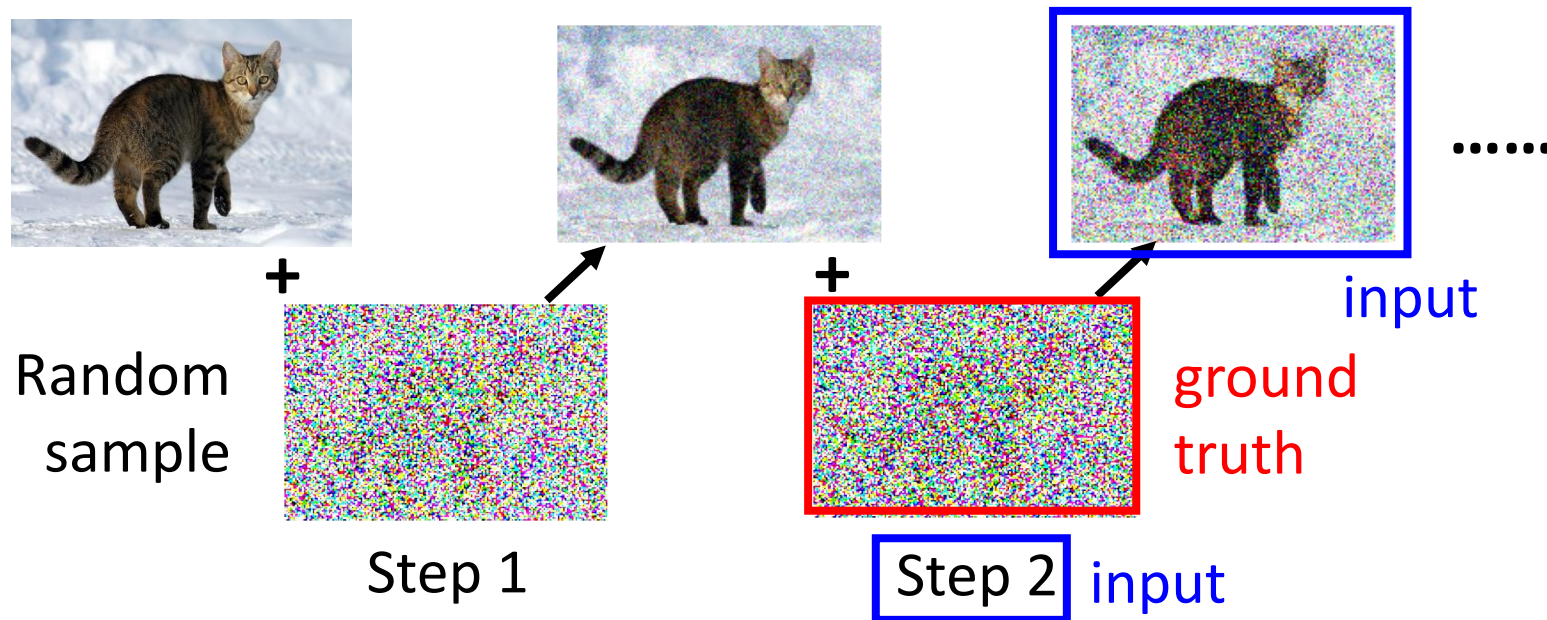


?????



$\epsilon$

想像中 ...

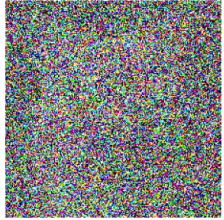


實際上 ...

$$\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \varepsilon = \text{input}$$

The equation shows the clear cat image  $x_0$  multiplied by  $\sqrt{\bar{\alpha}_t}$  (underlined in blue), plus the random noise  $\varepsilon$  (enclosed in a red box) multiplied by  $\sqrt{1 - \bar{\alpha}_t}$  (underlined in blue), equals the noisy "input" image (enclosed in a blue box). The text "ground truth" (in red) is positioned below the noise  $\varepsilon$ .

# Inference



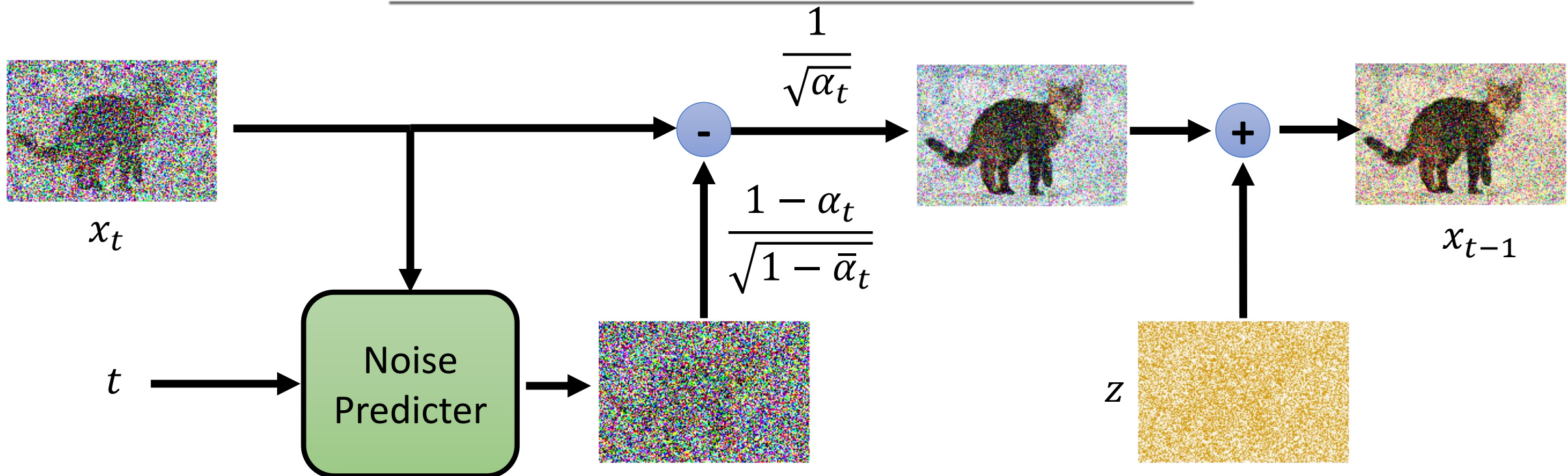
$x_T$

## Algorithm 2 Sampling

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

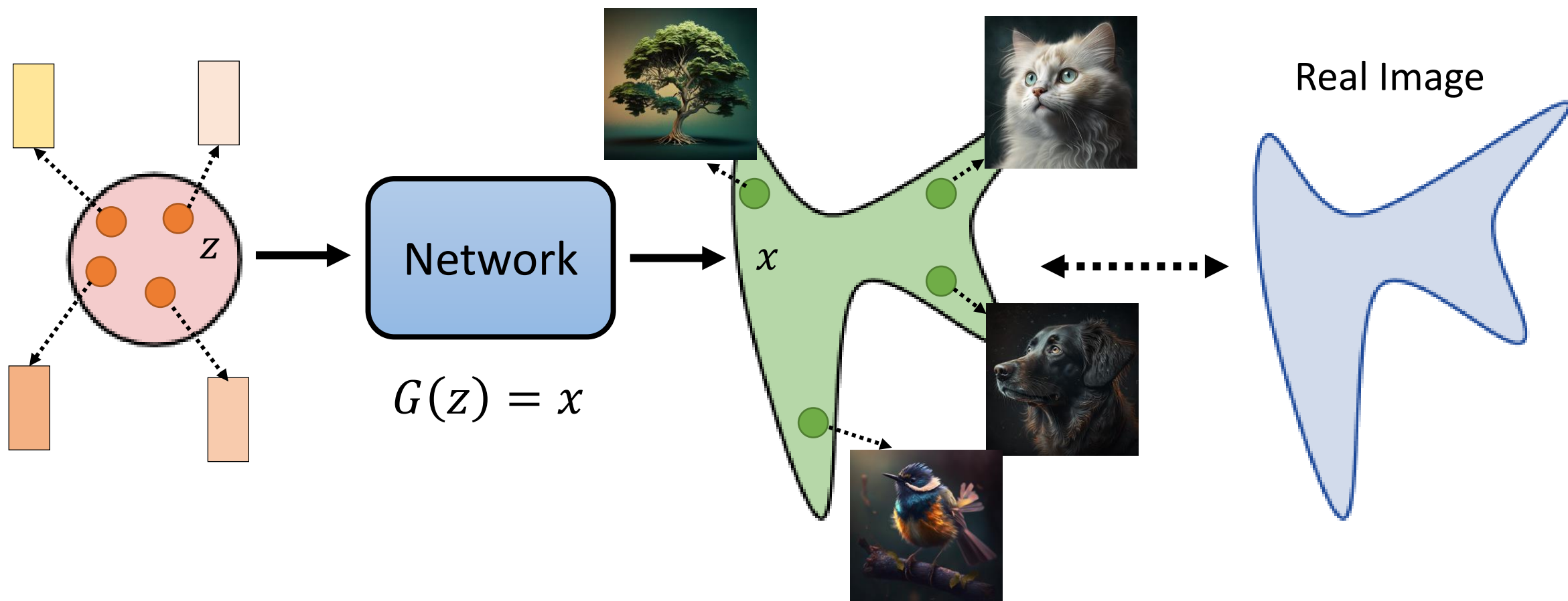
sample a noise?!

$\bar{\alpha}_1, \bar{\alpha}_2, \dots, \bar{\alpha}_T$   
 $\alpha_1, \alpha_2, \dots, \alpha_T$

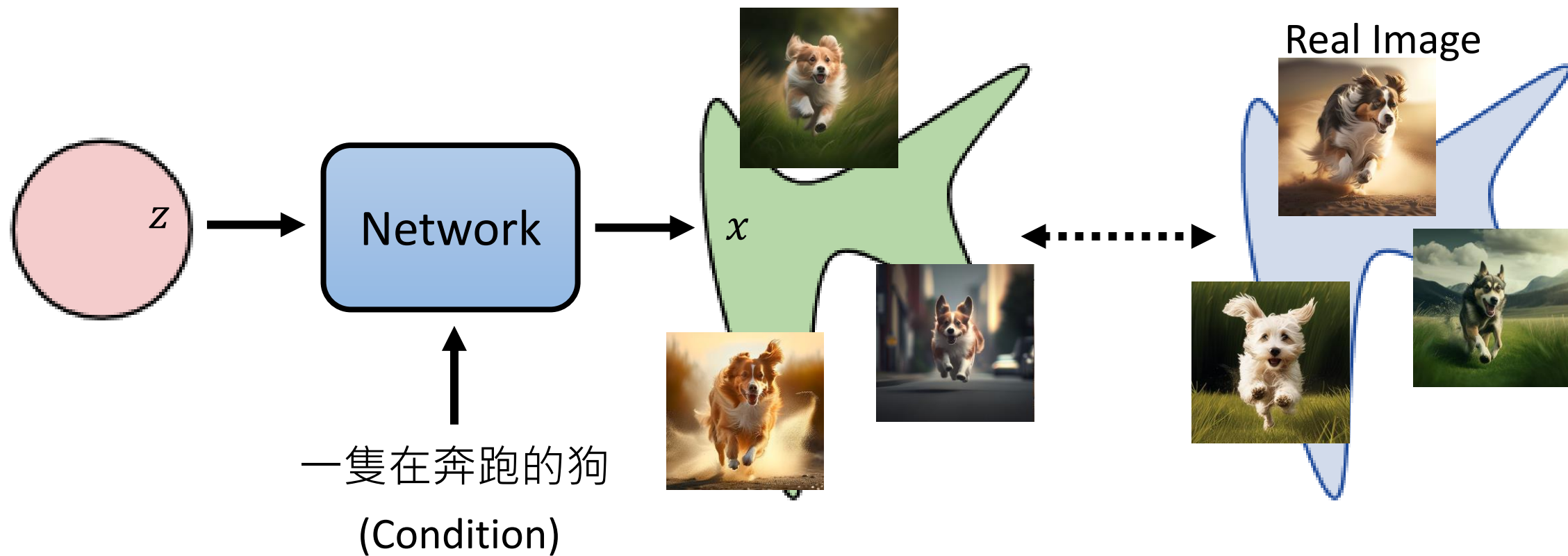




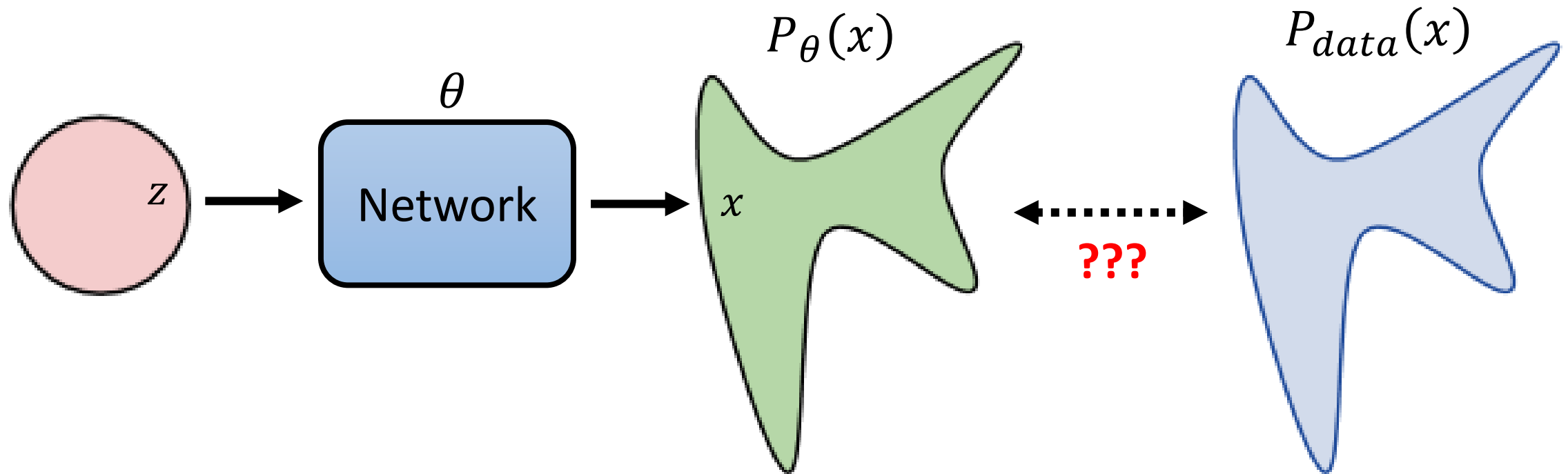
# 影像生成模型本質上的共同目標



# 影像生成模型本質上的共同目標



# Maximum Likelihood Estimation



Sample  $\{x^1, x^2, \dots, x^m\}$  from  $P_{data}(x)$

We can compute  $P_\theta(x^i)$

???

$$\theta^* = \arg \max_{\theta} \prod_{i=1}^m P_\theta(x^i)$$

Sample  $\{x^1, x^2, \dots, x^m\}$  from  $P_{data}(x)$

$$\theta^* = \arg \max_{\theta} \prod_{i=1}^m P_{\theta}(x^i) = \arg \max_{\theta} \log \prod_{i=1}^m P_{\theta}(x^i)$$

$$= \arg \max_{\theta} \sum_{i=1}^m \log P_{\theta}(x^i) \approx \arg \max_{\theta} E_{x \sim P_{data}} [\log P_{\theta}(x)]$$

$$= \arg \max_{\theta} \int_x P_{data}(x) \log P_{\theta}(x) dx - \int_x P_{data}(x) \log P_{data}(x) dx$$

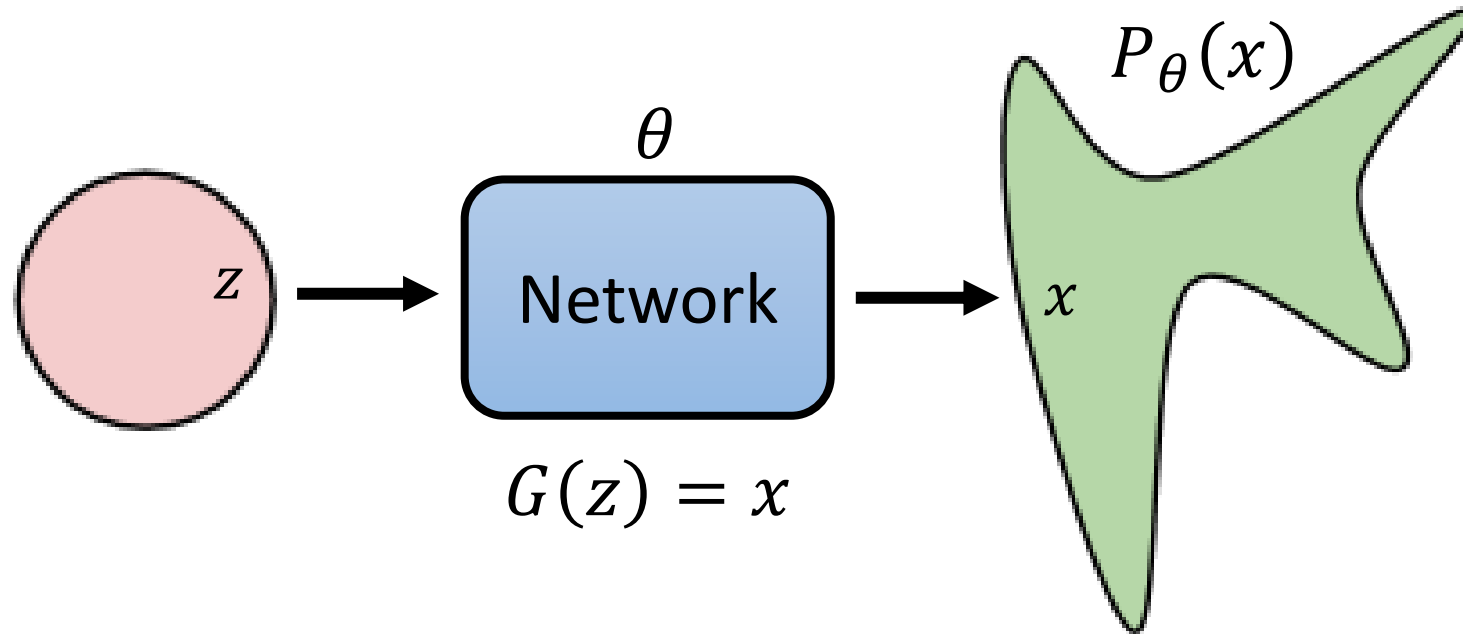
(not related to  $\theta$ )

$$= \arg \max_{\theta} \int_x P_{data}(x) \log \frac{P_{\theta}(x)}{P_{data}(x)} dx = \arg \min_{\theta} KL(P_{data} || P_{\theta})$$

Difference between  $P_{data}$  and  $P_{\theta}$

Maximum Likelihood = Minimize KL Divergence

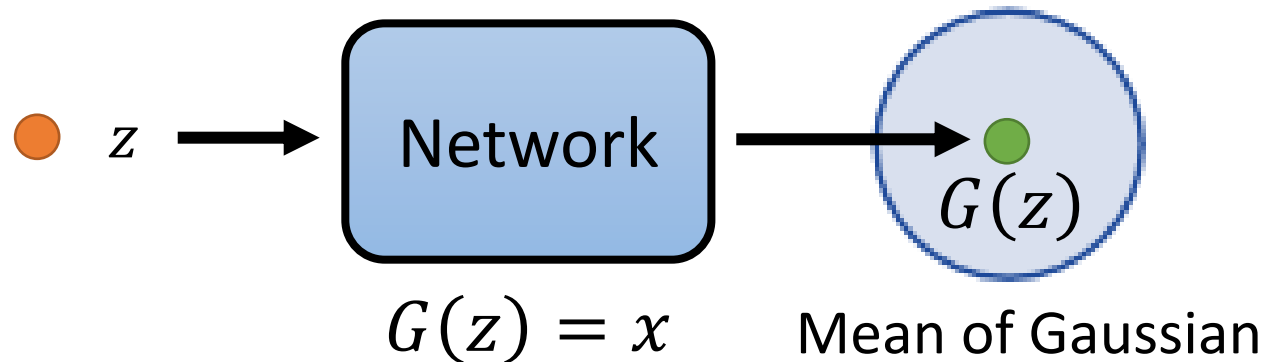
# VAE: Compute $P_{\theta}(x)$



$$P_{\theta}(x) = \int_z P(z)P_{\theta}(x|z)dz$$

$$P_{\theta}(x|z) = \begin{cases} 1, & G(z) = x \\ 0, & G(z) \neq x \end{cases}$$

可能會幾乎都是 0 ☹️



$$P_{\theta}(x|z) \propto \exp(-\|G(z) - x\|_2)$$

# VAE: Lower bound of $\log P(x)$

$$\log P_{\theta}(x) = \int_z q(z|x) \log P(x) dz$$

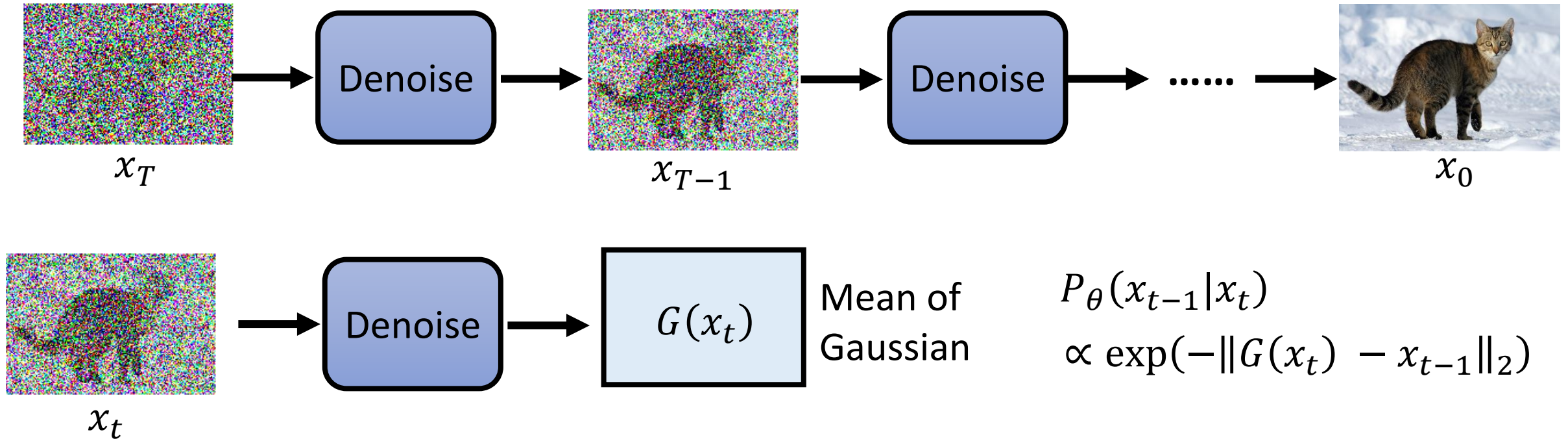
$q(z|x)$  can be any distribution

$$= \int_z q(z|x) \log \left( \frac{P(z, x)}{P(z|x)} \right) dz = \int_z q(z|x) \log \left( \frac{P(z, x) q(z|x)}{q(z|x) P(z|x)} \right) dz$$

$$= \int_z q(z|x) \log \left( \frac{P(z, x)}{q(z|x)} \right) dz + \underbrace{\int_z q(z|x) \log \left( \frac{q(z|x)}{P(z|x)} \right) dz}_{KL(q(z|x) || P(z|x))} \geq 0$$

$$\geq \int_z q(z|x) \log \left( \frac{P(z, x)}{q(z|x)} \right) dz = \underbrace{E_{q(z|x)} \left[ \log \left( \frac{P(x, z)}{q(z|x)} \right) \right]}_{\text{Encoder}} \quad \text{lower bound}$$

# DDPM: Compute $P_\theta(x)$



$$P_\theta(x_0) = \int_{x_1:x_T} P(x_T) P_\theta(x_{T-1}|x_T) \dots P_\theta(x_{t-1}|x_t) \dots P_\theta(x_0|x_1) dx_1:x_T$$

# DDPM: Lower bound of $\log P(x)$

**VAE**      Maximize  $\log P_\theta(\underline{x})$        $\longrightarrow$       Maximize  $\mathbb{E}_{\underline{q(z|x)}}$   $\left[ \log \left( \frac{P(\underline{x}, z)}{\underline{q(z|x)}} \right) \right]$

Encoder

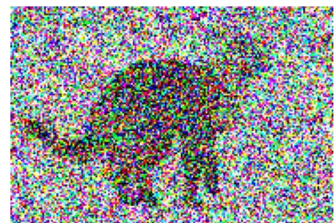
**DDPM**      Maximize  $\log P_\theta(\underline{x_0})$        $\longrightarrow$       Maximize  $\mathbb{E}_{\underline{q(x_1:x_T|x_0)}}$   $\left[ \log \left( \frac{P(\underline{x_0:x_T})}{\underline{q(x_1:x_T|x_0)}} \right) \right]$

Forward Process  
(Diffusion Process)

$$q(x_1:x_T|x_0) = q(x_1|x_0)q(x_2|x_1) \dots q(x_T|x_{T-1})$$



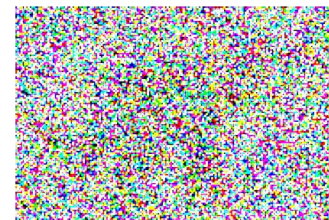
$$q(x_t|x_{t-1})$$

 $x_{t-1}$ 

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

 $x_t$ 

$$+ \sqrt{\beta_t}$$



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

 $\beta_1, \beta_2, \dots, \beta_T$ 

$$q(x_t|x_0)$$

 $x_0$ 

+



+



.....



+

 $x_t$



$x_1$

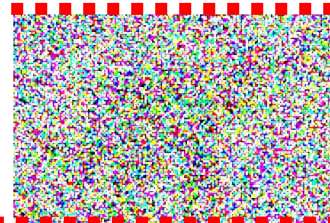
=

$$\sqrt{1 - \beta_1}$$



$x_0$

$$+ \sqrt{\beta_1}$$

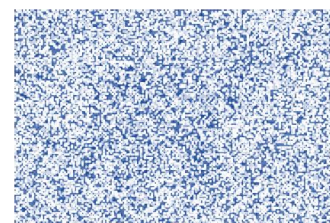


$$\sqrt{1 - \beta_2}$$



$x_1$

$$+ \sqrt{\beta_2}$$



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

Ind.

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



$x_2$

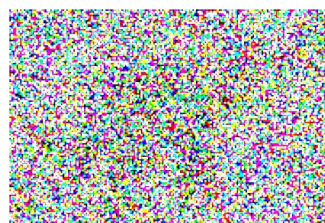
=

$$\sqrt{1 - \beta_2} \sqrt{1 - \beta_1}$$

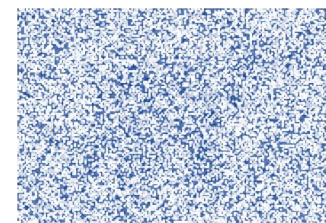


$x_0$

$$+ \sqrt{1 - \beta_2} \sqrt{\beta_1}$$



$$+ \sqrt{\beta_2}$$





$x_2$

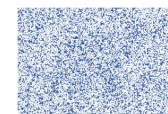
$$= \sqrt{1 - \beta_2} \sqrt{1 - \beta_1}$$



$x_0$



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



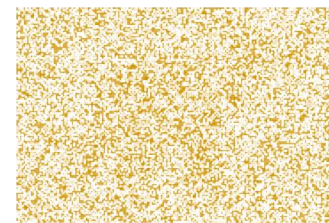
$\sim \mathcal{N}(\mathbf{0}, I)$

$$+ \sqrt{1 - \beta_2} \sqrt{\beta_1} \text{ [noisy image] } + \sqrt{\beta_2} \text{ [blue noisy image]}$$



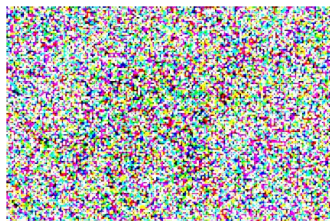
$\sim \mathcal{N}(\mathbf{0}, I)$

$$+ \sqrt{1 - (1 - \beta_2)(1 - \beta_1)}$$



$$q(x_t|x_0)$$

$$\beta_1, \beta_2, \dots, \beta_T$$



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

$$\alpha_t = 1 - \beta_t$$

$$\bar{\alpha}_t = \alpha_1 \alpha_2 \dots \alpha_t$$



⋮



⋮

$$= \sqrt{1 - \beta_1}$$

$$= \sqrt{1 - \beta_2}$$

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

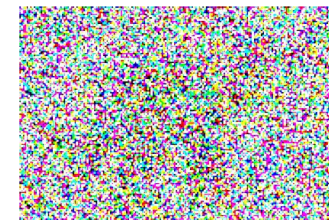
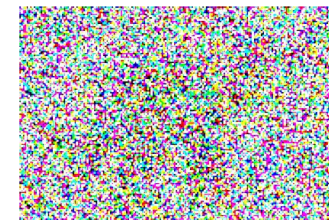
$$+ \sqrt{\beta_1}$$

$$+ \sqrt{\beta_2}$$

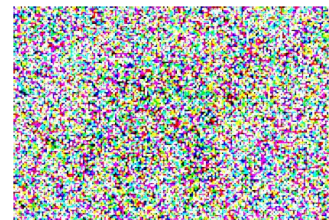
$$+ \sqrt{\beta_t}$$



⋮



⋮

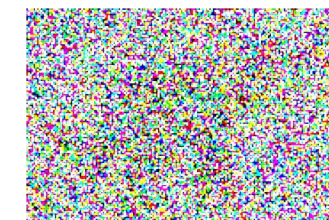


$$= \frac{\sqrt{1 - \beta_1} \dots \sqrt{1 - \beta_t}}{\sqrt{\bar{\alpha}_t}}$$



+

$$\frac{\sqrt{1 - (1 - \beta_1) \dots (1 - \beta_t)}}{\sqrt{1 - \bar{\alpha}_t}}$$



$$\log p(\mathbf{x}) \geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \quad (47)$$

$$= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right]$$

$$= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1) \prod_{t=2}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_1|\mathbf{x}_0) \prod_{t=2}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right]$$

$$= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1) \prod_{t=2}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_1|\mathbf{x}_0) \prod_{t=2}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)} \right]$$

$$= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ \log \frac{p_{\theta}(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)} \right] \quad (50)$$

$$= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}} \right] \quad (51)$$

$$= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}} \right] \quad (52)$$

$$= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{q(\mathbf{x}_T|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \quad (53)$$

$$= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_T|\mathbf{x}_0)} + \sum_{t=2}^T \log \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \quad (54)$$

$$= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_0)} \right] + \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ \log \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \quad (55)$$

$$= \mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_T|\mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_0)} \right] + \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t, \mathbf{x}_{t-1}|\mathbf{x}_0)} \left[ \log \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \quad (56)$$

$$= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{\text{prior matching term}} - \sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}} \quad (57)$$

Maximize  $\mathbb{E}_{q(\mathbf{x}_1:\mathbf{x}_T|\mathbf{x}_0)} \left[ \log \left( \frac{P(\mathbf{x}_0:\mathbf{x}_T)}{q(\mathbf{x}_1:\mathbf{x}_T|\mathbf{x}_0)} \right) \right]$

(50)

(51)

(52)

(53)

(54)

(55)

(56)

(57)

(58)

Understanding Diffusion Models:  
A Unified Perspective

<https://arxiv.org/pdf/2208.11970.pdf>

# DDPM: Lower bound of $\log P(x)$

$$\begin{aligned} \mathbb{E}_{q(x_1|x_0)}[\log P(x_0|x_1)] & -KL(q(x_T|x_0)||P(x_T)) \\ & - \sum_{t=2}^T \mathbb{E}_{q(x_t|x_0)}[KL(q(x_{t-1}|x_t, x_0)||P(x_{t-1}|x_t))] \end{aligned}$$

$$\mathbb{E}_{q(x_1|x_0)}[\log P(x_0|x_1)] - KL(q(x_T|x_0)||P(x_T))$$

$$- \sum_{t=2}^T \mathbb{E}_{q(x_t|x_0)} [KL(q(x_{t-1}|x_t, x_0)||P(x_{t-1}|x_t)))]$$

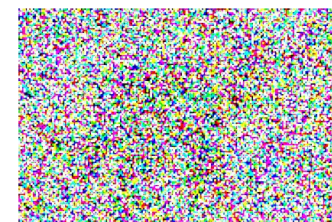
$$q(x_t|x_0)$$



$$= \sqrt{\bar{\alpha}_t}$$



$$+ \sqrt{1 - \bar{\alpha}_t}$$



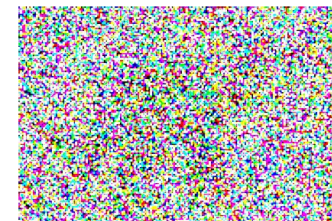
$$q(x_{t-1}|x_0)$$



$$= \sqrt{\bar{\alpha}_{t-1}}$$



$$+ \sqrt{1 - \bar{\alpha}_{t-1}}$$



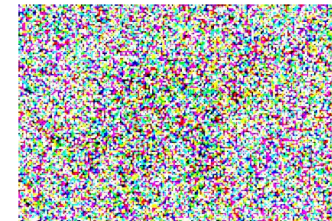
$$q(x_t|x_{t-1})$$



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



$$+ \sqrt{\beta_t}$$



$$E_{q(x_1|x_0)}[\log P(x_0|x_1)] - KL(q(x_T|x_0)||P(x_T))$$

$$- \sum_{t=2}^T E_{q(x_t|x_0)} [KL(q(x_{t-1}|x_t, x_0)||P(x_{t-1}|x_t)))]$$



$x_0$



$x_{t-1}$



$x_t$

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

$$q(x_t|x_0)$$

$$q(x_{t-1}|x_0)$$

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

已知 已知  
Gaussian Gaussian

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

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

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

已知  
Gaussian



$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)} \quad (71)$$

$$= \frac{\mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1-\alpha_t)\mathbf{I})\mathcal{N}(\mathbf{x}_{t-1}; \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0, (1-\bar{\alpha}_{t-1})\mathbf{I})}{\mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1-\bar{\alpha}_t)\mathbf{I})} \quad (72)$$

$$\propto \exp \left\{ - \left[ \frac{(\mathbf{x}_t - \sqrt{\alpha_t}\mathbf{x}_{t-1})^2}{2(1-\alpha_t)} + \frac{(\mathbf{x}_{t-1} - \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0)^2}{2(1-\bar{\alpha}_{t-1})} - \frac{(\mathbf{x}_t - \sqrt{\bar{\alpha}_t}\mathbf{x}_0)^2}{2(1-\bar{\alpha}_t)} \right] \right\} \quad (73)$$

$$= \exp \left\{ - \frac{1}{2} \left[ \frac{(\mathbf{x}_t - \sqrt{\alpha_t}\mathbf{x}_{t-1})^2}{1-\alpha_t} + \frac{(\mathbf{x}_{t-1} - \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0)^2}{1-\bar{\alpha}_{t-1}} - \frac{(\mathbf{x}_t - \sqrt{\bar{\alpha}_t}\mathbf{x}_0)^2}{1-\bar{\alpha}_t} \right] \right\} \quad (74)$$

$$= \exp \left\{ - \frac{1}{2} \left[ \frac{(-2\sqrt{\alpha_t}\mathbf{x}_t\mathbf{x}_{t-1} + \alpha_t\mathbf{x}_{t-1}^2)}{1-\alpha_t} + \frac{(\mathbf{x}_{t-1}^2 - 2\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_{t-1}\mathbf{x}_0)}{1-\bar{\alpha}_{t-1}} + C(\mathbf{x}_t, \mathbf{x}_0) \right] \right\} \quad (75)$$

$$\propto \exp \left\{ - \frac{1}{2} \left[ - \frac{2\sqrt{\alpha_t}\mathbf{x}_t\mathbf{x}_{t-1}}{1-\alpha_t} + \frac{\alpha_t\mathbf{x}_{t-1}^2}{1-\alpha_t} + \frac{\mathbf{x}_{t-1}^2}{1-\bar{\alpha}_{t-1}} - \frac{2\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_{t-1}\mathbf{x}_0}{1-\bar{\alpha}_{t-1}} \right] \right\} \quad (76)$$

$$= \exp \left\{ - \frac{1}{2} \left[ \left( \frac{\alpha_t}{1-\alpha_t} + \frac{1}{1-\bar{\alpha}_{t-1}} \right) \mathbf{x}_{t-1}^2 - 2 \left( \frac{\sqrt{\alpha_t}\mathbf{x}_t}{1-\alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0}{1-\bar{\alpha}_{t-1}} \right) \mathbf{x}_{t-1} \right] \right\} \quad (77)$$

$$= \exp \left\{ - \frac{1}{2} \left[ \frac{\alpha_t(1-\bar{\alpha}_{t-1}) + 1-\alpha_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} \mathbf{x}_{t-1}^2 - 2 \left( \frac{\sqrt{\alpha_t}\mathbf{x}_t}{1-\alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0}{1-\bar{\alpha}_{t-1}} \right) \mathbf{x}_{t-1} \right] \right\} \quad (78)$$

$$= \exp \left\{ - \frac{1}{2} \left[ \frac{\alpha_t - \bar{\alpha}_t + 1 - \alpha_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} \mathbf{x}_{t-1}^2 - 2 \left( \frac{\sqrt{\alpha_t}\mathbf{x}_t}{1-\alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0}{1-\bar{\alpha}_{t-1}} \right) \mathbf{x}_{t-1} \right] \right\} \quad (79)$$

$$= \exp \left\{ - \frac{1}{2} \left[ \frac{1 - \bar{\alpha}_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} \mathbf{x}_{t-1}^2 - 2 \left( \frac{\sqrt{\alpha_t}\mathbf{x}_t}{1-\alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0}{1-\bar{\alpha}_{t-1}} \right) \mathbf{x}_{t-1} \right] \right\} \quad (80)$$

$$= \exp \left\{ - \frac{1}{2} \left( \frac{1 - \bar{\alpha}_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} \right) \left[ \mathbf{x}_{t-1}^2 - 2 \frac{\left( \frac{\sqrt{\alpha_t}\mathbf{x}_t}{1-\alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0}{1-\bar{\alpha}_{t-1}} \right)}{\frac{1-\bar{\alpha}_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}} \mathbf{x}_{t-1} \right] \right\} \quad (81)$$

$$= \exp \left\{ - \frac{1}{2} \left( \frac{1 - \bar{\alpha}_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} \right) \left[ \mathbf{x}_{t-1}^2 - 2 \frac{\left( \frac{\sqrt{\alpha_t}\mathbf{x}_t}{1-\alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0}{1-\bar{\alpha}_{t-1}} \right) (1-\alpha_t)(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} \mathbf{x}_{t-1} \right] \right\} \quad (82)$$

$$= \exp \left\{ - \frac{1}{2} \left( \frac{1}{\frac{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}} \right) \left[ \mathbf{x}_{t-1}^2 - 2 \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1-\alpha_t)\mathbf{x}_0}{1-\bar{\alpha}_t} \mathbf{x}_{t-1} \right] \right\} \quad (83)$$

$$\propto \mathcal{N}(\mathbf{x}_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1-\alpha_t)\mathbf{x}_0}{1-\bar{\alpha}_t}}_{\mu_q(\mathbf{x}_t, \mathbf{x}_0)}, \underbrace{\frac{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}}_{\Sigma_q(t)} \mathbf{I}) \quad (84)$$

$$\mathbb{E}_{q(x_1|x_0)}[\log P(x_0|x_1)] - KL(q(x_T|x_0)||P(x_T))$$

$$- \sum_{t=2}^T \mathbb{E}_{q(x_t|x_0)} [KL(q(x_{t-1}|x_t, x_0)||P(x_{t-1}|x_t)))]$$



Gaussian



Mean

Variance

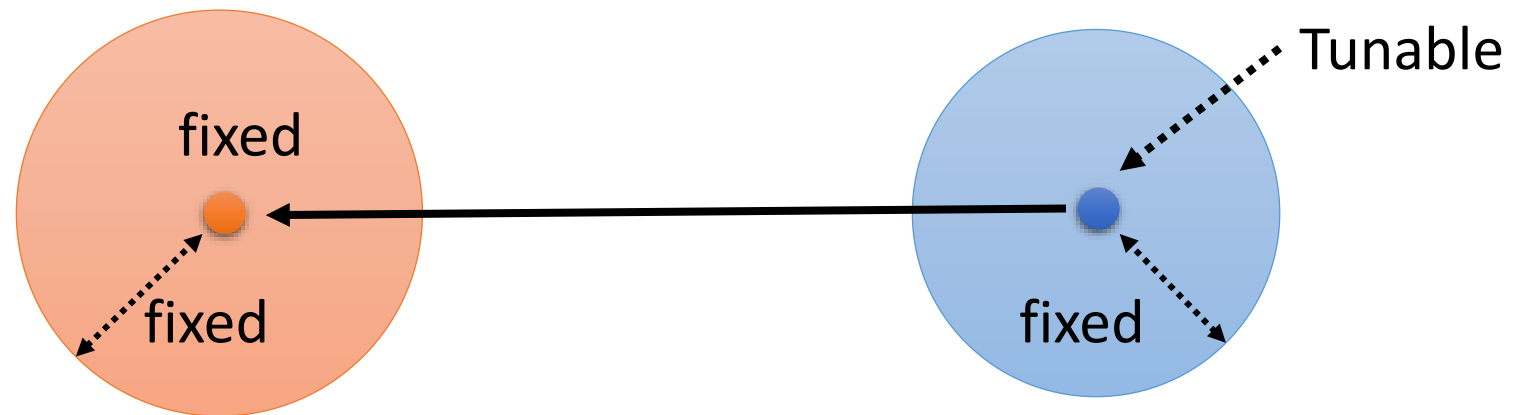
$$\frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t x_0 + \sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t}{1 - \bar{\alpha}_t}$$

$$\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t I$$

$$\mathbb{E}_{q(x_1|x_0)}[\log P(x_0|x_1)] - KL(q(x_T|x_0)||P(x_T))$$

$$- \sum_{t=2}^T \mathbb{E}_{q(x_t|x_0)} [KL(q(x_{t-1}|x_t, x_0)||P(x_{t-1}|x_t)))]$$

How to minimize  
KL divergence?



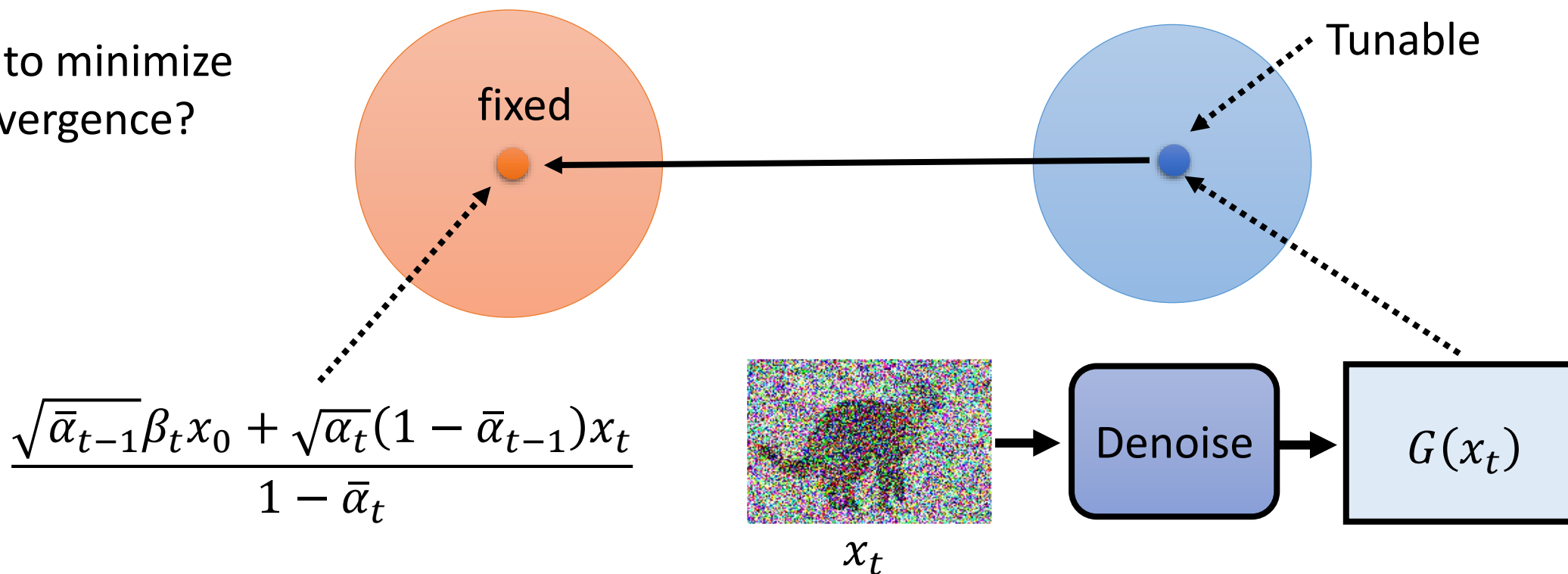
Recall that the KL Divergence between two Gaussian distributions is:

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

$$\mathbb{E}_{q(x_1|x_0)}[\log P(x_0|x_1)] - KL(q(x_T|x_0)||P(x_T))$$

$$- \sum_{t=2}^T \mathbb{E}_{q(x_t|x_0)} [KL(q(x_{t-1}|x_t, x_0)||P(x_{t-1}|x_t)))]$$

How to minimize  
KL divergence?



$$\mathbb{E}_{q(x_1|x_0)}[\log P(x_0|x_1)] - KL(q(x_T|x_0)||P(x_T))$$

$$- \sum_{t=2}^T \mathbb{E}_{q(x_t|x_0)}[KL(q(x_{t-1}|x_t, x_0)||P(x_{t-1}|x_t)))]$$



Sample  $x_0$



Sample  $x_t$



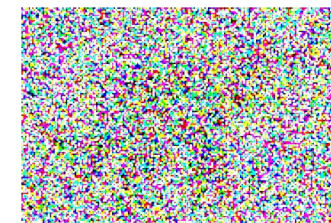
$x_t$

$$= \sqrt{\bar{\alpha}_t}$$



$x_0$

$$+ \sqrt{1 - \bar{\alpha}_t}$$



$\epsilon$

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon$$

$$\mathbb{E}_{q(x_1|x_0)}[\log P(x_0|x_1)] - KL(q(x_T|x_0)||P(x_T))$$

$$- \sum_{t=2}^T \mathbb{E}_{q(x_t|x_0)} [KL(q(x_{t-1}|x_t, x_0)||P(x_{t-1}|x_t)))]$$



$x_0$

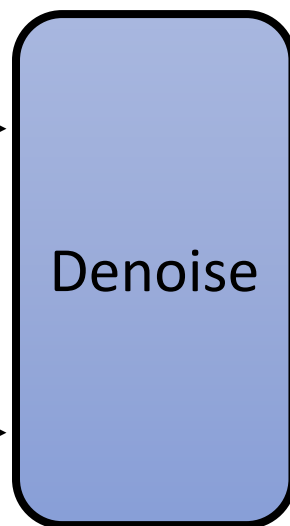


$x_t$



Sample  $x_t$

$t$



?



$$\frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t x_0 + \sqrt{\bar{\alpha}_t}(1 - \bar{\alpha}_{t-1})x_t}{1 - \bar{\alpha}_t}$$

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\varepsilon$$



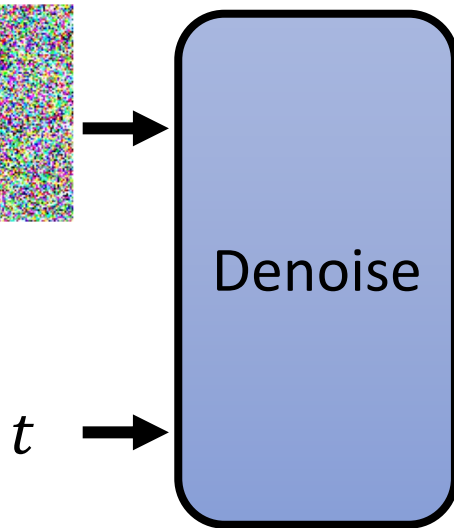
$x_0$



$x_t$



Sample  $x_t$



$$\frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t x_0 + \sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t}{1 - \bar{\alpha}_t}$$

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\varepsilon$$

$$x_t - \sqrt{1 - \bar{\alpha}_t}\varepsilon = \sqrt{\bar{\alpha}_t}x_0$$

$$\frac{x_t - \sqrt{1 - \bar{\alpha}_t}\varepsilon}{\sqrt{\bar{\alpha}_t}} = x_0$$

$$= \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t \frac{x_t - \sqrt{1 - \bar{\alpha}_t}\varepsilon}{\sqrt{\bar{\alpha}_t}} + \sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t}{1 - \bar{\alpha}_t}$$

$$= \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \varepsilon \right)$$



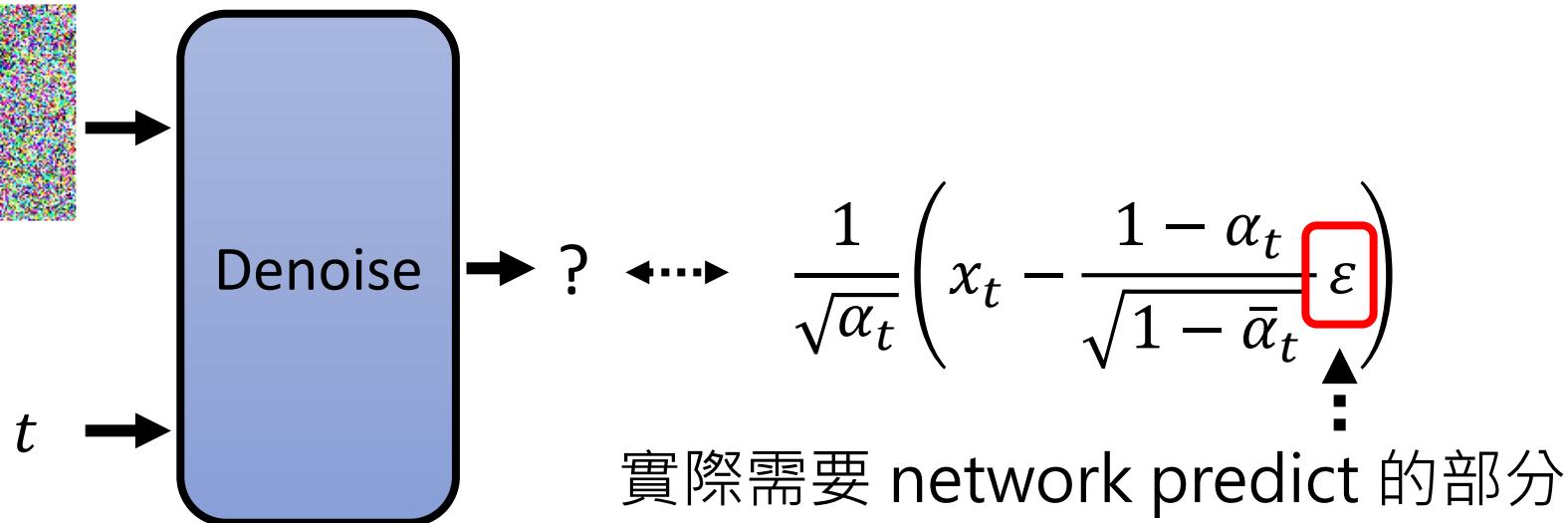
$x_0$



$x_t$



Sample  $x_t$




---

## Algorithm 2 Sampling

---

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\varepsilon$$

$$x_t - \sqrt{1 - \bar{\alpha}_t}\varepsilon = \sqrt{\bar{\alpha}_t}x_0$$

$$\frac{x_t - \sqrt{1 - \bar{\alpha}_t}\varepsilon}{\sqrt{\bar{\alpha}_t}} = x_0$$

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

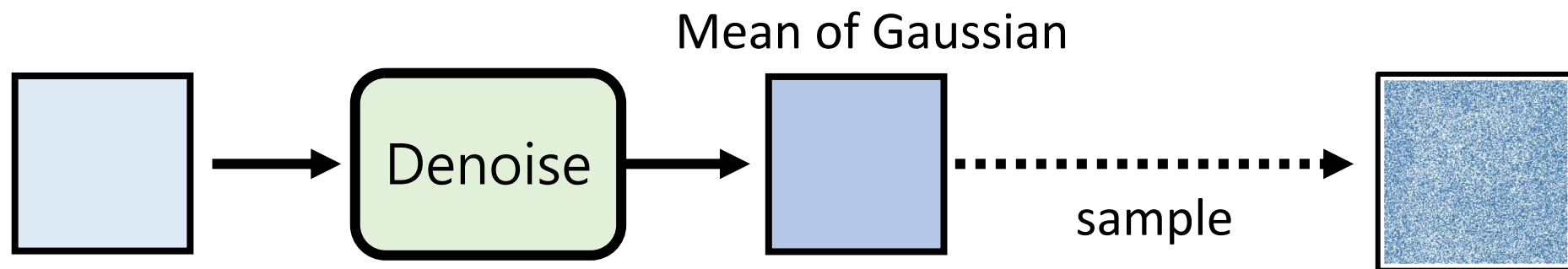


---

## Algorithm 2 Sampling

---

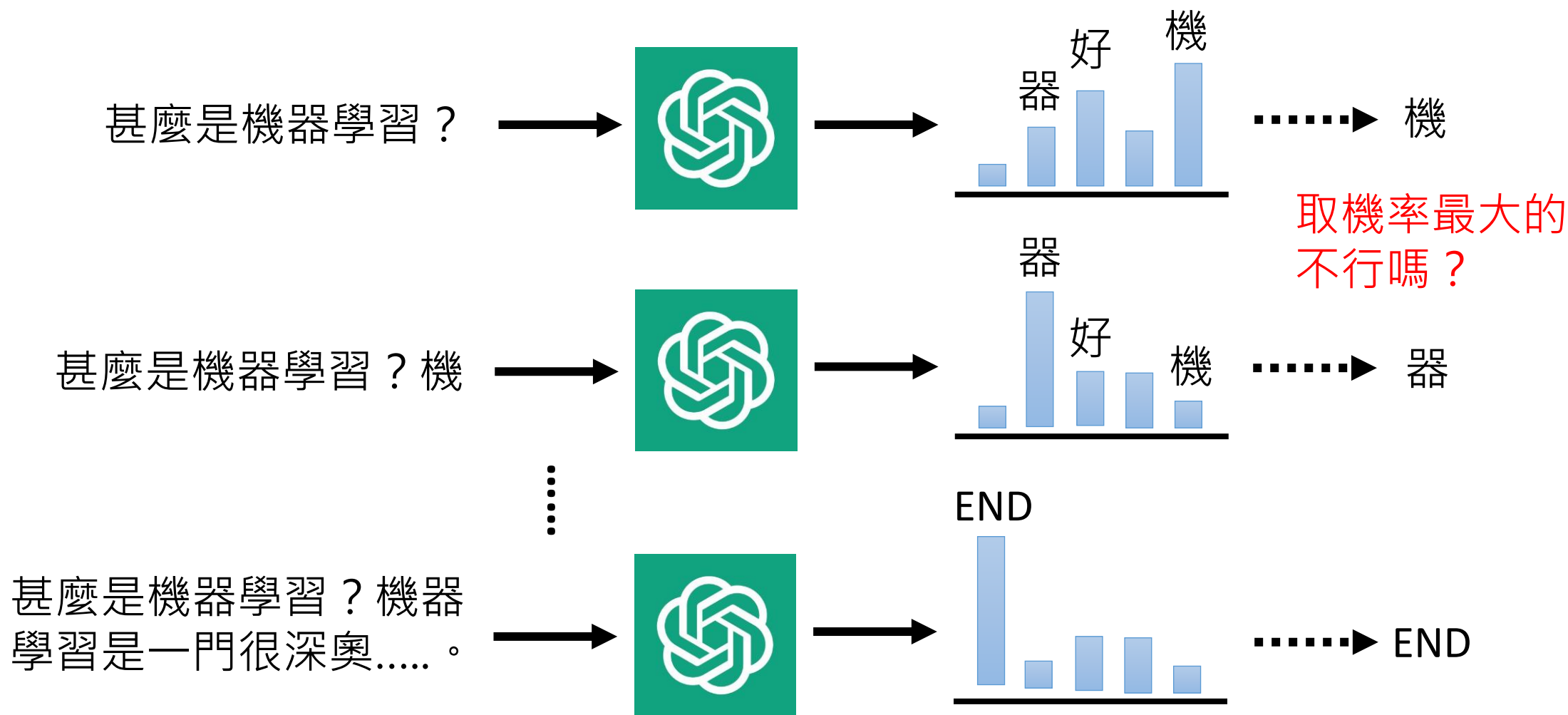
- 1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
  - 2: **for**  $t = T, \dots, 1$  **do**
  - 3:  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$
  - 4:  $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}}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



為什麼不直接取 Mean ?

免責聲明：以下只是猜測

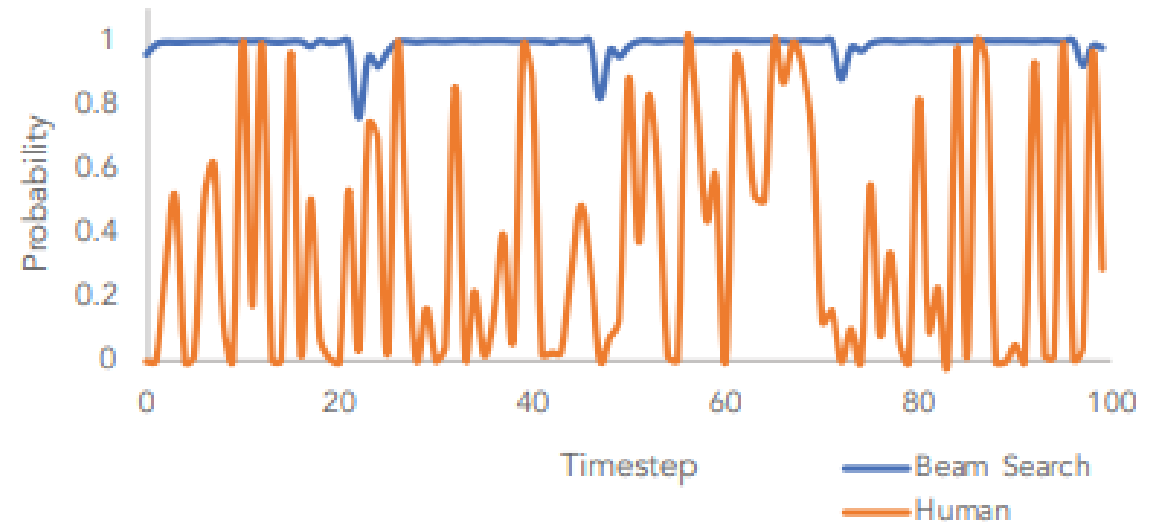
# 為什麼生成文句時需要 Sample ？





<https://arxiv.org/abs/1904.09751>

## Beam Search Text is Less Surprising



### Beam Search

...to provide an overview of the current state-of-the-art in the field of computer vision and machine learning, and to provide an overview of the current state-of-the-art in the field of computer vision and machine learning, and to provide an overview of the current state-of-the-art in the field of computer vision and machine learning, and to provide an overview of the current state-of-the-art in the field of computer vision and machine learning, and to provide an overview of the current state-of-the-art in the field of computer vision and machine learning, and...

### Human

...which grant increased life span and three years warranty. The Antec HCG series consists of five models with capacities spanning from 400W to 900W. Here we should note that we have already tested the HCG-620 in a previous review and were quite satisfied With its performance. In today's review we will rigorously test the Antec HCG-520, which as its model number implies, has 520W capacity and contrary to Antec's strong beliefs in multi-rail PSUs is equipped...

# 語音合成也需要 Sampling !

<https://arxiv.org/abs/1712.05884>

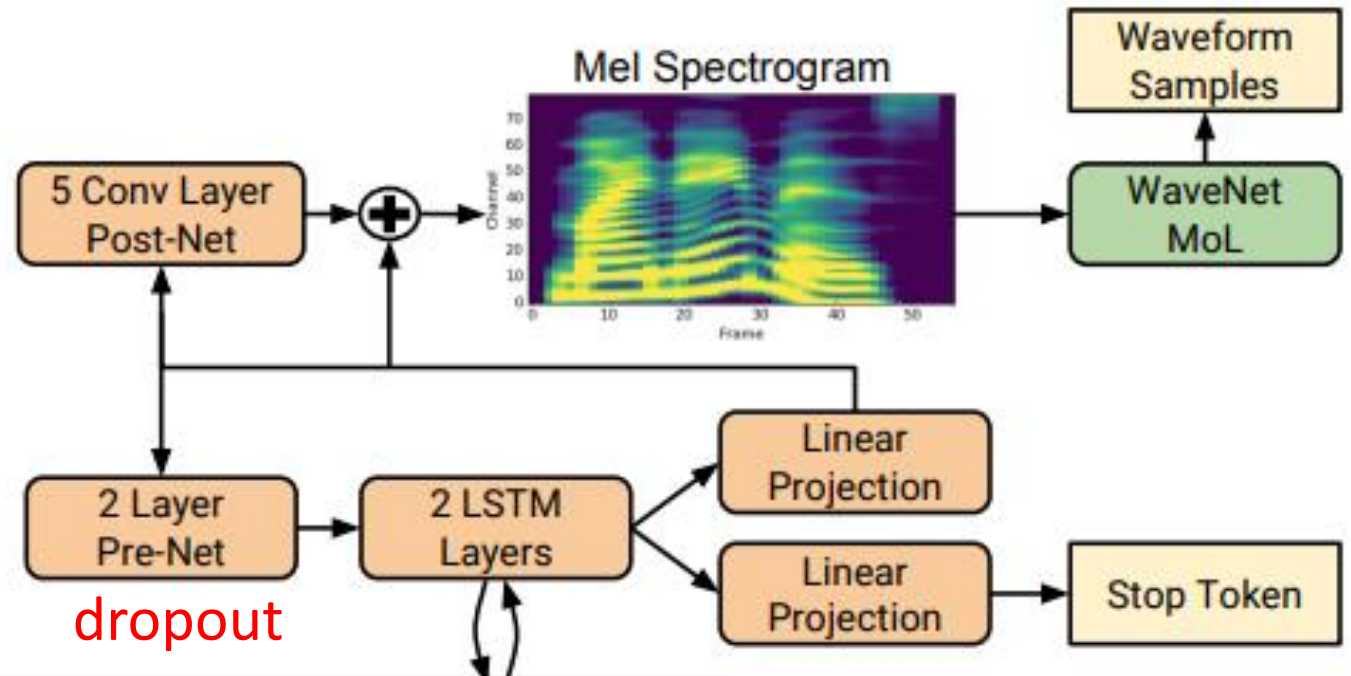


with  
dropout

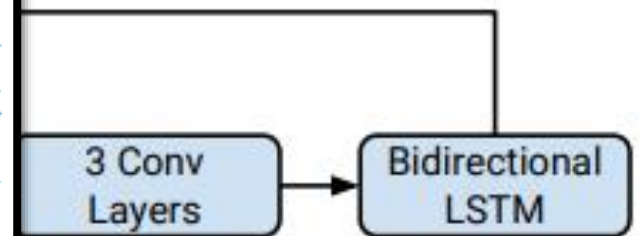


without  
dropout

感謝杜濤同學提供實驗結果

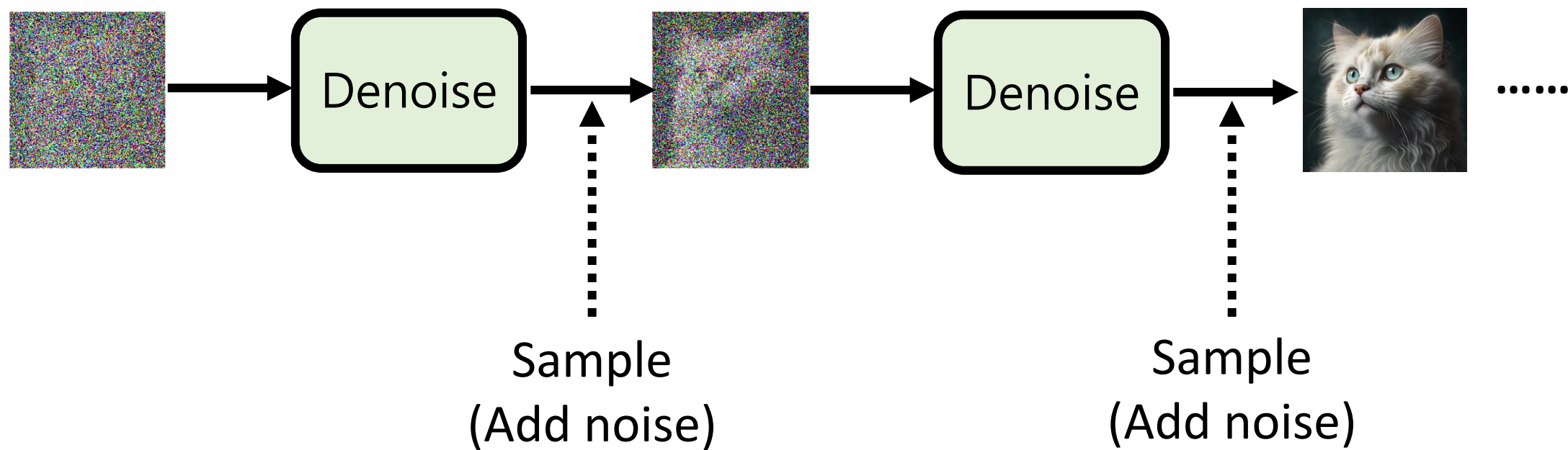


The convolutional layers in the network are regularized using dropout [25] with probability 0.5, and LSTM layers are regularized using zoneout [26] with probability 0.1. In order to introduce output variation at inference time, dropout with probability 0.5 is applied only to layers in the pre-net of the autoregressive decoder.



# Diffusion Model 是一種 Autoregressive

「一次到位」改成「N次到位」



---

## Algorithm 2 Sampling

---

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

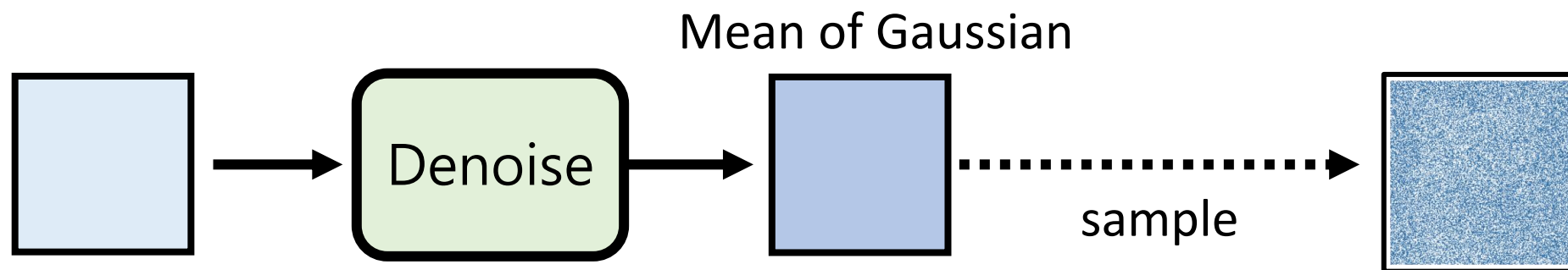
$\sigma_t$  as paper



$\sigma_t = 0$



感謝伏宇寬助  
教提供結果



為什麼不直接取 Mean ?



# Denoising Diffusion Probabilistic Models

---

## 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  
$$\nabla_{\theta} \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t) \right\|^2$$
  - 6: **until** converged
- 

---

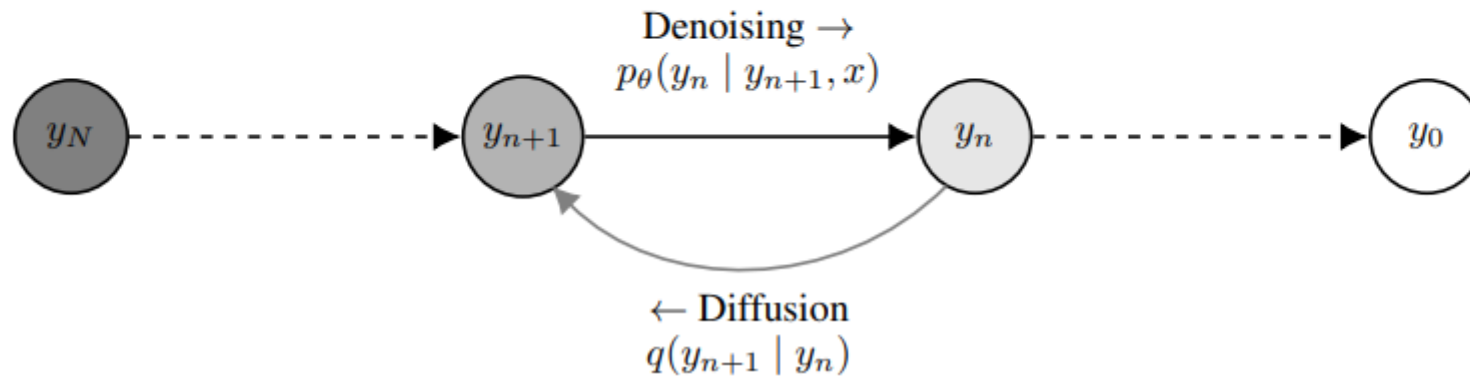
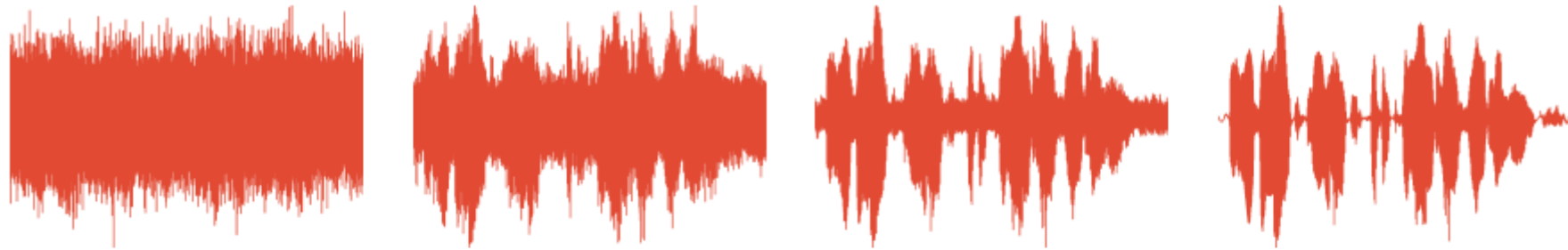
## Algorithm 2 Sampling

---

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

# Diffusion Model for Speech

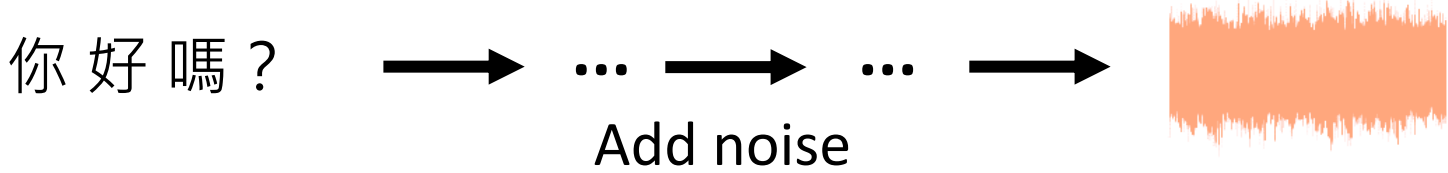
- WaveGrad



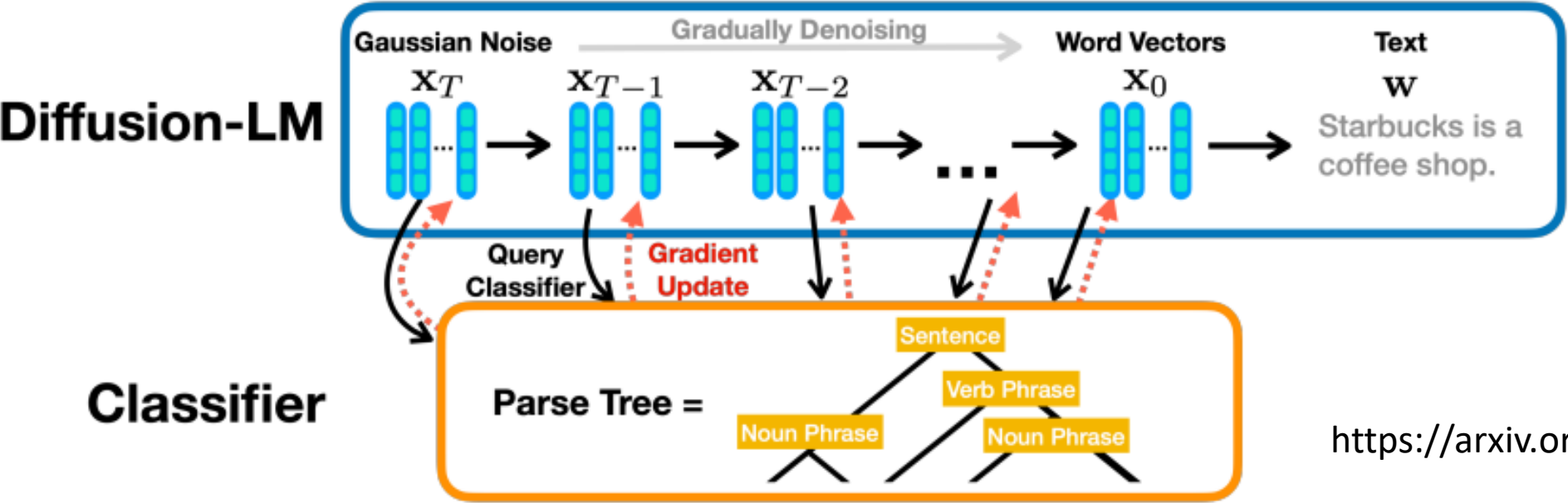
<https://arxiv.org/abs/2009.00713>

# Diffusion Model for Text

- Difficulty:



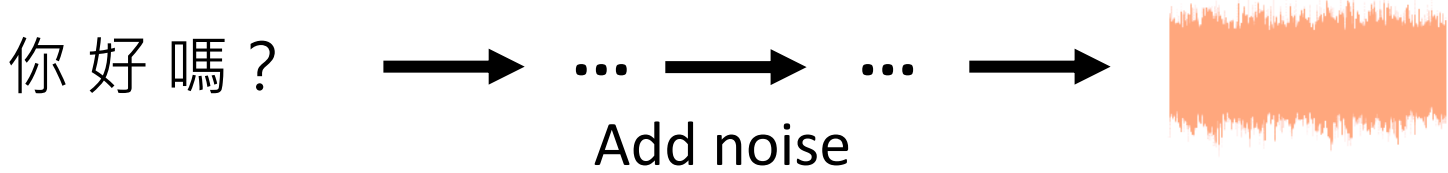
- Solution: Noise on latent space



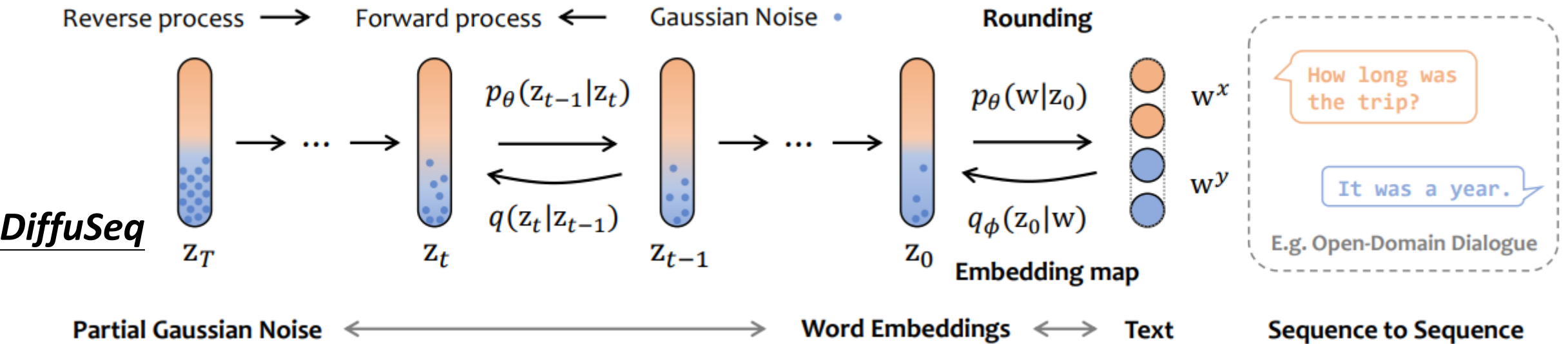
<https://arxiv.org/abs/2205.14217>

# Diffusion Model for Text

- Difficulty:



- Solution: Noise on latent space

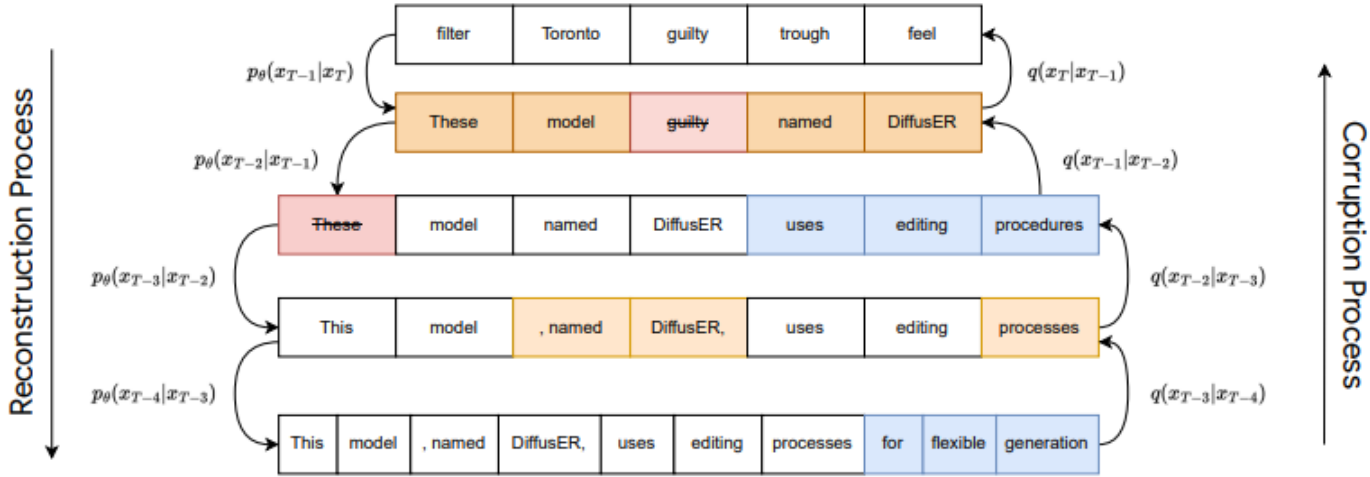


# Diffusion Model for Text

- Solution: Don't add Gaussian noise

<https://arxiv.org/abs/2210.16886>

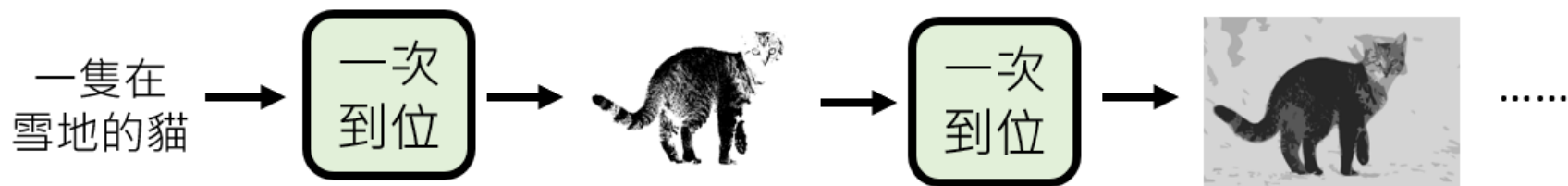
Diffusion via Edit-based Reconstruction (DiffusER)



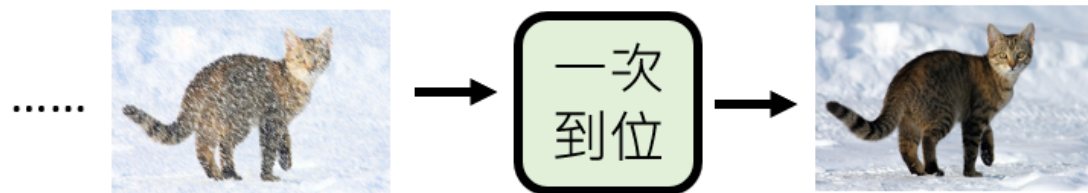
`t = 128 [MASK] [MASK] [MASK] [MASK] [MASK] [MASK]...`  
`t = 25 In response [MASK] the demands , [MASK] [MASK]y Workers union said [MASK] backflow fund [MASK]s would face further investigation and a fine.`  
`t = 0 In response to the demands , the Community Workers union said the backflow fund managers would face further investigation and a fine .`

<https://arxiv.org/abs/2107.03006>

# 各個擊破 + 一次到位



「一次到位」改成  
「N次到位」



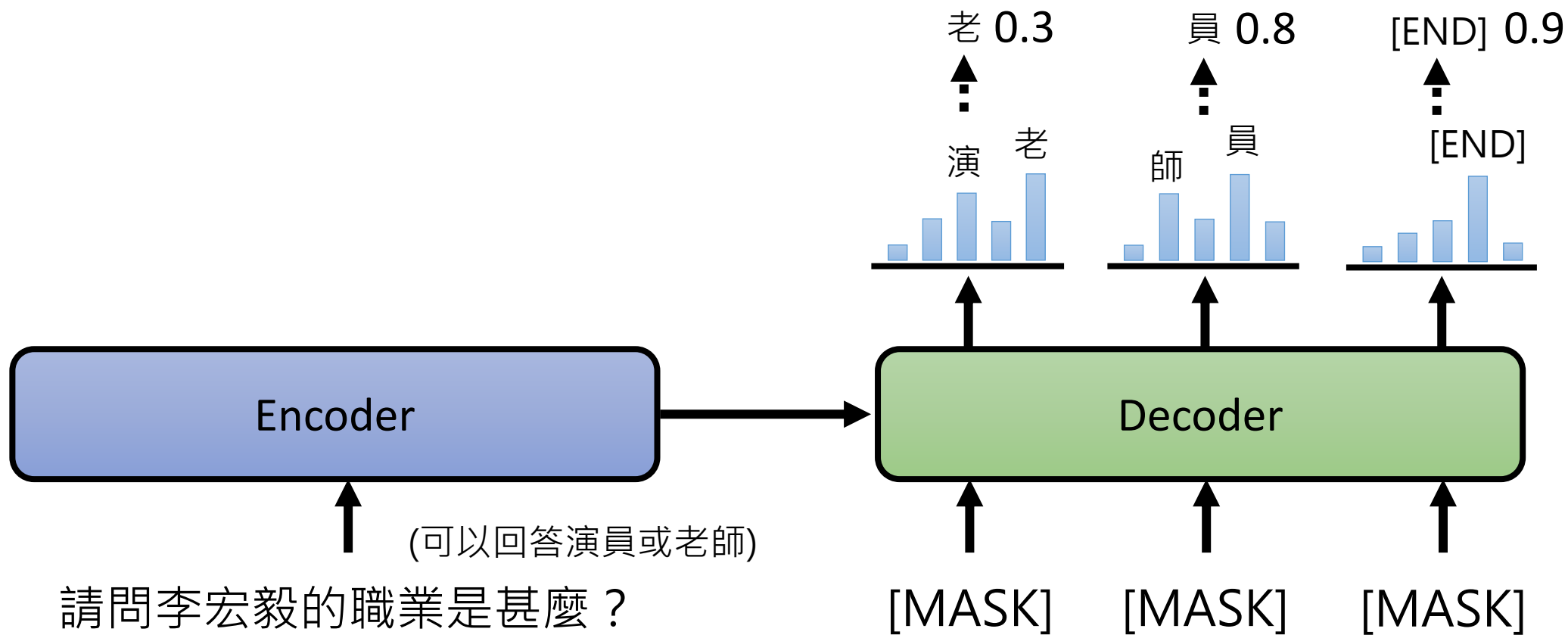
疑？這聽起來很像 Diffusion Model



之前上課投影片

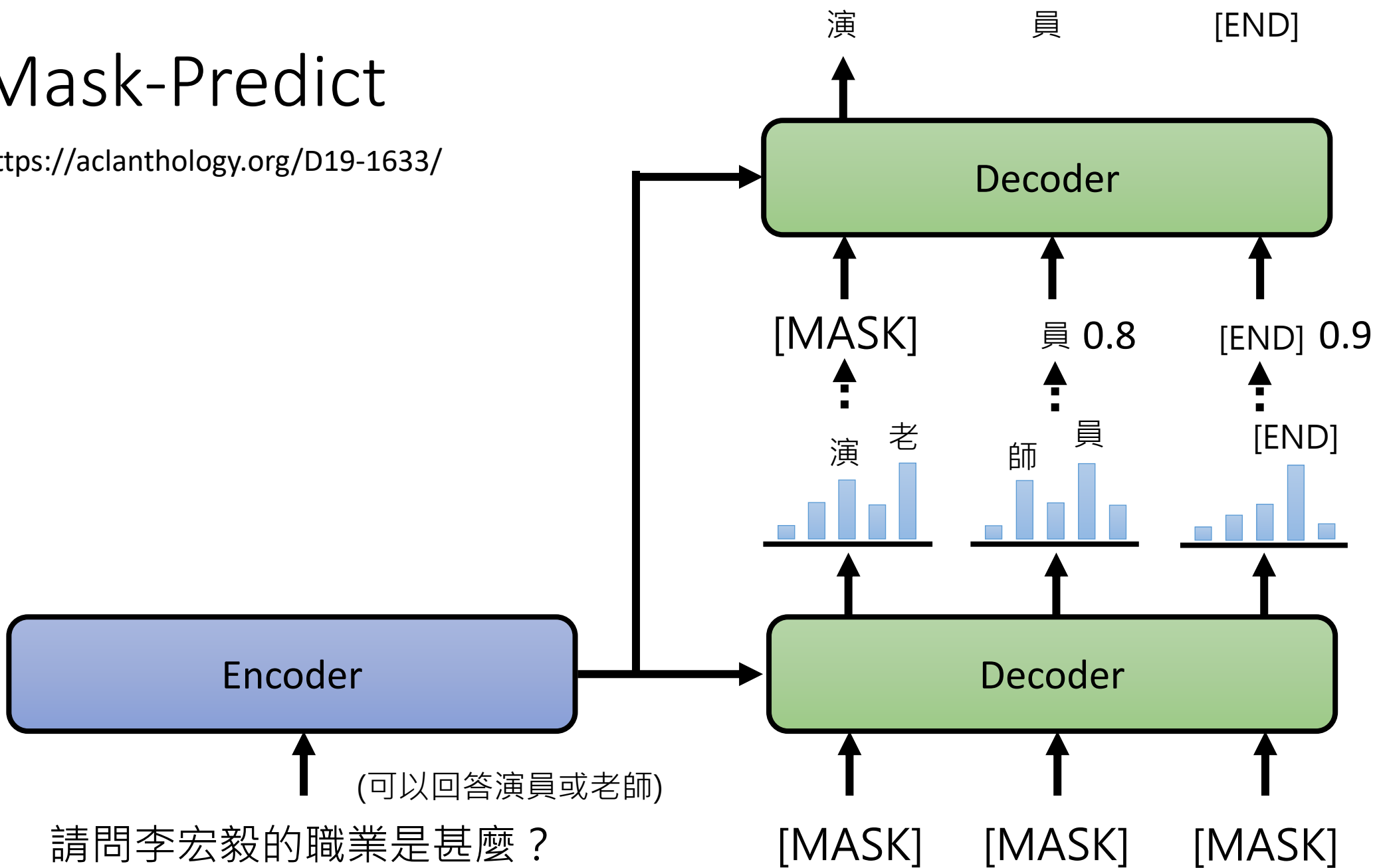
# Mask-Predict

<https://aclanthology.org/D19-1633/>



# Mask-Predict

<https://aclanthology.org/D19-1633/>

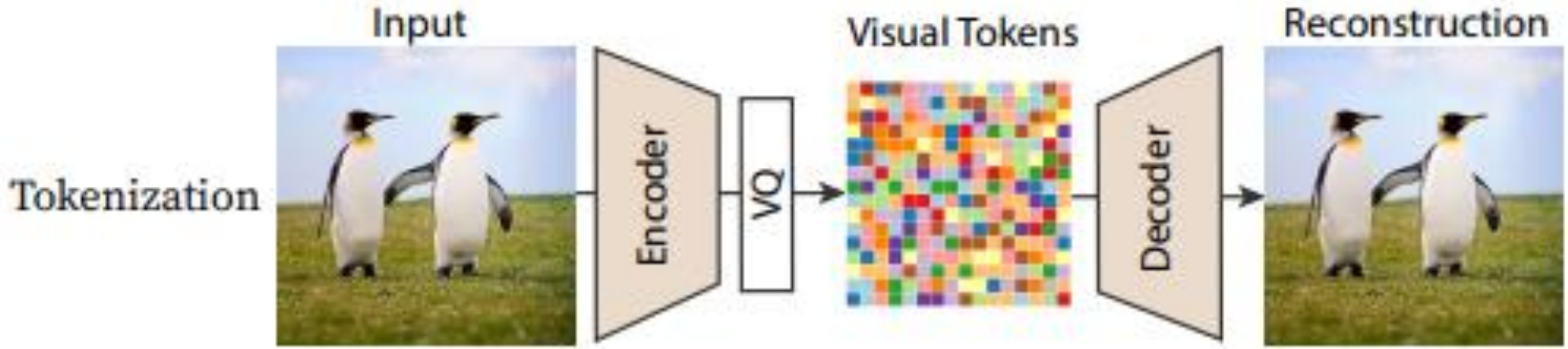




<https://arxiv.org/abs/2202.04200>

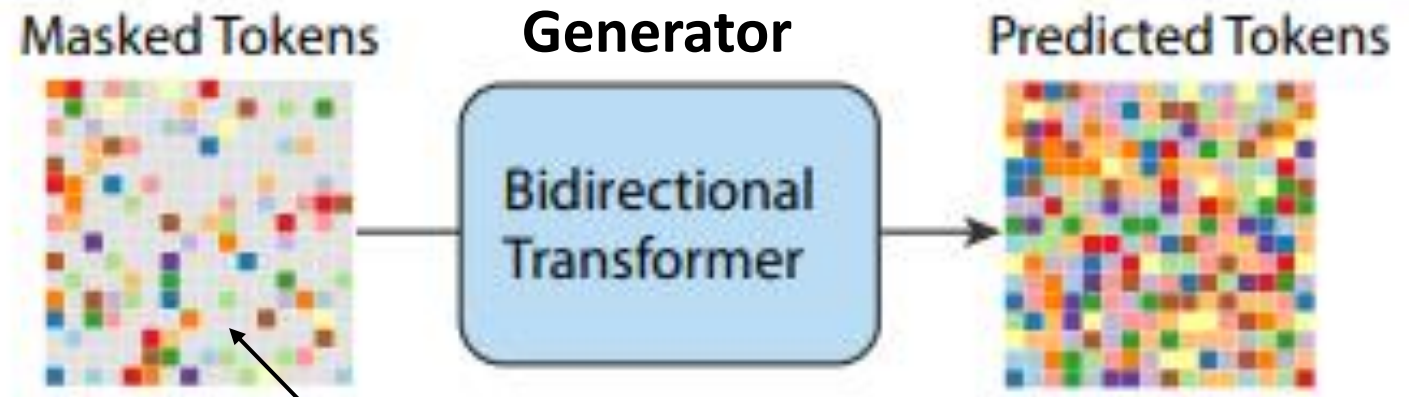
<https://arxiv.org/abs/2301.00704>

# Mask-Predict

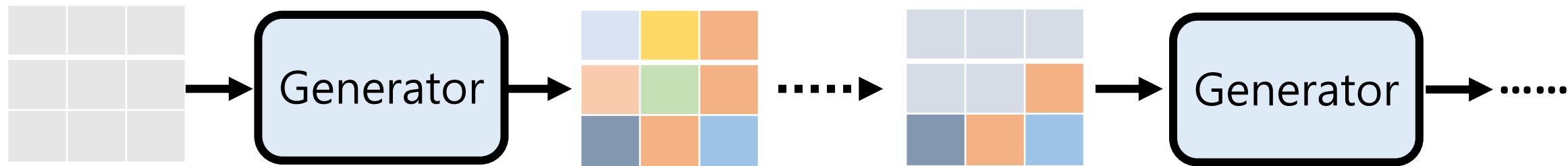


## Training

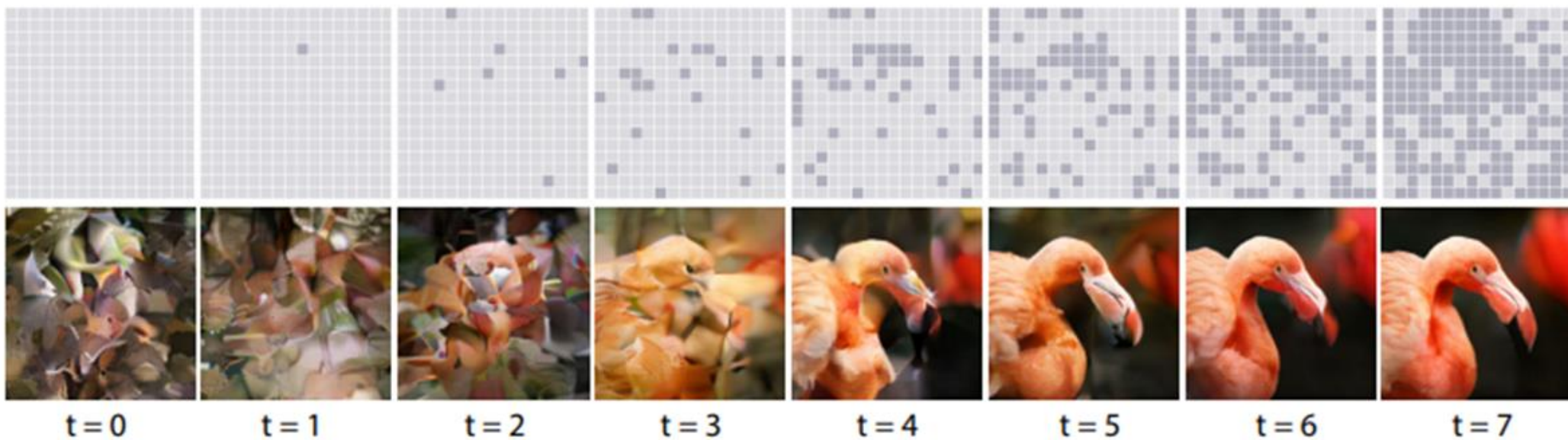
Masked Visual Token Modeling (MVTM)



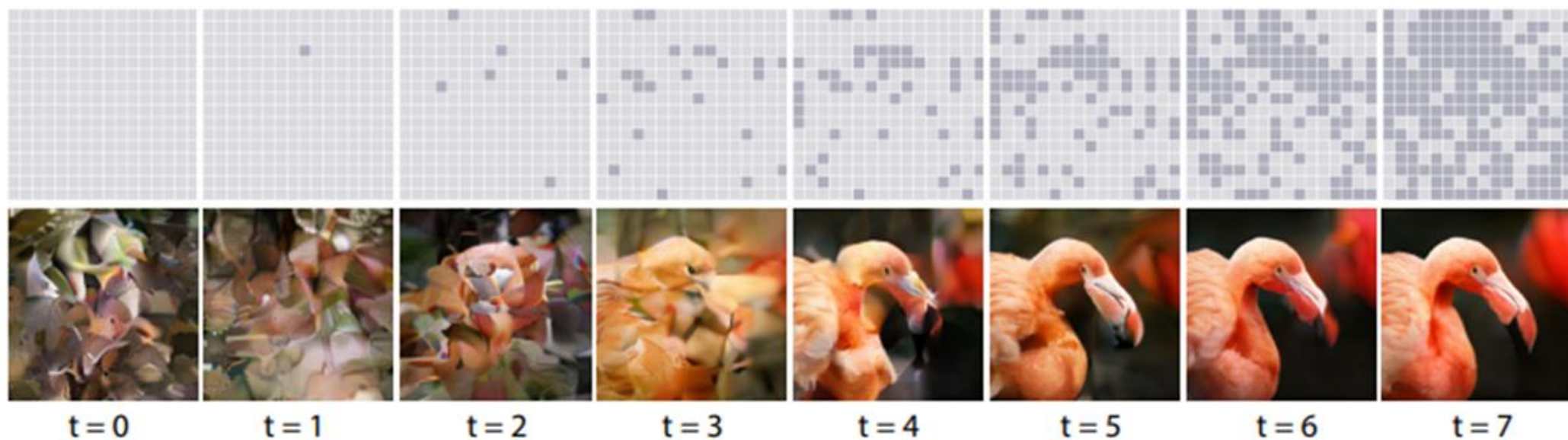
Gray: [mask] token



Scheduled  
Parallel  
Decoding  
with MaskGIT



Scheduled  
Parallel  
Decoding  
with MaskGIT



Sequential  
Decoding  
with Autoregressive  
Transformers

