

Adv ML for Gen-AI

Diffusion

<https://ml-graph.github.io/spring-2026/>

Yu Wang, Ph.D.

Assistant Professor

Department of Computer Science

University of Oregon

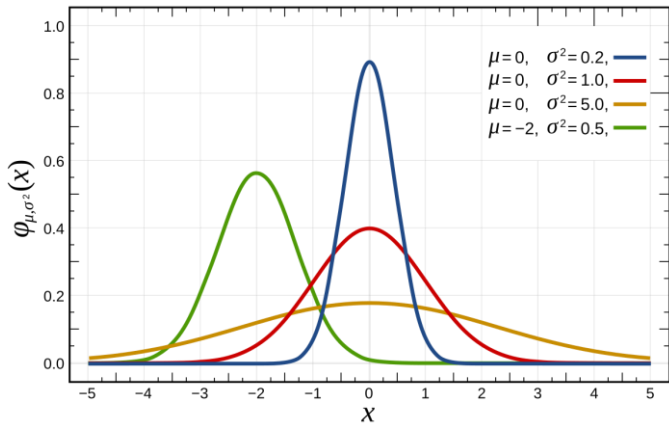
Personal: <https://yuwang0103.github.io/>

Lab: <https://kindlab-fly.github.io/>



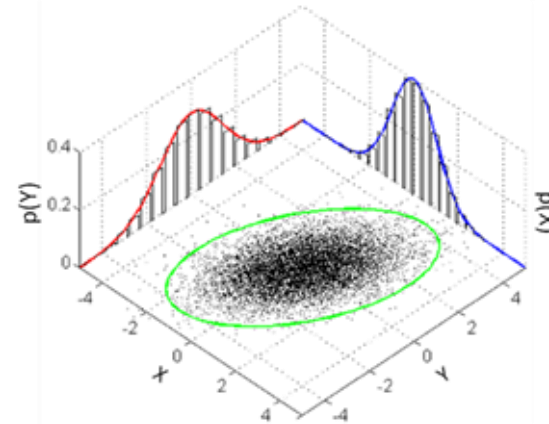
Summary

1D Gaussian Distribution

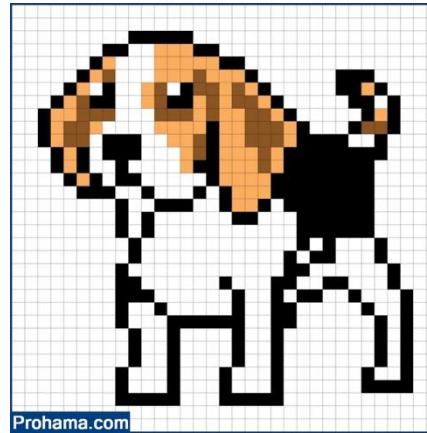
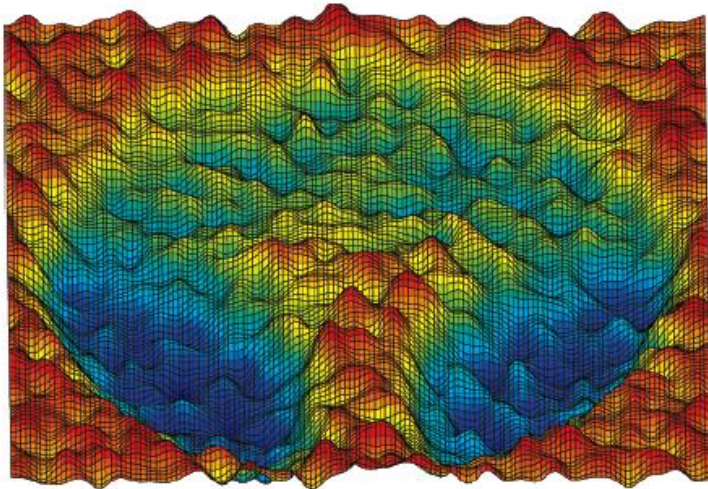


\mathbb{R}

2D Gaussian Distribution



\mathbb{R}^2



$\mathbb{R}^{256 \times 256}$

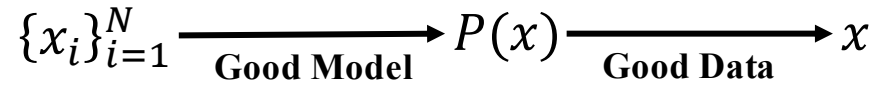


Summary

Probability distribution of the **objective** based on the **observed data**

- **Machine Learning Methods**

- Gaussian Kernel Density Estimation
- Gaussian Mixture Models



Using **existing function** to estimate what you do not know that can best fit your observation

- **Deep Learning Methods**

- Auto-Encoder (AE)
- Variational AE (VAE)
- Generative Adversarial Network (GAN)
- Diffusion Model
-

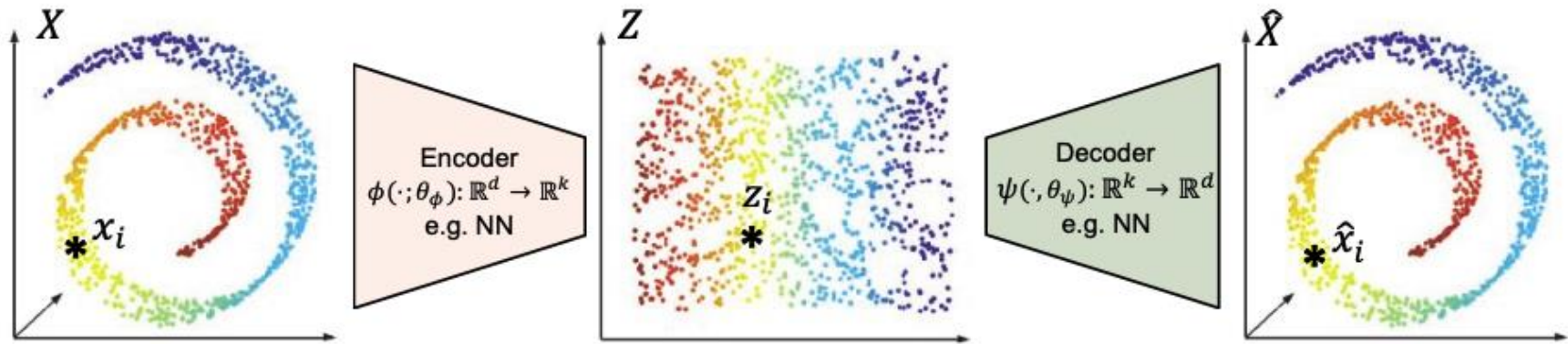
Using **learnable function** to estimate what you do not know that can best fit your observation



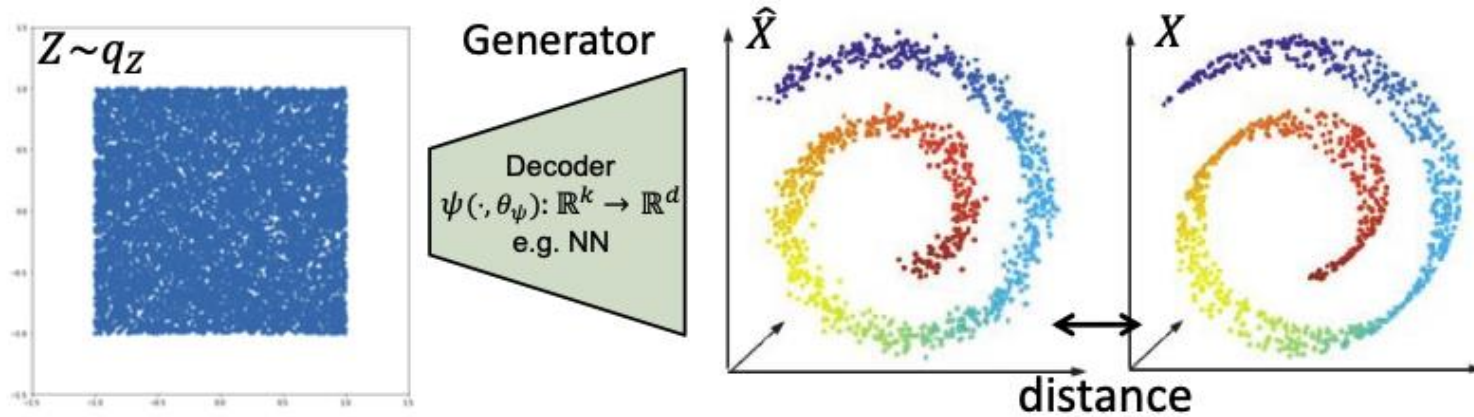


Summary

Auto-encoder based generative modeling



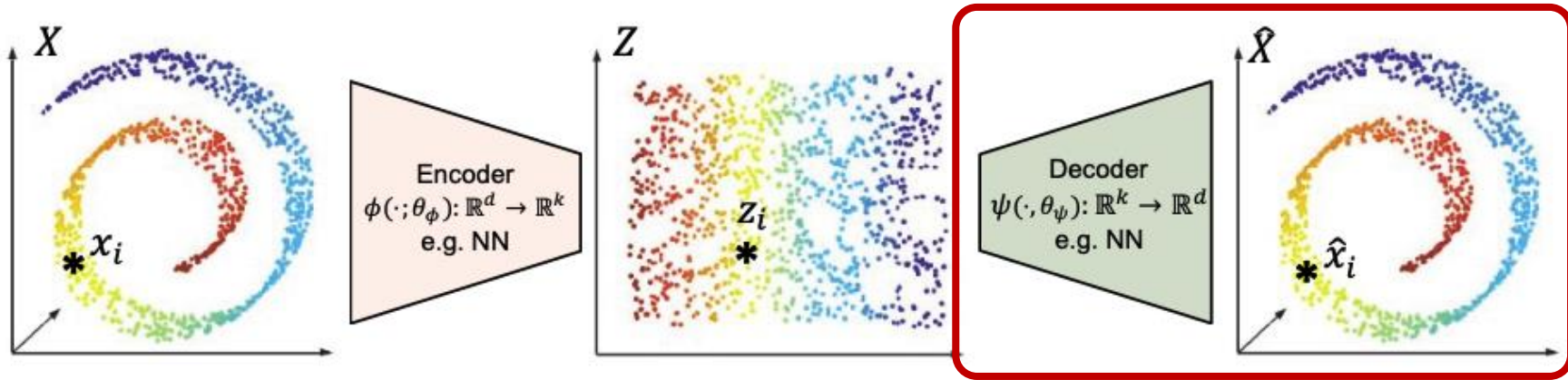
Encoder-less Generative Modeling



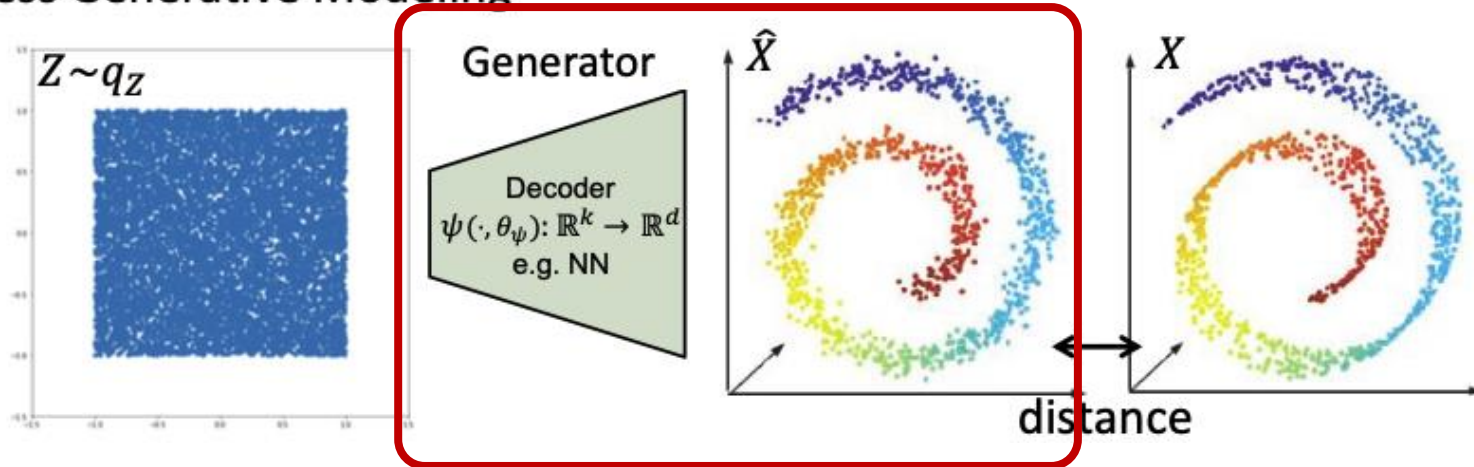


Summary – GAN/AE(VAE) shares the same decoder architecture

Auto-encoder based generative modeling



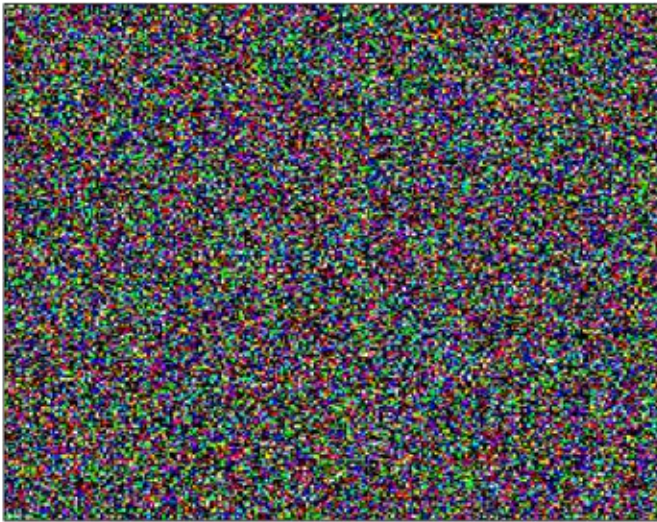
Encoder-less Generative Modeling



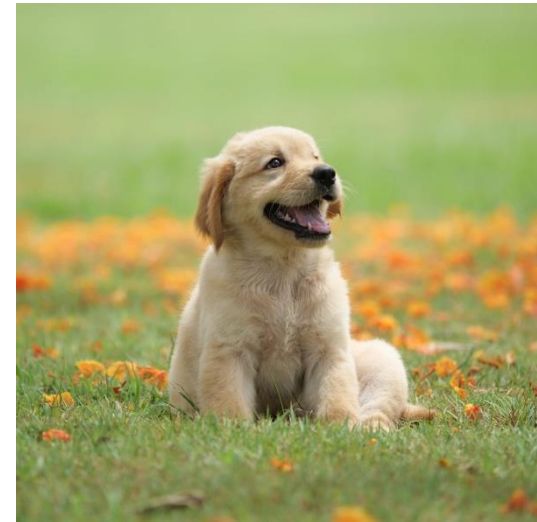


One-step Decoder-based Generation is too Difficult

Gaussian Noise 256x256 (3-channel)



**One-shot
Generation**

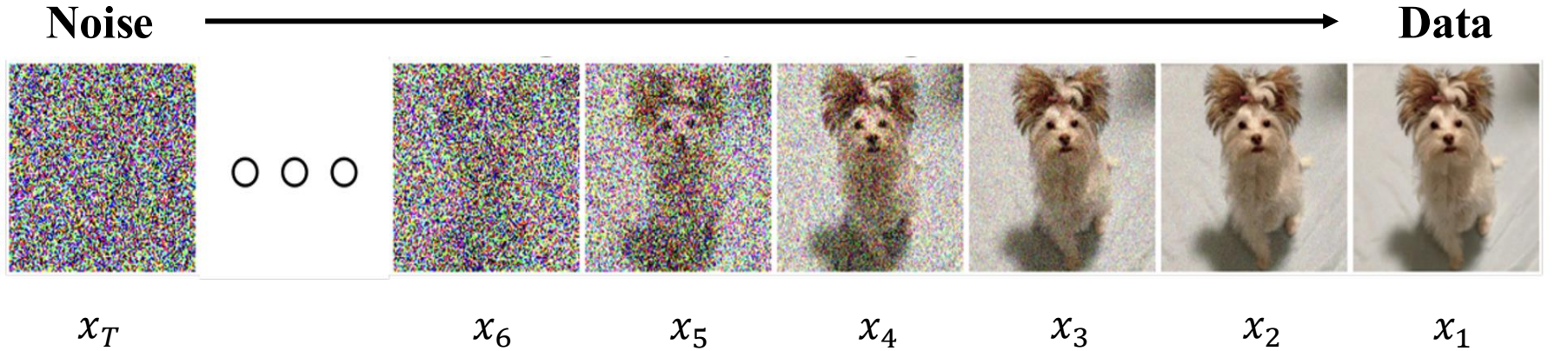


Standard Gaussian $\xrightarrow{\text{Too difficulty}}$ **Target Distribution**



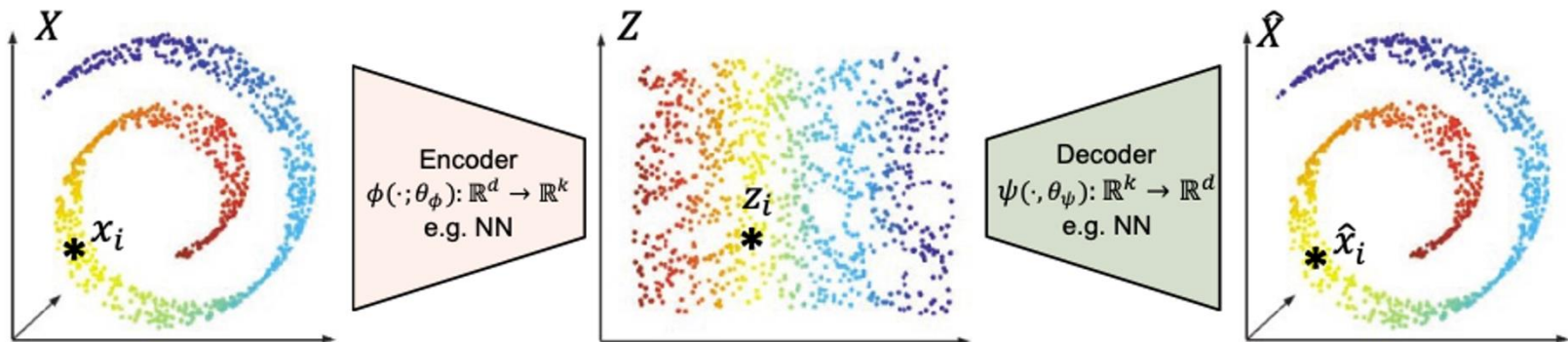
Motivation of Diffusion

Generation by denoising step by step



Predicted Ground-truth
 $\mathcal{F}(x) \quad \|\hat{x} - x\|_2^2$

Generation by one step





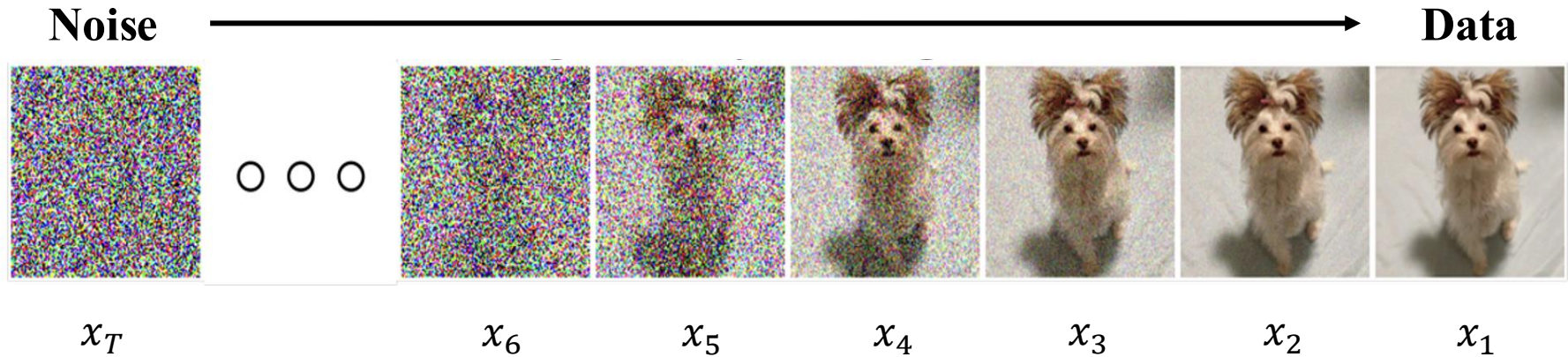
- **Motivation and Introduction of Diffusion**
- **Forward process of Diffusion**
- **Backward process of Diffusion**
 - **Backward progressive:** $x_t \rightarrow x_{t-1}$
 - **Backward one step:** $x_t \rightarrow x_0$
 - **Backward noise:** $x_t \rightarrow \epsilon$





Motivation of Diffusion

Generation by denoising step by step



Predicted

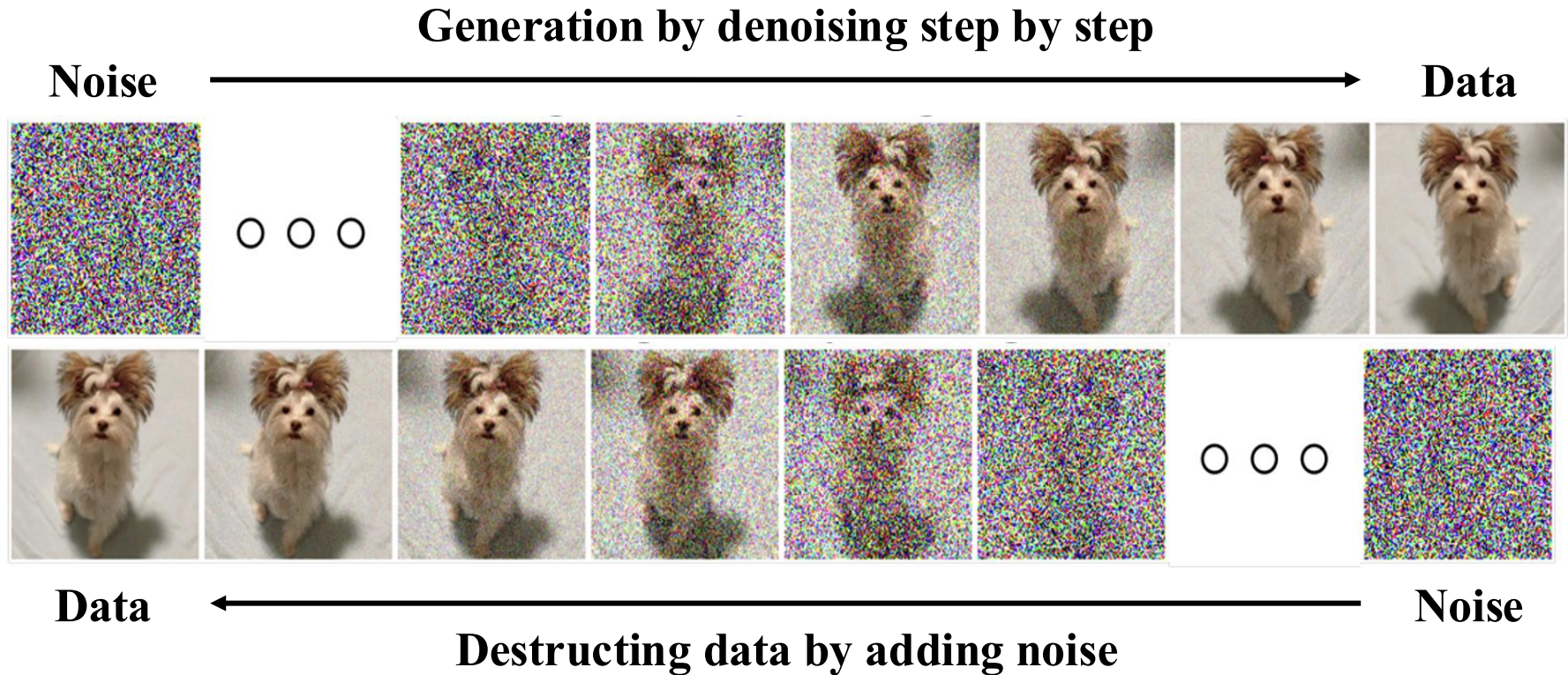
Ground-truth

$$\|\hat{x} - x_{t-1}\|_2^2 = \|\mathcal{F}(x_t, t) - x_{t-1}\|_2^2$$

However, how do we obtain $\{x_t\}_{t=1}^T$?



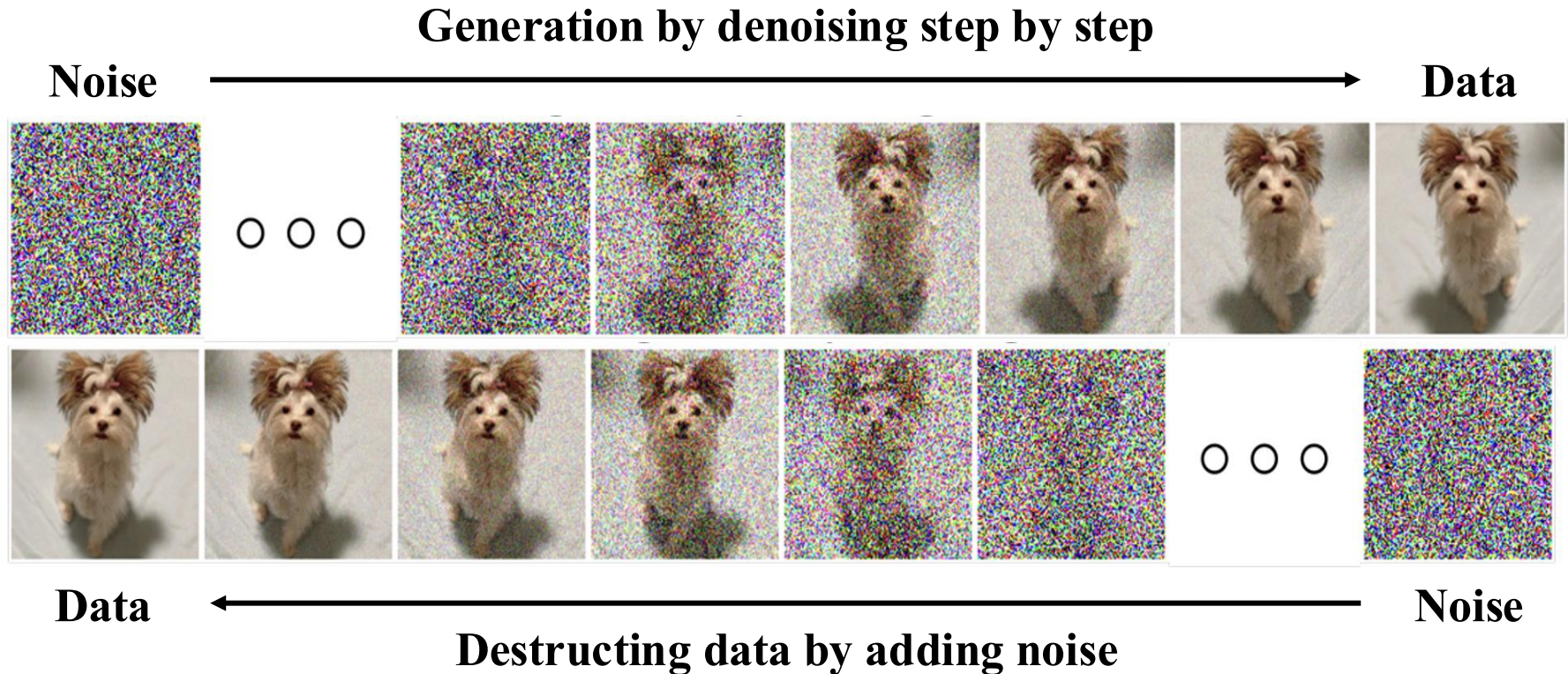
Motivation of Diffusion



However, how do we obtain $\{x_t\}_{t=1}^T$? – **Forward Process**



Motivation of Diffusion



- **Forward Process (Data - Noise):** Gradually **add noise** to generate intermediate data
- **Backward Process (Noise - Data):** **Learn** to denoise the intermediate data and recover back

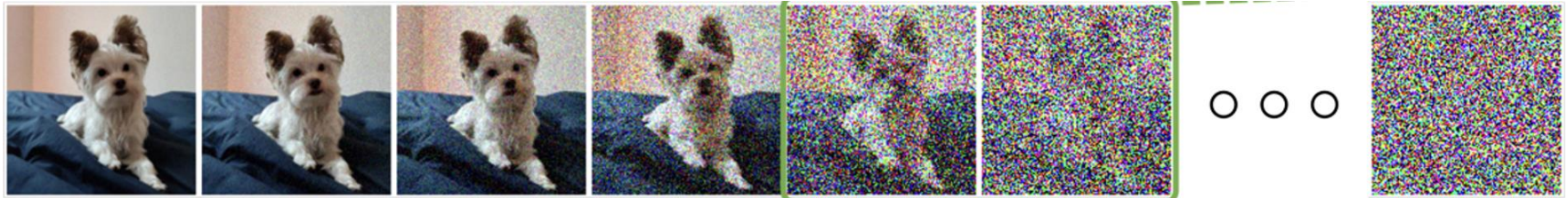
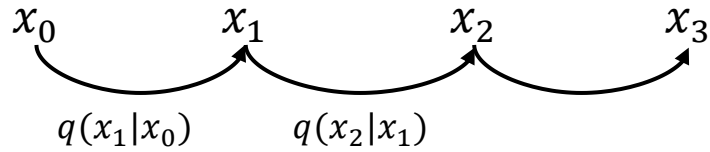
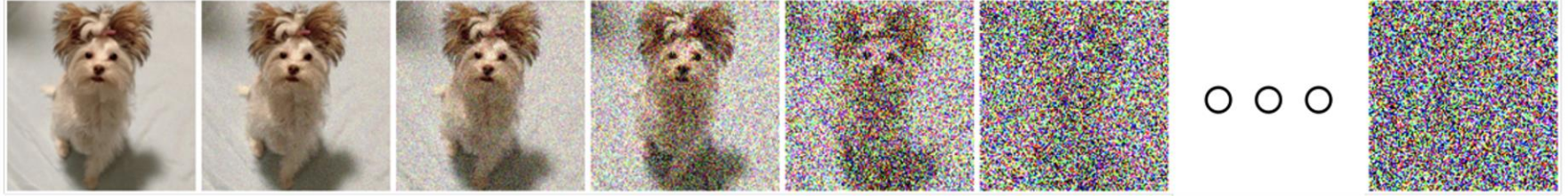


- **Motivation and Introduction of Diffusion**
- **Forward process of Diffusion**
- **Backward process of Diffusion**
 - **Backward progressive:** $x_t \rightarrow x_{t-1}$
 - **Backward one step:** $x_t \rightarrow x_0$
 - **Backward noise:** $x_t \rightarrow \epsilon$



Forward Process

Data $\xrightarrow{\hspace{2cm}}$ Destructing data by adding noise $\xrightarrow{\hspace{2cm}}$ Noise



$q(x_0)$ $q(x_1)$ $q(x_2)$ $q(x_3)$ $q(x_4)$ $q(x_5)$ $q(x_T)$

$q(x_t)$

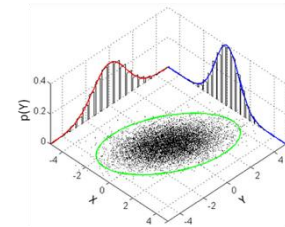
Data distribution at time step t

$q(x_0)$

Start

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

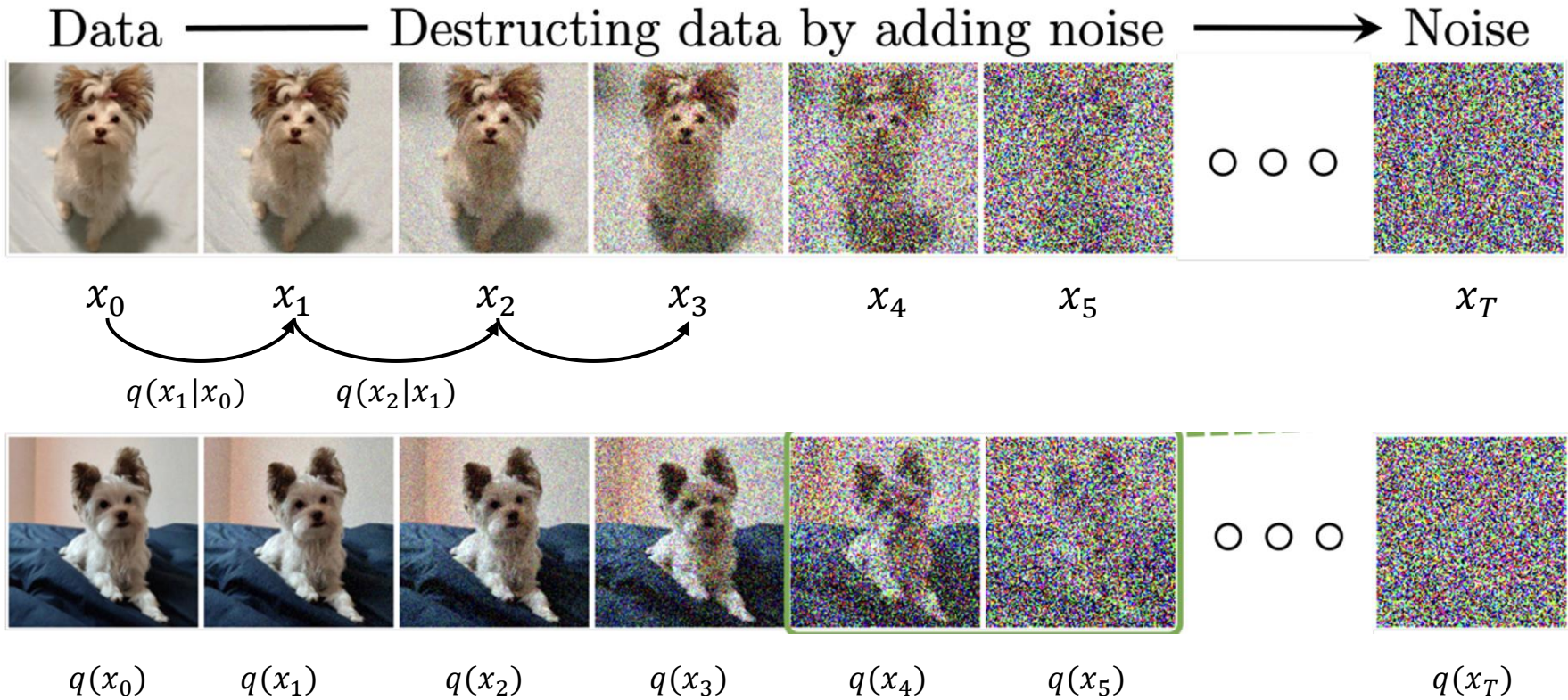
Transitional distribution at time step t



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



Forward Process



$$x_1 = \sqrt{\alpha_1}x_0 + \sqrt{(1 - \alpha_1)}\epsilon_1$$

$$\epsilon_1 \sim \mathcal{N}(0, I) \quad \epsilon_2 \sim \mathcal{N}(0, I)$$

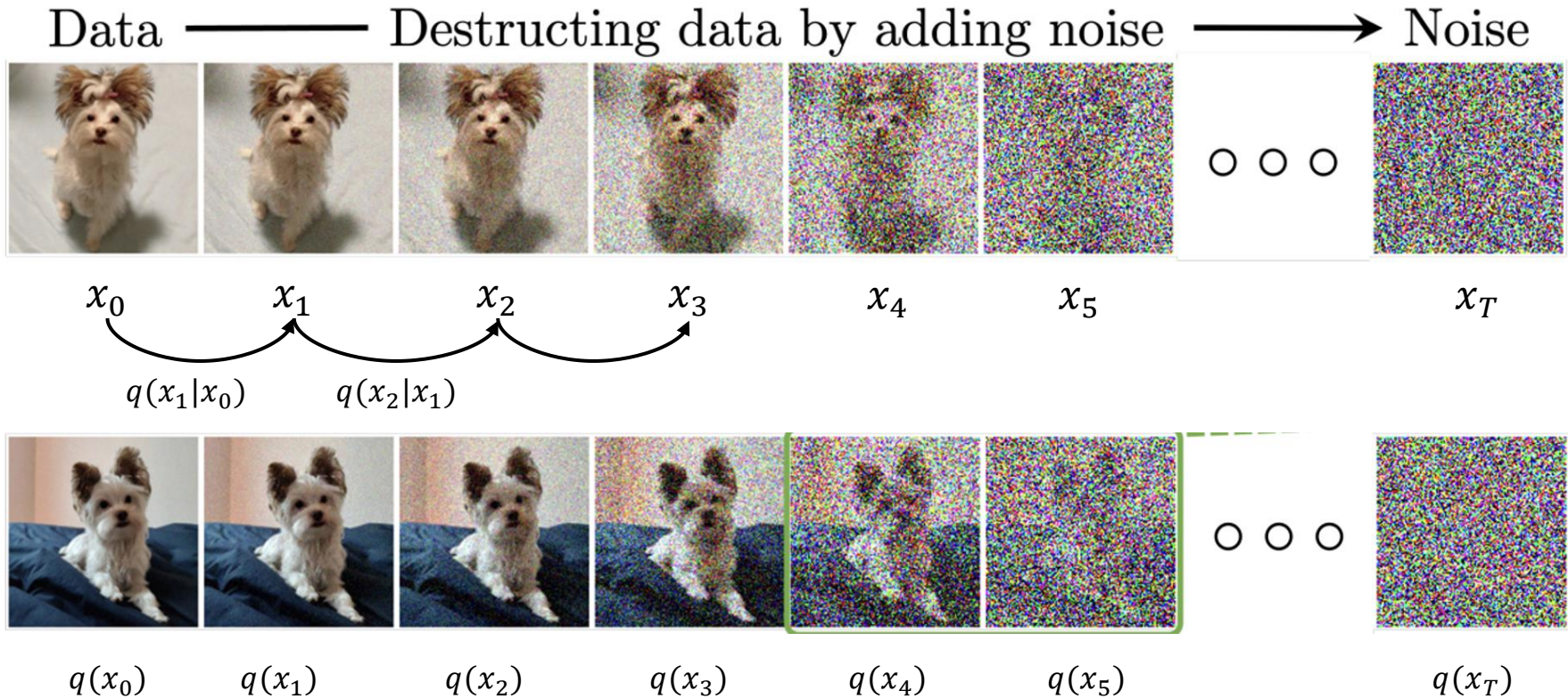
$$x_2 = \sqrt{\alpha_2}x_1 + \sqrt{(1 - \alpha_2)}\epsilon_2$$

$$A\epsilon_1 + B\epsilon_2 \sim \mathcal{N}(0, (A^2 + B^2)I)$$

$$x_2 = \sqrt{\alpha_2}(\sqrt{\alpha_1}x_0 + \sqrt{(1 - \alpha_1)}\epsilon_1) + \sqrt{(1 - \alpha_2)}\epsilon_2 = \sqrt{\alpha_2\alpha_1}x_0 + \sqrt{\alpha_2}\sqrt{1 - \alpha_1}\epsilon_1 + \sqrt{(1 - \alpha_2)}\epsilon_2$$



Forward Process



$$x_2 = \sqrt{\alpha_2}(\sqrt{\alpha_1}x_0 + \sqrt{(1 - \alpha_1)}\epsilon_1) + \sqrt{(1 - \alpha_2)}\epsilon_2 = \sqrt{\alpha_2\alpha_1}x_0 + \sqrt{\alpha_2}\sqrt{1 - \alpha_1}\epsilon_1 + \sqrt{(1 - \alpha_2)}\epsilon_2$$

$$\epsilon_1 \sim \mathcal{N}(0, I) \quad \epsilon_2 \sim \mathcal{N}(0, I)$$

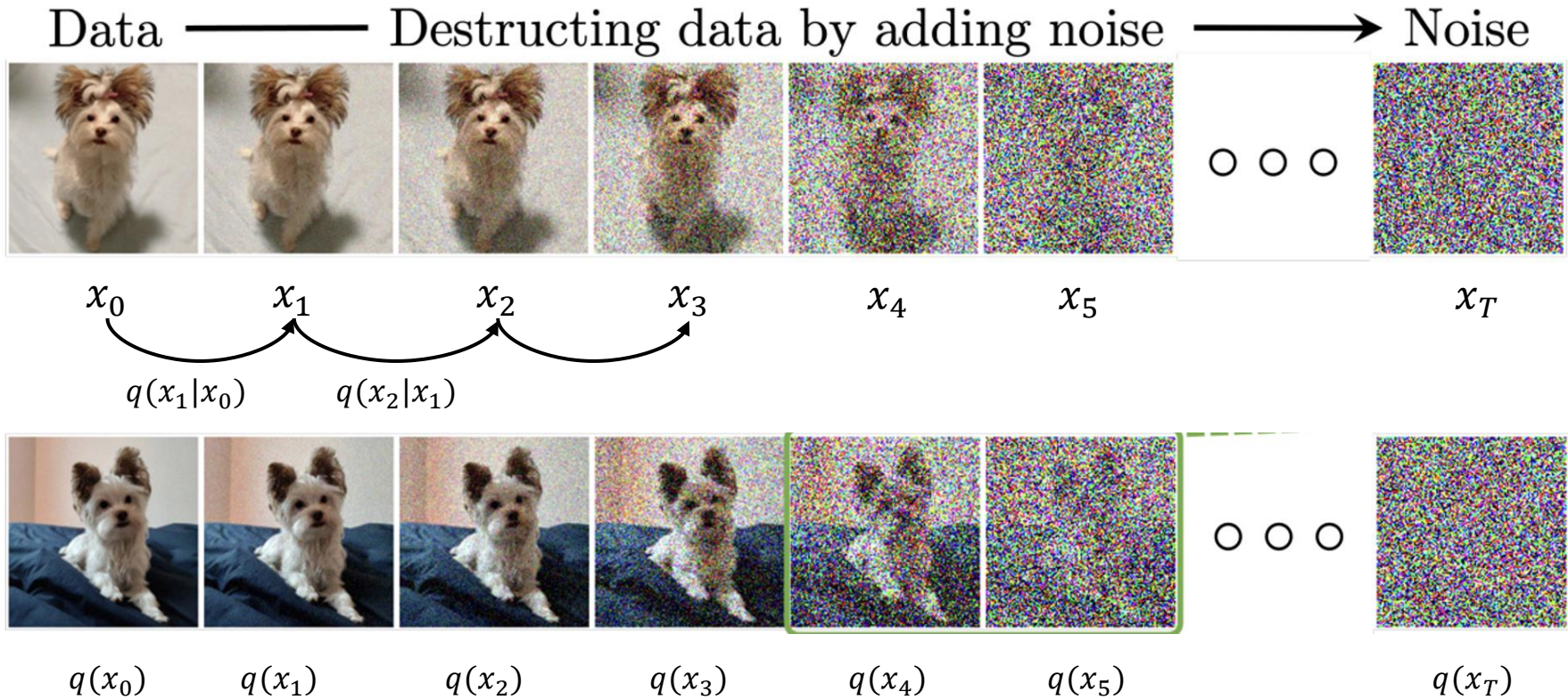
$$\sqrt{\alpha_2}\sqrt{1 - \alpha_1}\epsilon_1 + \sqrt{(1 - \alpha_2)}\epsilon_2 \sim \mathcal{N}(0, (1 - \alpha_1\alpha_2)I)$$

$$A\epsilon_1 + B\epsilon_2 \sim \mathcal{N}(0, (A^2 + B^2)I)$$

$$x_2 \sim \mathcal{N}(\sqrt{\alpha_2\alpha_1}x_0, (1 - \alpha_1\alpha_2)I)$$



Forward Process



$$x_2 = \sqrt{\alpha_2}(\sqrt{\alpha_1}x_0 + \sqrt{(1-\alpha_1)}\epsilon_1) + \sqrt{(1-\alpha_2)}\epsilon_2 = \sqrt{\alpha_2\alpha_1}x_0 + \sqrt{\alpha_2}\sqrt{1-\alpha_1}\epsilon_1 + \sqrt{(1-\alpha_2)}\epsilon_2$$

$$\epsilon_1 \sim \mathcal{N}(0, I) \quad \epsilon_2 \sim \mathcal{N}(0, I) \quad \sqrt{\alpha_2}\sqrt{1-\alpha_1}\epsilon_1 + \sqrt{(1-\alpha_2)}\epsilon_2 \sim \mathcal{N}(0, (1-\alpha_1\alpha_2)I)$$

$$A\epsilon_1 + B\epsilon_2 \sim \mathcal{N}(0, (A^2 + B^2)I) \quad x_2 \sim \mathcal{N}(\sqrt{\alpha_2\alpha_1}x_0, (1-\alpha_1\alpha_2)I) \quad x_3 \sim \mathcal{N}(\sqrt{\alpha_3\alpha_2\alpha_1}x_0, (1-\alpha_1\alpha_2\alpha_3)I)$$



Forward Process

$$q(x_0) \quad x_0 \sim \mathcal{N}(\mathbf{0}, I)$$

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

$$q(x_t|x_{t-1}) \quad x_t \sim \mathcal{N}(\sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)I)$$

$$\bar{\alpha}_t = \prod_{s=1}^t \alpha_s = \prod_{s=1}^t (1 - \beta_s)$$

$$\beta_s \in [0,1]$$

Forward process parameters

```
def linear_beta_schedule(T, beta_start=1e-4, beta_end=2e-2):
    return torch.linspace(beta_start, beta_end, T)
```

```
T = 200
```

```
betas = linear_beta_schedule(T).to(device)
```

```
alphas = 1.0 - betas
```

```
alpha_bar = torch.cumprod(alphas, dim=0)
```

```
print('T:', T)
```

```
print('beta range:', float(betas.min()), '->', float(betas.max()))
```

```
print('alpha_bar[0]:', float(alpha_bar[0]), 'alpha_bar[-1]:', float(alpha_bar[-1]))
```

Closed-form sampling from $q(x_t | x_0)$

```
def q_sample(x0, t, noise=None):
```

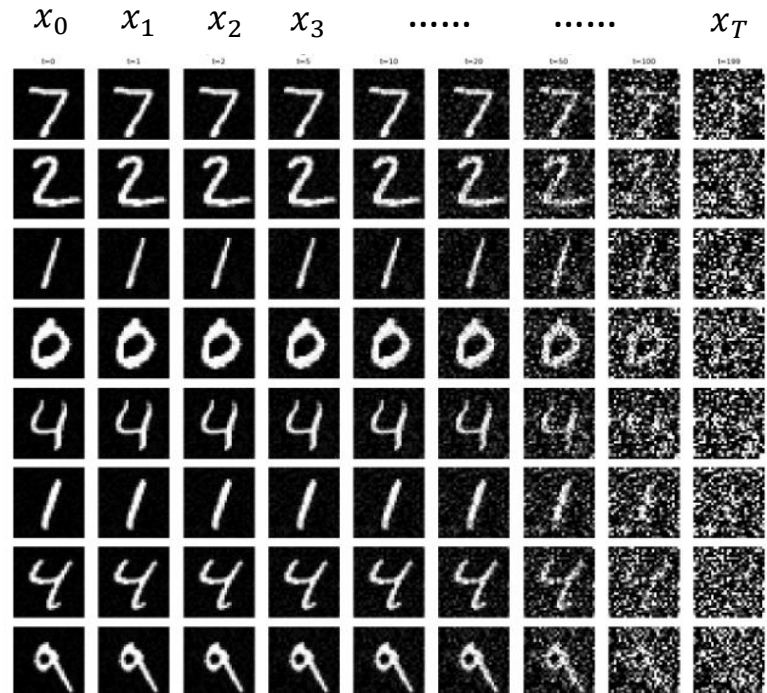
```
    # Sample x_t given x_0 at discrete time index t (0..T-1)
```

```
    if noise is None:
```

```
        noise = torch.randn_like(x0)
```

```
    a_bar = alpha_bar[t].view(-1, 1, 1, 1) # (B,1,1,1)
```

```
    return torch.sqrt(a_bar) * x0 + torch.sqrt(1.0 - a_bar) * noise
```





Forward Process

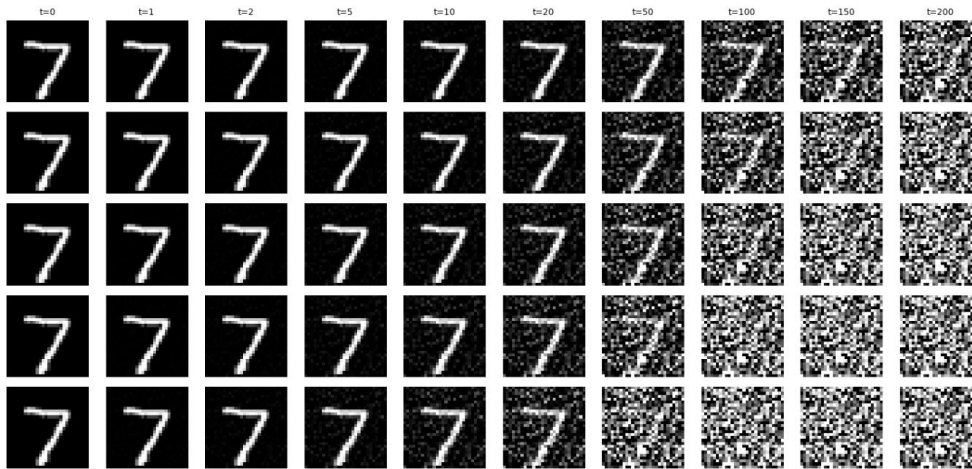
$$q(x_0) \quad x_0 \sim \mathcal{N}(\mathbf{0}, I)$$

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

$$q(x_t|x_{t-1}) \quad x_t \sim \mathcal{N}(\sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)I)$$

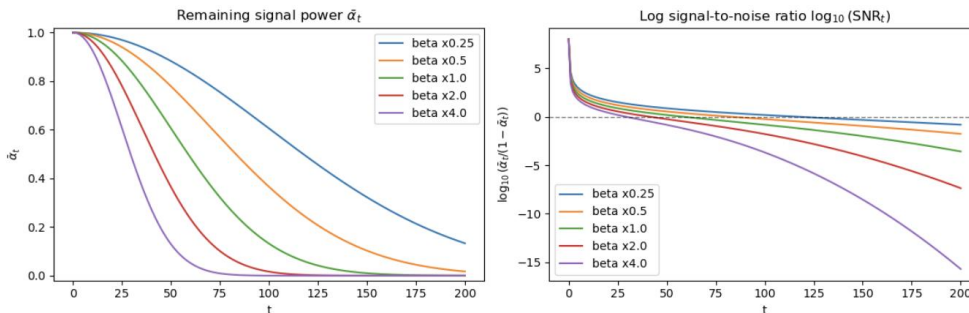
$$\bar{\alpha}_t = \prod_{s=1}^T \alpha_s = \prod_{s=1}^T (1 - \beta_s)$$

$$\beta_s \in [0,1]$$



$$\text{SNR}_t = \frac{\bar{\alpha}_t}{1 - \bar{\alpha}_t}$$

Signal Noise Ratio





Forward Process

$$q(x_0) \quad x_0 \sim \mathcal{N}(\mathbf{0}, I)$$

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

$$q(x_t|x_{t-1}) \quad x_t \sim \mathcal{N}(\sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)I)$$

$$\bar{\alpha}_t = \prod_{s=1}^T \alpha_s = \prod_{s=1}^T (1 - \beta_s)$$

$$\beta_s \in [0,1]$$

```
# Forward process parameters
```

```
def linear_beta_schedule(T, beta_start=1e-4, beta_end=2e-2):  
    return torch.linspace(beta_start, beta_end, T)
```

```
T = 200
```

```
betas = linear_beta_schedule(T).to(device)
```

```
alphas = 1.0 - betas
```

```
alpha_bar = torch.cumprod(alphas, dim=0)
```

```
print('T:', T)
```

```
print('beta range:', float(betas.min()), '->', float(betas.max()))
```

```
print('alpha_bar[0]:', float(alpha_bar[0]), 'alpha_bar[-1]:', float(alpha_bar[-1]))
```

```
# Closed-form sampling from  $q(x_t | x_0)$ 
```

```
def q_sample(x0, t, noise=None):
```

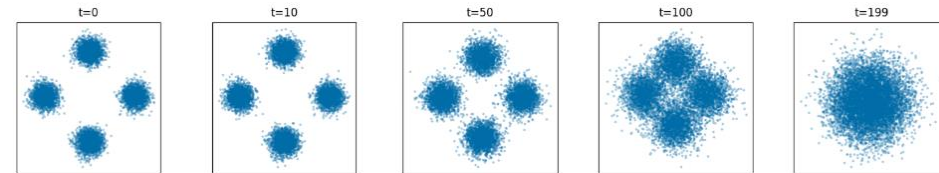
```
    # Sample  $x_t$  given  $x_0$  at discrete time index  $t$  (0..T-1)
```

```
    if noise is None:
```

```
        noise = torch.randn_like(x0)
```

```
    a_bar = alpha_bar[t].view(-1, 1, 1, 1) # (B,1,1,1)
```

```
    return torch.sqrt(a_bar) * x0 + torch.sqrt(1.0 - a_bar) * noise
```





Forward Process

$$q(x_0) \quad x_0 \sim \mathcal{N}(\mathbf{0}, I)$$

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

$$q(x_t|x_{t-1}) \quad x_t \sim \mathcal{N}(\sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)I)$$

$$\bar{\alpha}_t = \prod_{s=1}^T \alpha_s = \prod_{s=1}^T (1 - \beta_s)$$

$$\beta_s \in [0,1]$$

Forward process parameters

```
def linear_beta_schedule(T, beta_start=1e-4, beta_end=2e-2):  
    return torch.linspace(beta_start, beta_end, T)
```

```
T = 200
```

```
betas = linear_beta_schedule(T).to(device)
```

```
alphas = 1.0 - betas
```

```
alpha_bar = torch.cumprod(alphas, dim=0)
```

```
print('T:', T)
```

```
print('beta range:', float(betas.min()), '->', float(betas.max()))
```

```
print('alpha_bar[0]:', float(alpha_bar[0]), 'alpha_bar[-1]:', float(alpha_bar[-1]))
```

Closed-form sampling from $q(x_t | x_0)$

```
def q_sample(x0, t, noise=None):
```

```
    # Sample  $x_t$  given  $x_0$  at discrete time index  $t$  (0..T-1)
```

```
    if noise is None:
```

```
        noise = torch.randn_like(x0)
```

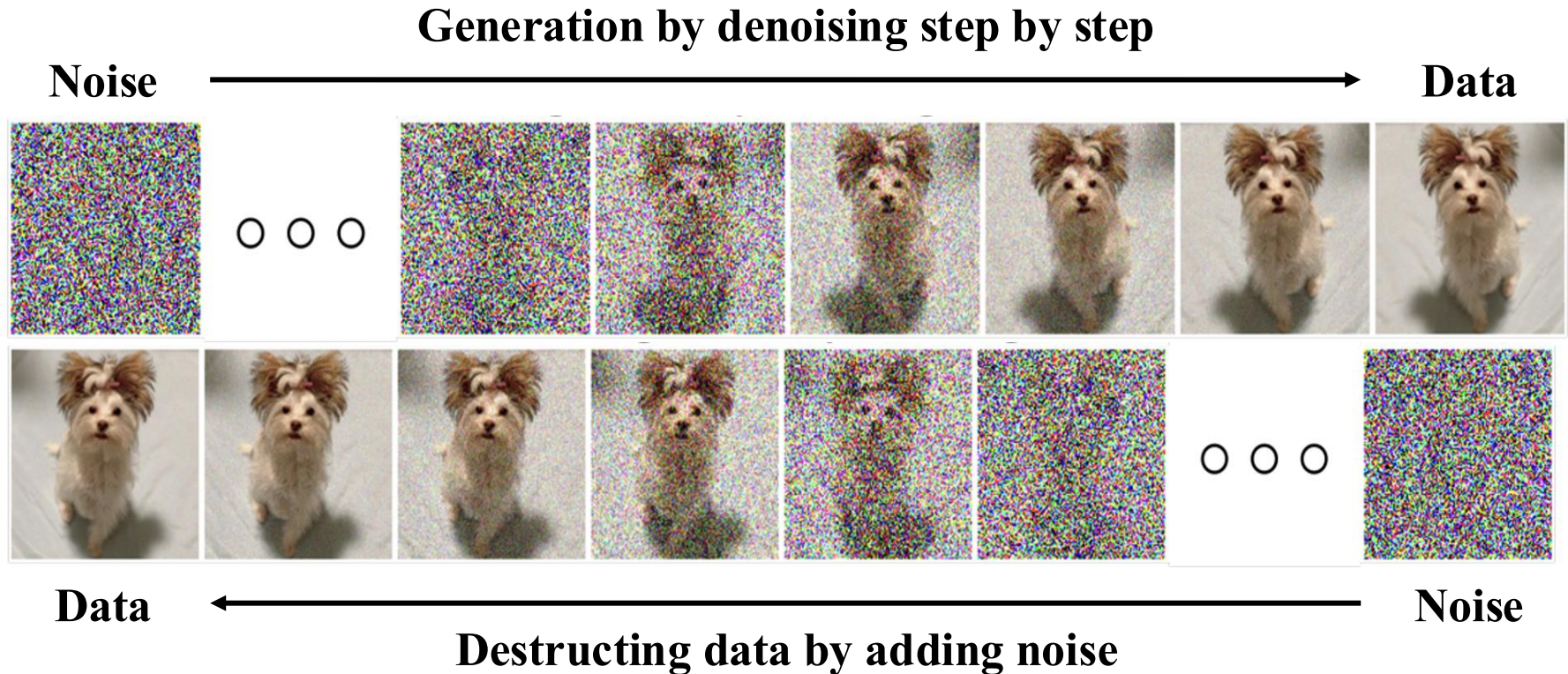
```
    a_bar = alpha_bar[t].view(-1, 1, 1, 1) # (B,1,1,1)
```

```
    return torch.sqrt(a_bar) * x0 + torch.sqrt(1.0 - a_bar) * noise
```

Dog Image



Forward/Backward Process of Diffusion

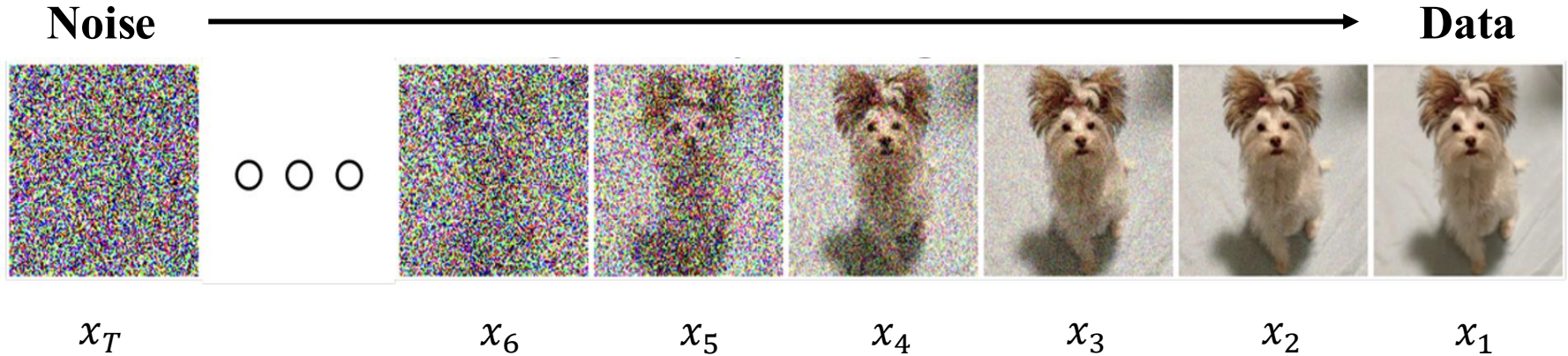


- **Forward Process (Data - Noise):** Gradually **add noise** to generate intermediate data ✓
- **Backward Process (Noise - Data):** **Learn** to denoise the intermediate data and recover back



We now can use forward process to obtain data pairs

Generation by denoising step by step



Predicted

Ground-truth



$$\|\hat{x} - x_{t-1}\|_2^2 = \|\mathcal{F}(x_t, t) - x_{t-1}\|_2^2$$

x_{t-1}

x_t



$$\mathbb{E}_{x_0 \sim p_{data}(x), t \sim U(1, T), x_{t-1} \sim q(x_{t-1} | x_0), x_t \sim q(x_t | x_{t-1})} \|\mathcal{F}(x_t, t) - x_{t-1}\|_2^2$$

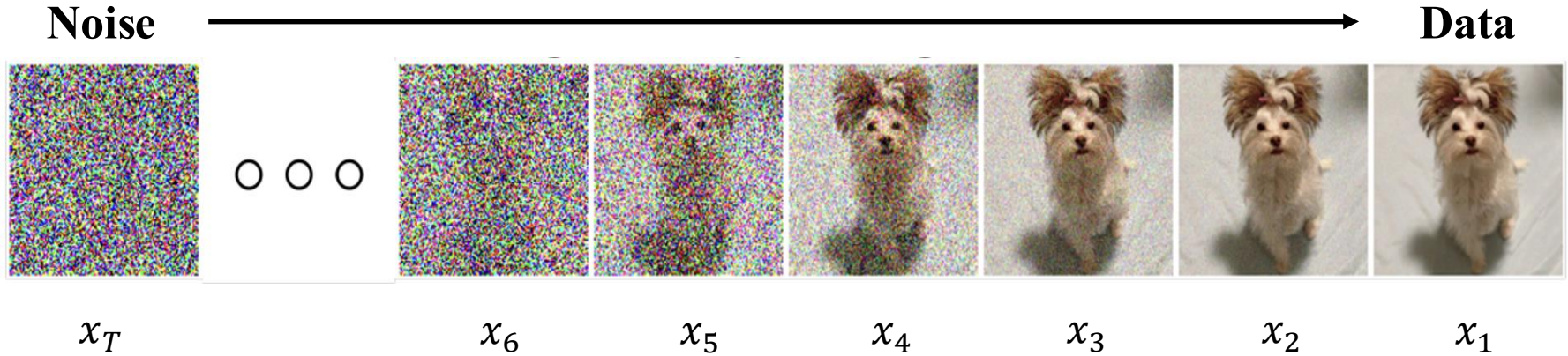


- **Motivation and Introduction of Diffusion**
- **Forward process of Diffusion**
- **Backward process of Diffusion**
 - **Backward progressive:** $x_t \rightarrow x_{t-1}$
 - **Backward one step:** $x_t \rightarrow x_0$
 - **Backward noise:** $x_t \rightarrow \epsilon$



Backward Optimization

Generation by denoising step by step



Predicted

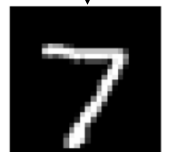
Ground-truth

$$\|\hat{x} - x_{t-1}\|_2^2 = \|\mathcal{F}(x_t, t) - x_{t-1}\|_2^2$$

$p_{\text{data}}(x)$



Sample



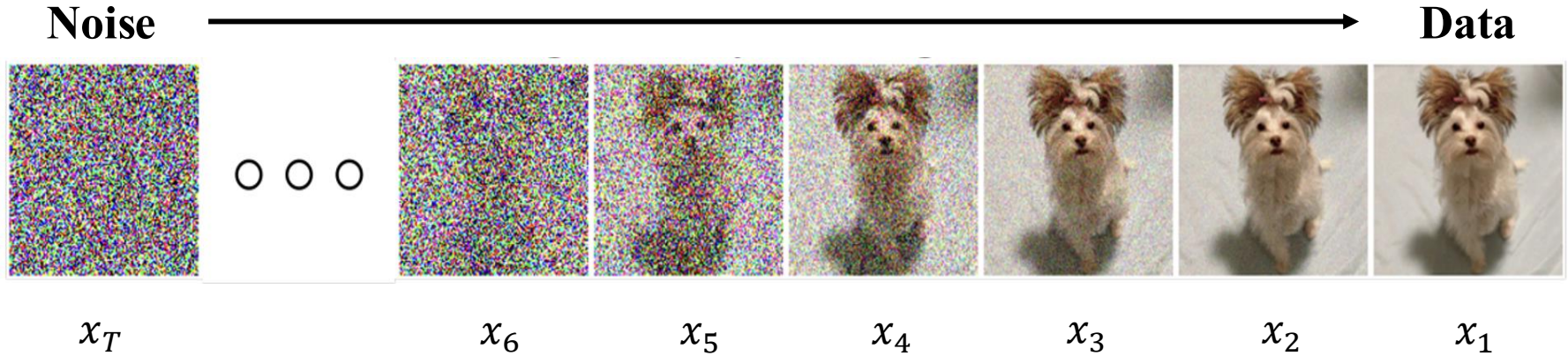
x_0

$$\mathbb{E}_{x_0 \sim p_{\text{data}}(x), t \sim U(1, T), x_{t-1} \sim q(x_{t-1} | x_0), x_t \sim q(x_t | x_{t-1})} \|\mathcal{F}(x_t, t) - x_{t-1}\|_2^2$$



Backward Optimization

Generation by denoising step by step

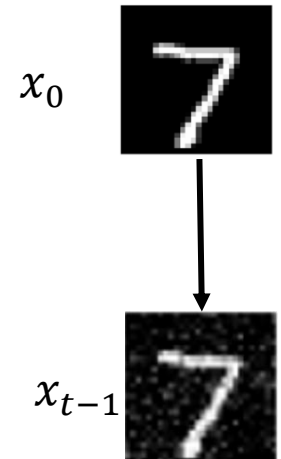


Predicted

Ground-truth

$$\|\hat{x} - x_{t-1}\|_2^2 = \|\mathcal{F}(x_t, t) - x_{t-1}\|_2^2$$

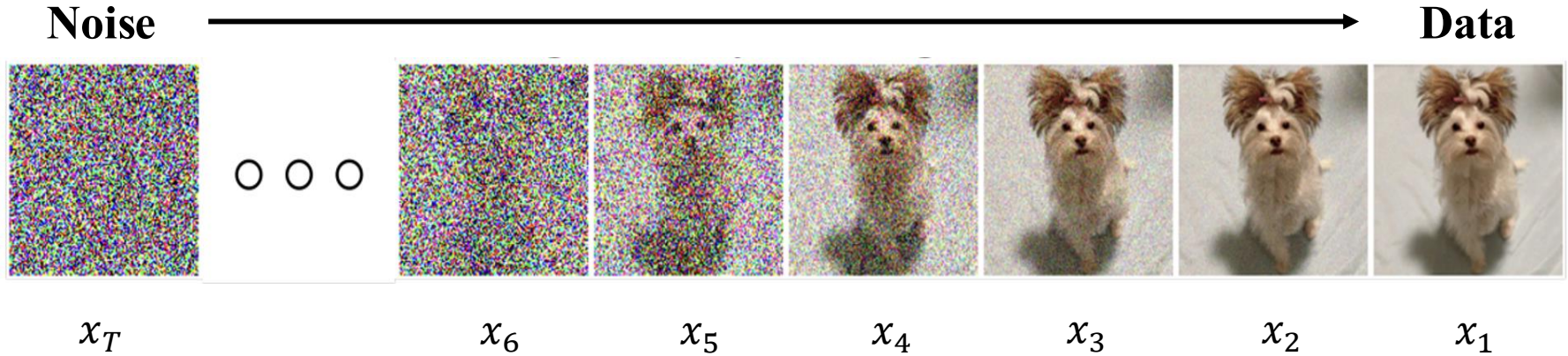
$$\mathbb{E}_{x_0 \sim p_{data}(x), t \sim U(1, T), x_{t-1} \sim q(x_{t-1} | x_0), x_t \sim q(x_t | x_{t-1})} \|\mathcal{F}(x_t, t) - x_{t-1}\|_2^2$$





Backward Optimization

Generation by denoising step by step

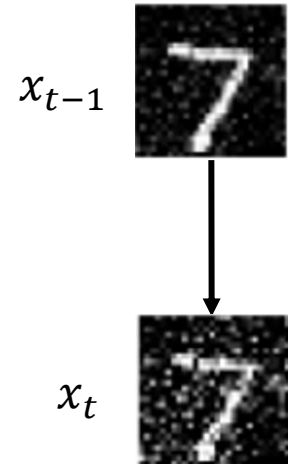


Predicted

Ground-truth

$$\|\hat{x} - x_{t-1}\|_2^2 = \|\mathcal{F}(x_t, t) - x_{t-1}\|_2^2$$

$$\mathbb{E}_{x_0 \sim p_{data}(x), t \sim U(1, T), x_{t-1} \sim q(x_{t-1} | x_0), x_t \sim q(x_t | x_{t-1})} \|\mathcal{F}(x_t, t) - x_{t-1}\|_2^2$$



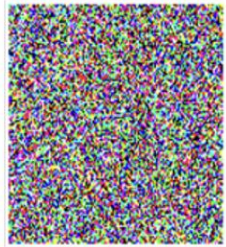


Backward Optimization

Generation by denoising step by step

Noise

Data



x_T

x_6

x_5

x_4

x_3

x_2

x_1

Predicted

Ground-truth



$$\|\hat{x} - x_{t-1}\|_2^2 = \|\mathcal{F}(x_t, t) - x_{t-1}\|_2^2$$

x_{t-1}

x_t

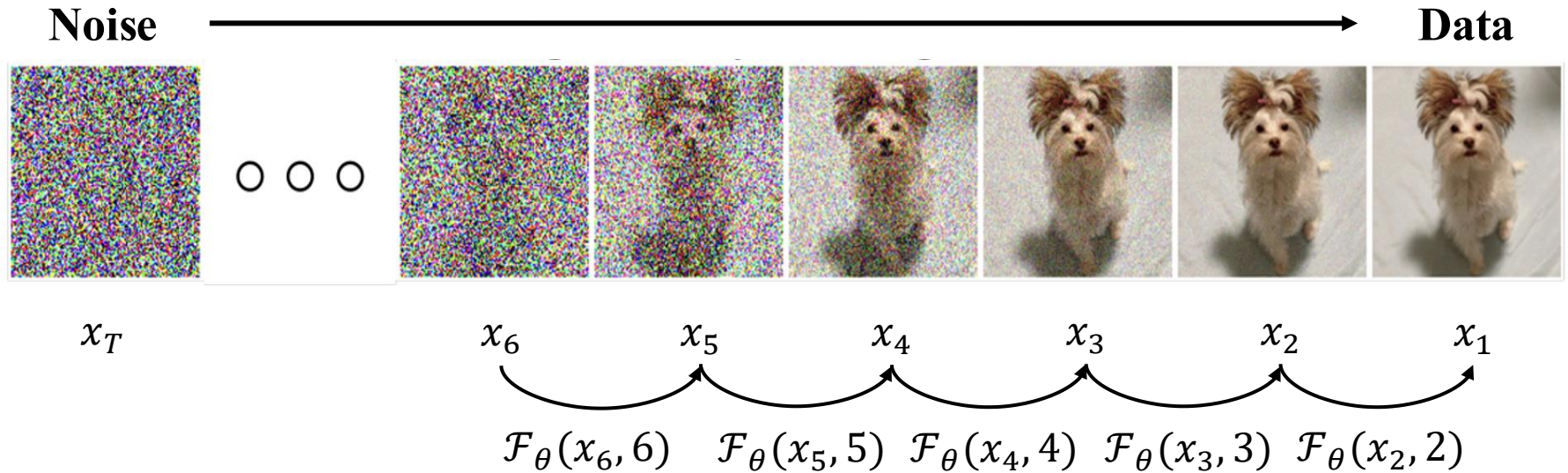


$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{data}(x), t \sim U(1, T), x_{t-1} \sim q(x_{t-1} | x_0), x_t \sim q(x_t | x_{t-1})} \|\mathcal{F}_\theta(x_t, t) - x_{t-1}\|_2^2$$

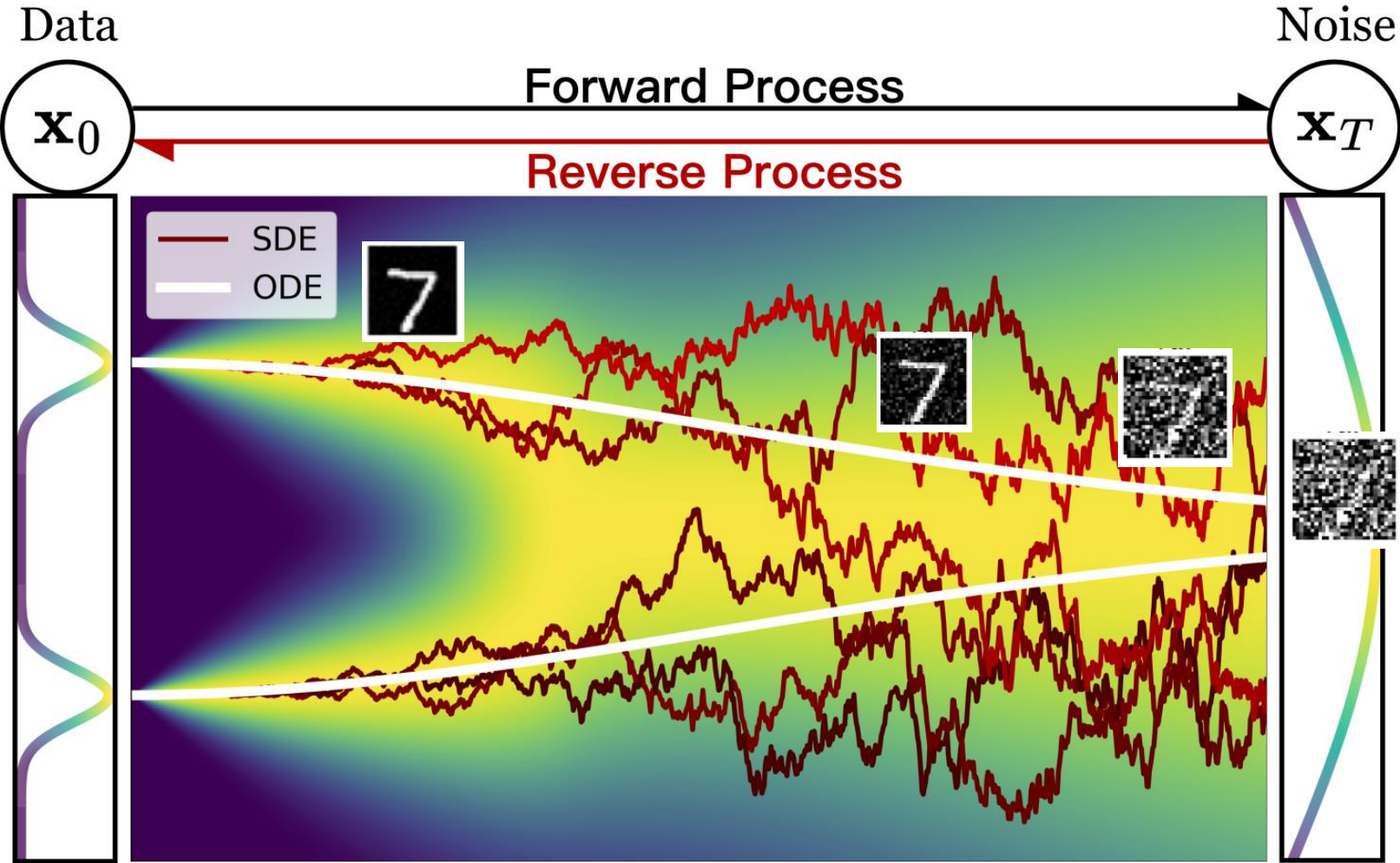


Generation

Generation by denoising step by step



Recursively generate back the original image



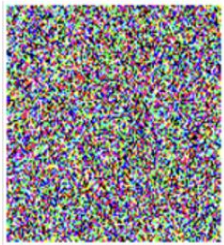


Generation - Problem

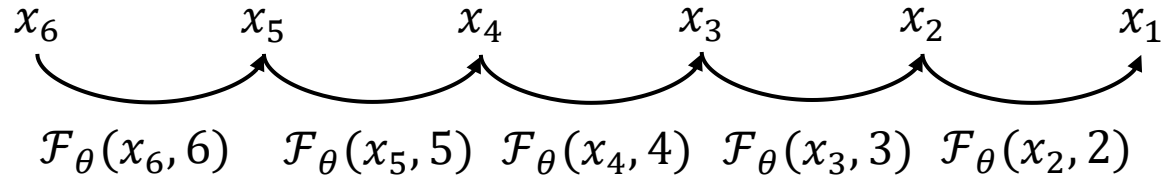
Generation by denoising step by step

Noise

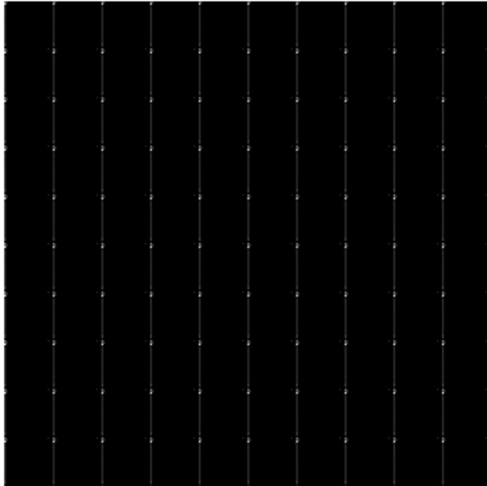
Data



x_T



Naive generation (deterministic): repeatedly predict x_{t-1}

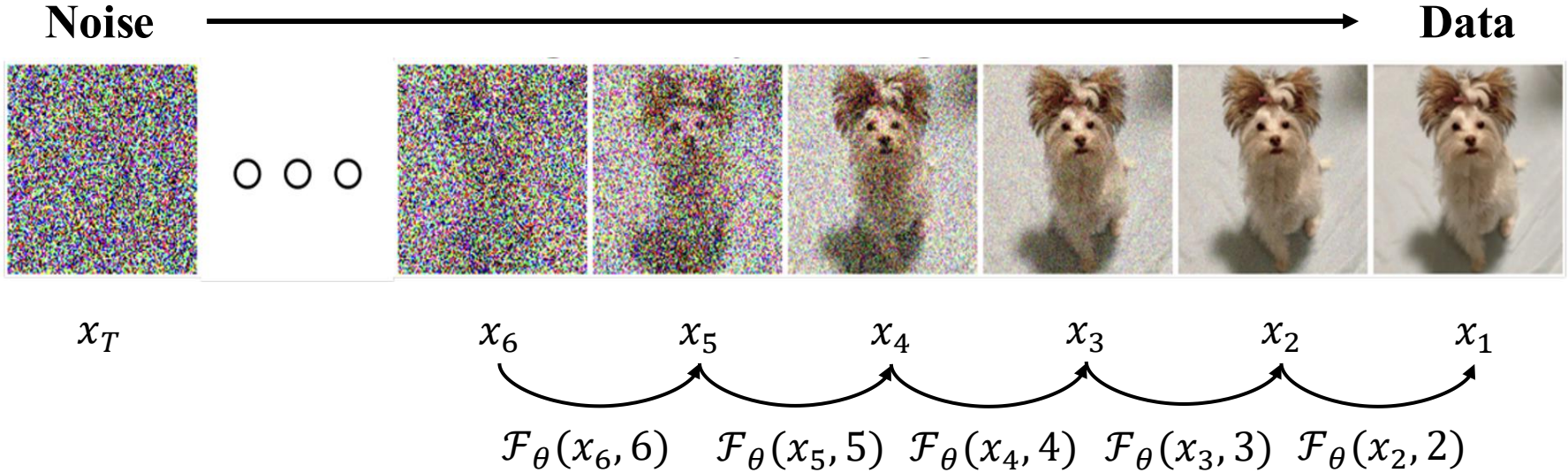


$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{data}(x), t \sim U(1, T), x_{t-1} \sim q(x_{t-1} | x_0), x_t \sim q(x_t | x_{t-1})} \|\mathcal{F}_\theta(x_t, t) - x_{t-1}\|_2^2$$

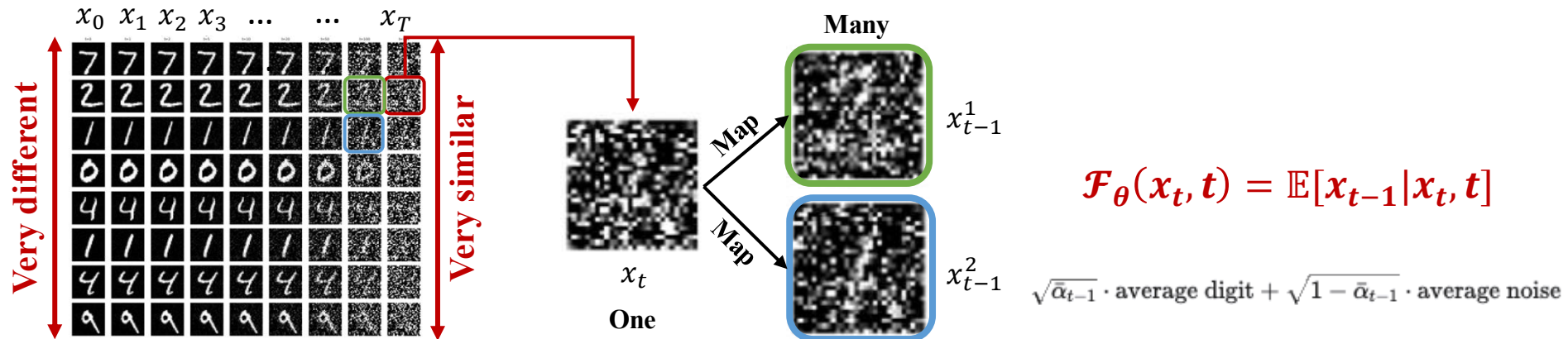


Generation - Problem

Generation by denoising step by step



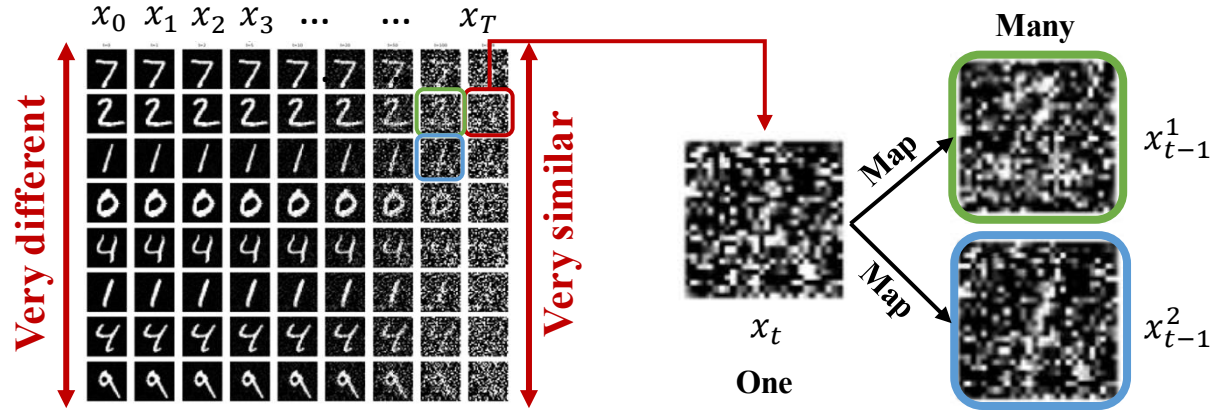
$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{\text{data}}(x), t \sim U(1, T), x_{t-1} \sim q(x_{t-1} | x_0), x_t \sim q(x_t | x_{t-1})} \|\mathcal{F}_\theta(x_t, t) - x_{t-1}\|_2^2$$





Generation - Problem

$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{\text{data}}(x), t \sim U(1, T), x_{t-1} \sim q(x_{t-1} | x_0), x_t \sim q(x_t | x_{t-1})} \| \mathcal{F}_\theta(x_t, t) - x_{t-1} \|_2^2$$



$$\mathcal{F}_\theta(x_t, t) = \mathbb{E}[x_{t-1} | x_t, t]$$

$$\frac{1}{K} \sum_i \epsilon^{(i)} \sim \mathcal{N}\left(0, \frac{1}{K} I\right)$$

$$\frac{1}{K} \sum_{i=1}^K x_{t-1}^{(i)} = \sqrt{\bar{\alpha}_{t-1}} x_0 + \sqrt{1 - \bar{\alpha}_{t-1}} \frac{1}{K} \sum_{i=1}^K \epsilon^{(i)}$$

early t:

average looks like a somewhat meaningful blurry digit

middle t:

average looks like a faint mixture of digit shapes

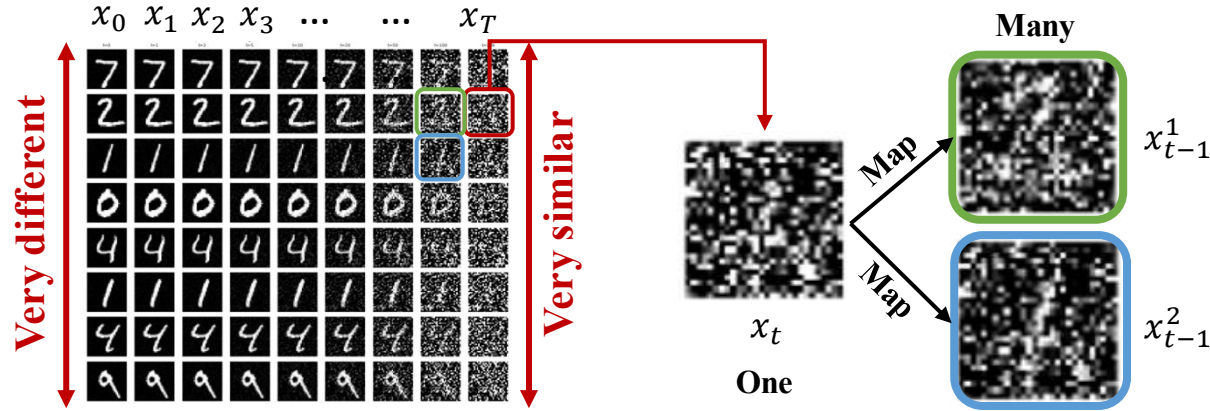
late t:

average approaches gray / zero, not a new Gaussian sample



Generation - Problem

$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{\text{data}}(x), t \sim U(1, T), x_{t-1} \sim q(x_{t-1} | x_0), x_t \sim q(x_t | x_{t-1})} \| \mathcal{F}_\theta(x_t, t) - x_{t-1} \|_2^2$$



$$\mathcal{F}_\theta(x_t, t) = \mathbb{E}[x_{t-1} | x_t, t]$$

$$\frac{1}{K} \sum_i \epsilon^{(i)} \sim \mathcal{N}\left(0, \frac{1}{K} I\right)$$

$$\frac{1}{K} \sum_{i=1}^K x_{t-1}^{(i)} = \sqrt{\bar{\alpha}_{t-1}} x_0 + \sqrt{1 - \bar{\alpha}_{t-1}} \frac{1}{K} \sum_{i=1}^K \epsilon^{(i)}$$

early t:

average looks like a somewhat meaningful blurry digit

middle t:

average looks like a faint mixture of digit shapes

late t:

average approaches gray / zero, not a new Gaussian sample



Generation - Problem

Clean

t=20



t=100



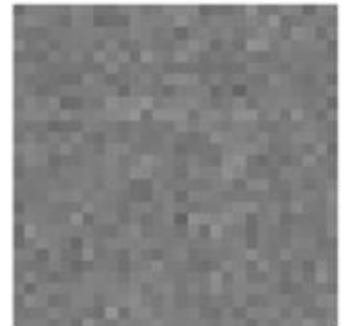
t=190



x_{t-1}



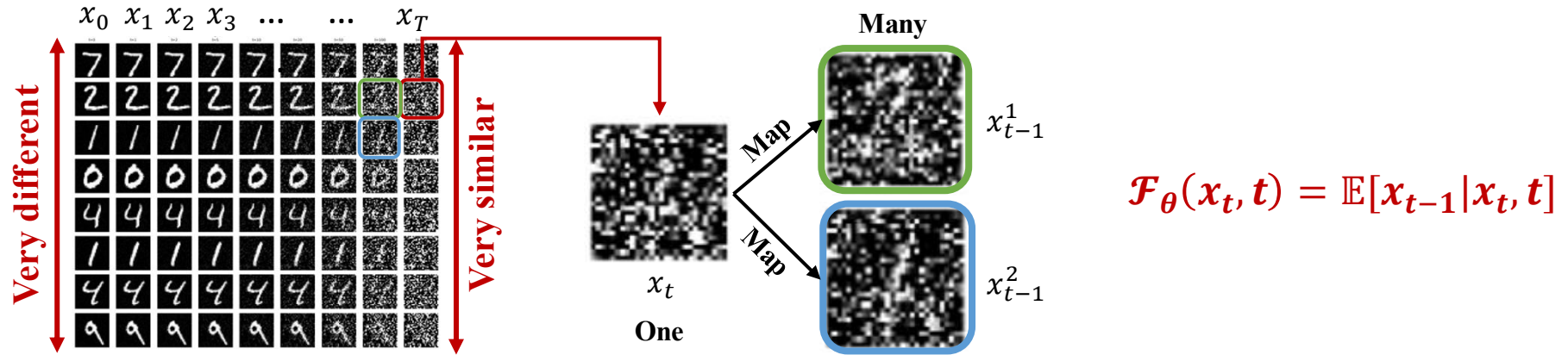
$$\sum_{i=1-256} x_{t-1}^i$$





Generation - Problem

$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{data}(x), t \sim U(1, T), x_{t-1} \sim q(x_{t-1} | x_0), x_t \sim q(x_t | x_{t-1})} \|\mathcal{F}_\theta(x_t, t) - x_{t-1}\|_2^2$$



Train a neural network to predict the nearly grey image



- **Motivation and Introduction of Diffusion**
- **Forward process of Diffusion**
- **Backward process of Diffusion**
 - **Backward progressive:** $x_t \rightarrow x_{t-1}$
 - **Backward one step:** $x_t \rightarrow x_0$
 - **Backward noise:** $x_t \rightarrow \epsilon$

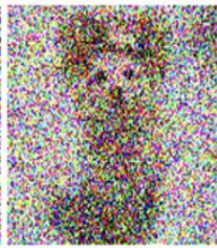
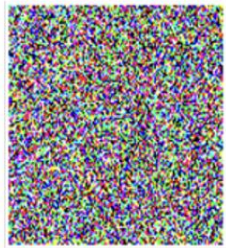


Backward Optimization 2

Generation by denoising step by step

Noise

Data



x_T

x_6

x_5

x_4

x_3

x_2

x_1

Predicted

Ground-truth

$$\|\hat{x} - x_{t-1}\|_2^2 = \|\mathcal{F}(x_t, t) - x_{t-1}\|_2^2$$

$$\|\hat{x} - \mathbf{x}_0\|_2^2 = \|\mathcal{F}(x_t, t) - \mathbf{x}_0\|_2^2$$

No randomness

x_0

x_t



$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{\text{data}}(x), t \sim U(1, T), x_t \sim q(x_t | x_0)} \|\mathcal{F}_\theta(x_t, t) - x_0\|_2^2$$

Given this noisy image and noise level, recover the clean digit.

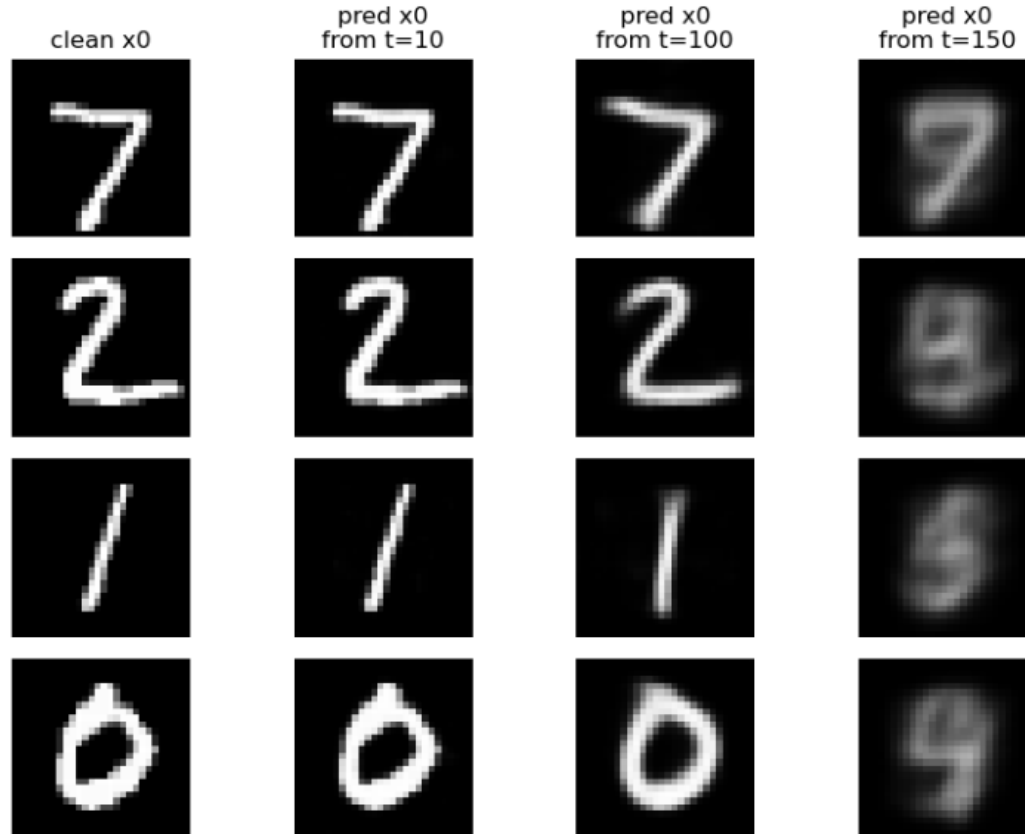


Backward Optimization 2

$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{\text{data}}(x), t \sim U(1, T), x_t \sim q(x_t | x_0)} \|\mathcal{F}_\theta(x_t, t) - x_0\|_2^2$$

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

$$\mathcal{F}_{\theta^*}(x_T, t) = \bar{x}_0$$



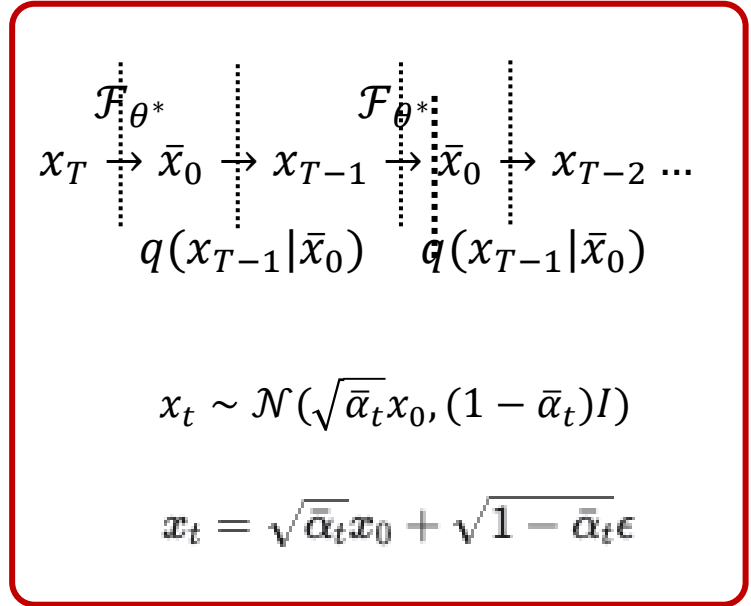
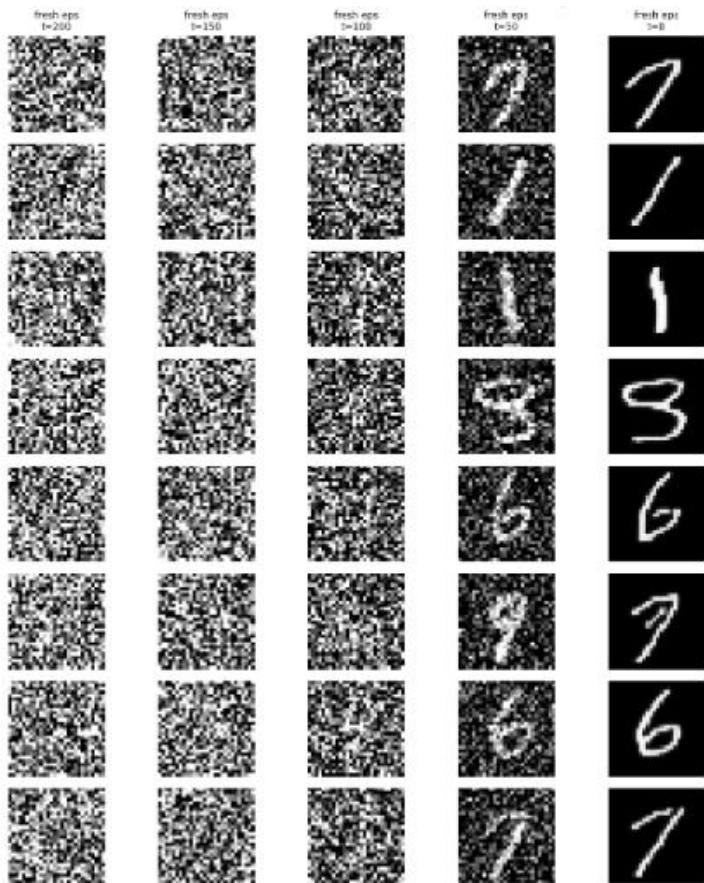
Predict from later time step gives more noisy images



Backward Optimization 2

$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{\text{data}}(x), t \sim U(1, T), x_t \sim q(x_t | x_0)} \|\mathcal{F}_\theta(x_t, t) - x_0\|_2^2$$

$$x_T \sim \mathcal{N}(0, I) \quad \mathcal{F}_{\theta^*}(x_T, t) = \bar{x}_0$$



Inconsistency

since

resample ϵ

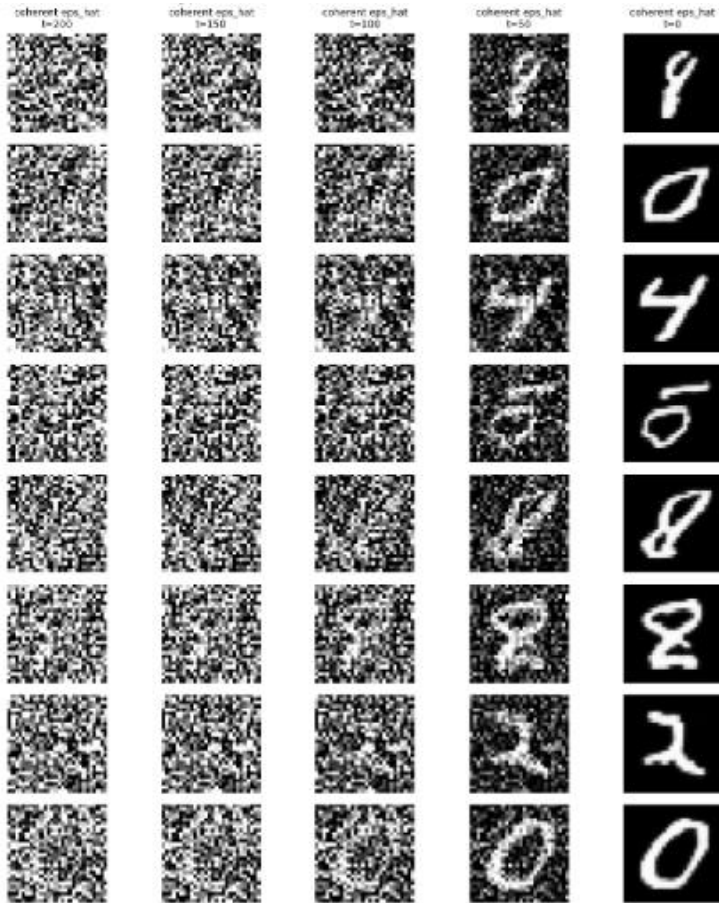


Backward Optimization 2

$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{\text{data}}(x), t \sim U(1, T), x_t \sim q(x_t | x_0)} \|\mathcal{F}_\theta(x_t, t) - x_0\|_2^2$$

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

$$\mathcal{F}_{\theta^*}(x_T, t) = \bar{x}_0$$



$$x_T \rightarrow \bar{x}_0 \rightarrow \epsilon \rightarrow x_{T-1} \rightarrow$$

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

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

$$\hat{\epsilon} = \frac{x_t - \sqrt{\bar{\alpha}_t} \hat{x}_0}{\sqrt{1 - \bar{\alpha}_t}}$$

$$x_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \hat{x}_0 + \sqrt{1 - \bar{\alpha}_{t-1}} \hat{\epsilon}$$

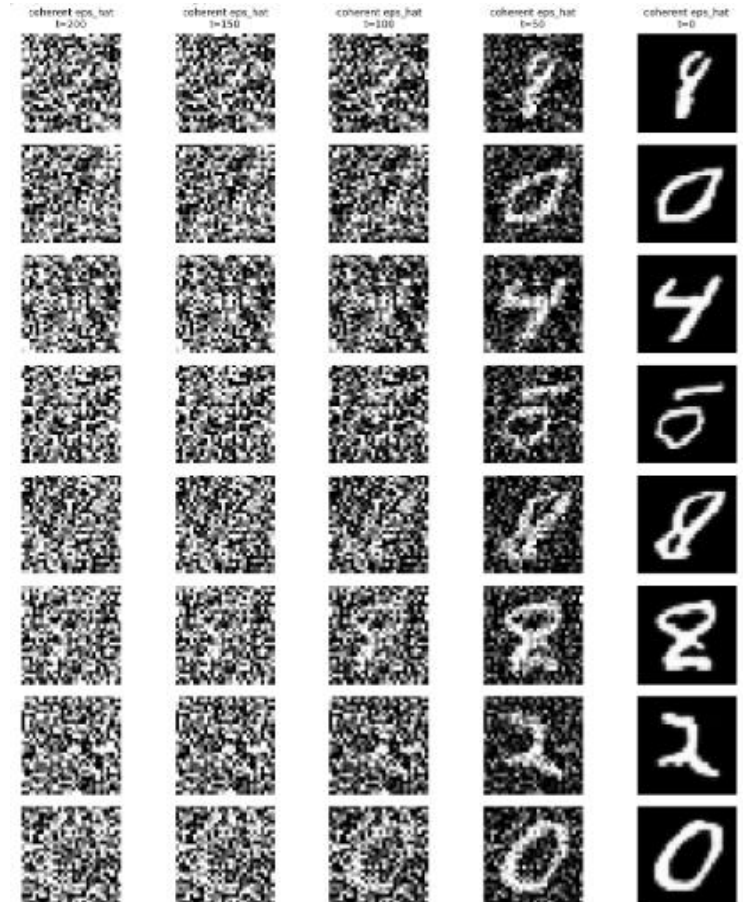
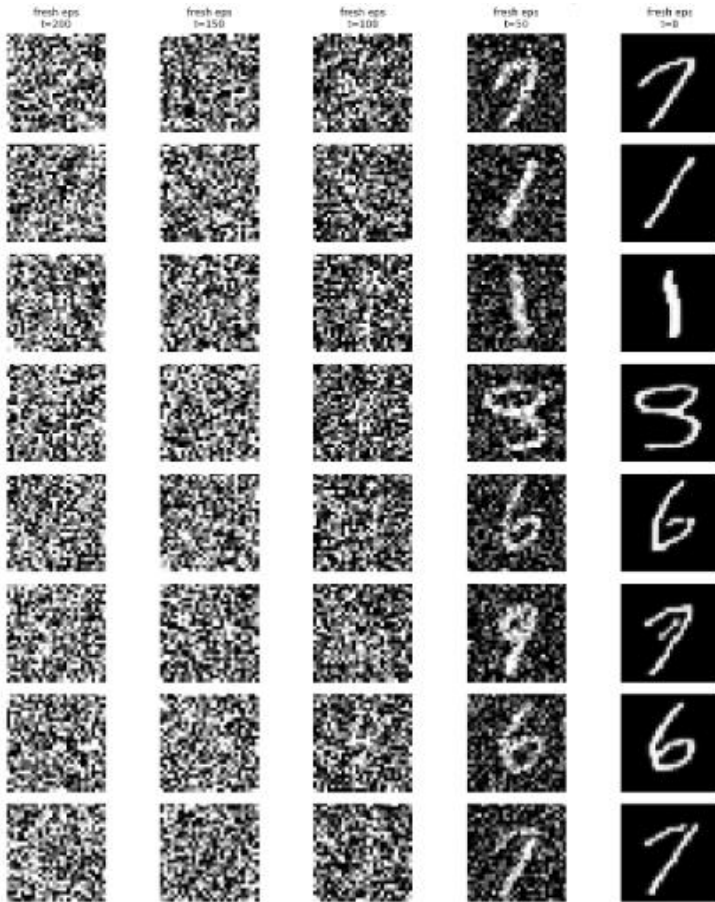
Consistency

since

Solve ϵ



Backward Optimization 2



$$x_T \rightarrow \bar{x}_0 \rightarrow x_{T-1} \rightarrow \bar{x}_0 \rightarrow x_{T-2} \dots$$

Inconsistency

$$x_T \rightarrow \bar{x}_0 \rightarrow \epsilon \rightarrow x_{T-1} \rightarrow$$

Consistency

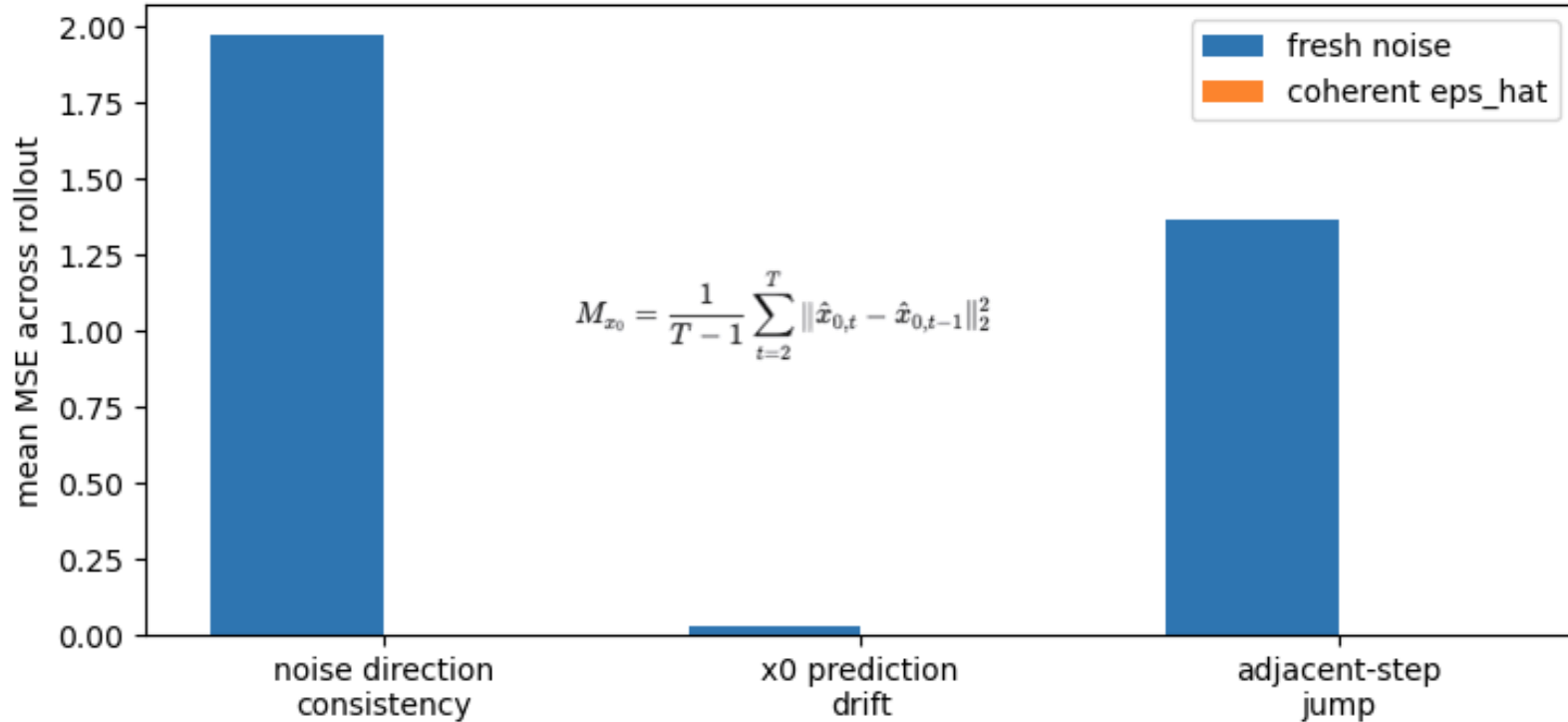


Backward Optimization 2

$$M_\epsilon = \frac{1}{T-1} \sum_{t=2}^T \|\hat{\epsilon}_{t-1}^{\text{implied}} - \hat{\epsilon}_t\|_2^2$$

$$M_{\Delta x} = \frac{1}{T} \sum_{t=1}^T \|x_{t-1} - x_t\|_2^2$$

Empirical trajectory advantage of the coherent update



$$M_{x_0} = \frac{1}{T-1} \sum_{t=2}^T \|\hat{x}_{0,t} - \hat{x}_{0,t-1}\|_2^2$$

$x_T \rightarrow \bar{x}_0 \rightarrow x_{T-1} \rightarrow \bar{x}_0 \rightarrow x_{T-2} \dots$

$x_T \rightarrow \bar{x}_0 \rightarrow \epsilon \rightarrow x_{T-1} \rightarrow$

Both are good

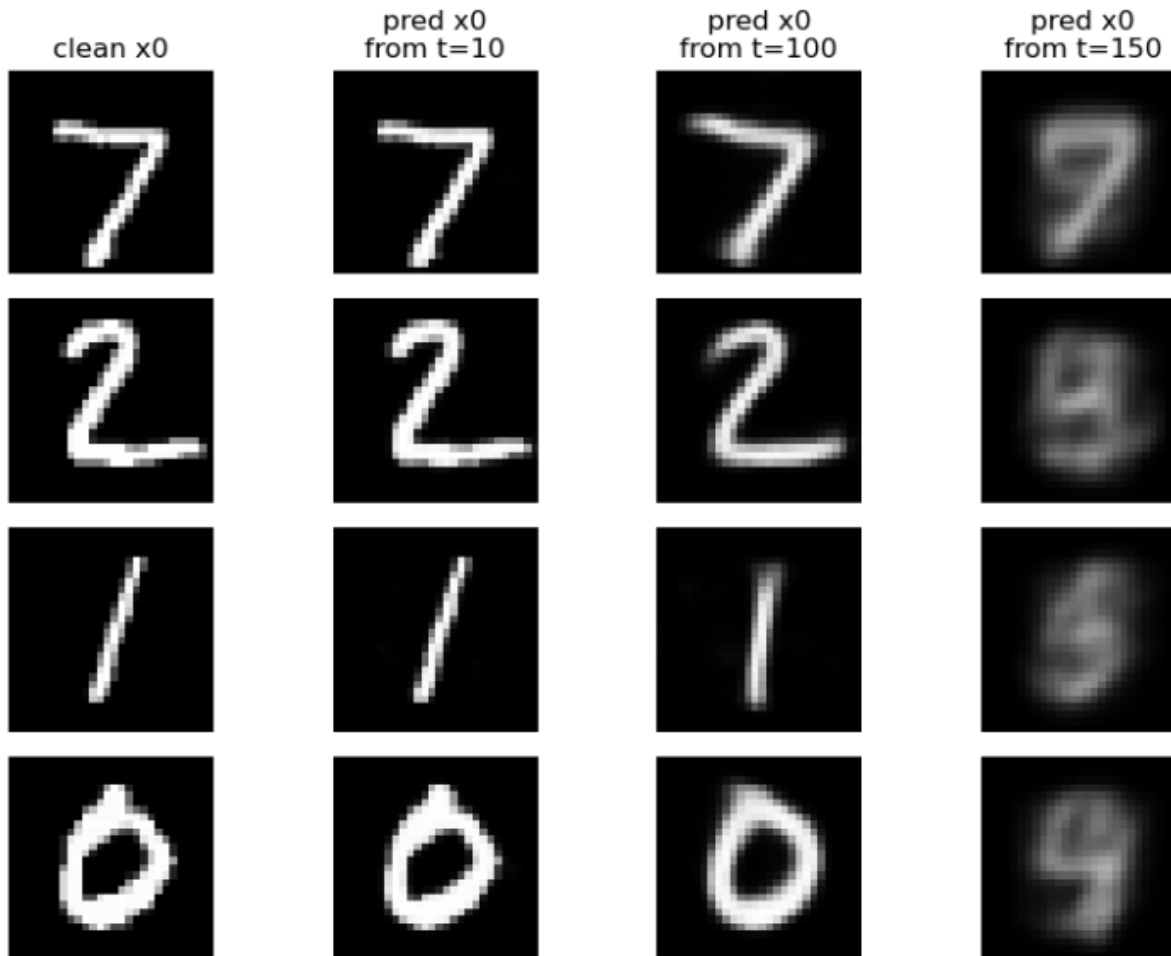


Backward Optimization 2 - Problem

$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{\text{data}}(x), t \sim U(1, T), x_t \sim q(x_t | x_0)} \|\mathcal{F}_\theta(x_t, t) - x_0\|_2^2$$

t small, easy

t large, easy



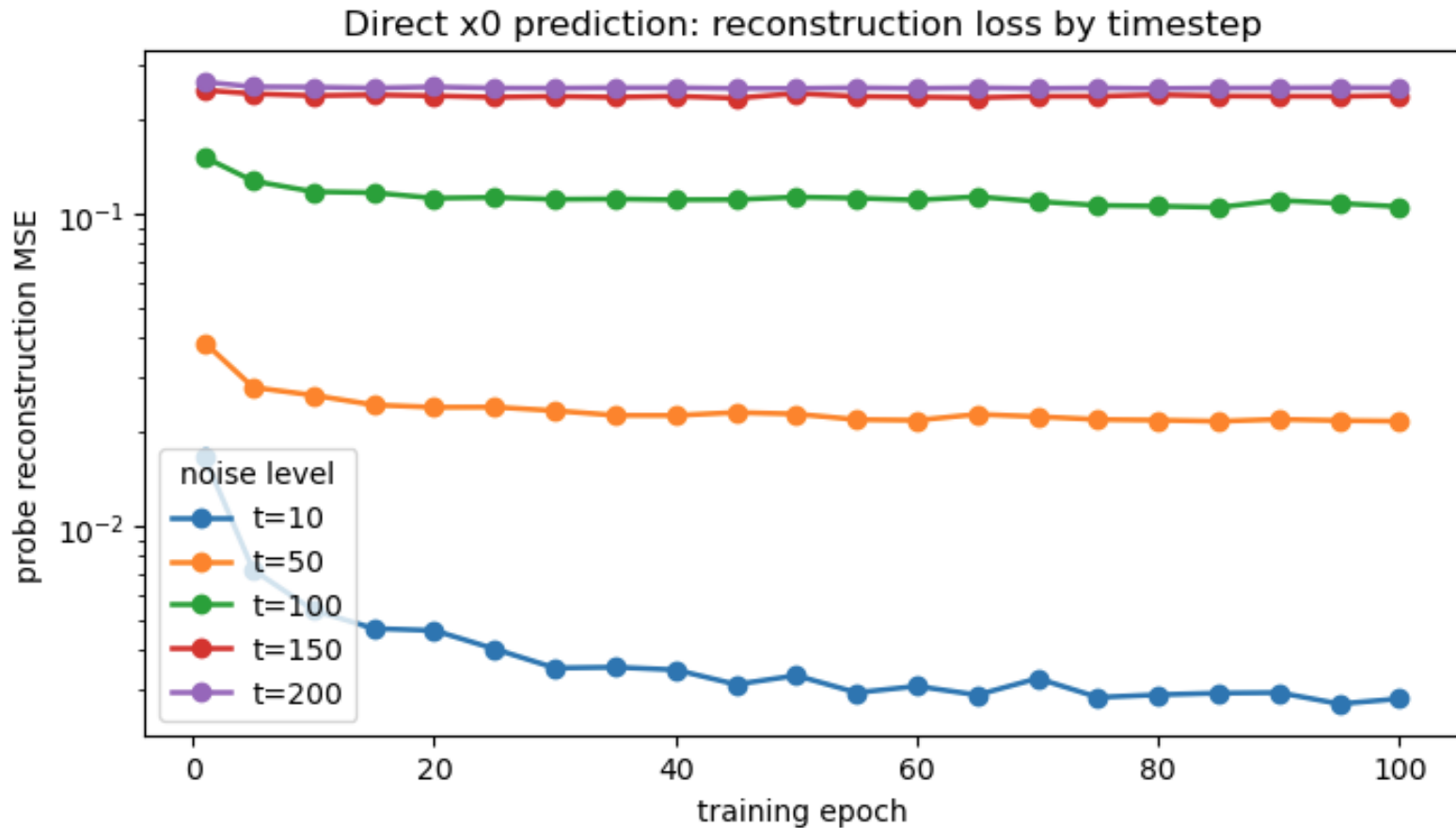


Backward Optimization 2 - Problem

$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{\text{data}}(x), t \sim U(1, T), x_t \sim q(x_t | x_0)} \|\mathcal{F}_\theta(x_t, t) - x_0\|_2^2$$

t small, easy

t large, easy





- **Motivation and Introduction of Diffusion**
- **Forward process of Diffusion**
- **Backward process of Diffusion**
 - **Backward progressive:** $x_t \rightarrow x_{t-1}$
 - **Backward one step:** $x_t \rightarrow x_0$
 - **Backward noise:** $x_t \rightarrow \epsilon$



Backward Optimization 3

$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{\text{data}}(x), t \sim U(1, T), x_t \sim q(x_t | x_0)} \|\mathcal{F}_\theta(x_t, t) - x_0\|_2^2$$

$$\mathcal{F}_{\theta^*}(x_T, t) = \bar{x}_0 \quad \longrightarrow \quad \hat{\epsilon} = \frac{x_t - \sqrt{\bar{\alpha}_t} \hat{x}_0}{\sqrt{1 - \bar{\alpha}_t}}$$

Can we instead directly predict noise?

$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{\text{data}}(x), t \sim U(1, T), \epsilon \sim \mathcal{N}(0, I)} \|\epsilon_\theta(x_t, t) - \epsilon\|_2^2$$

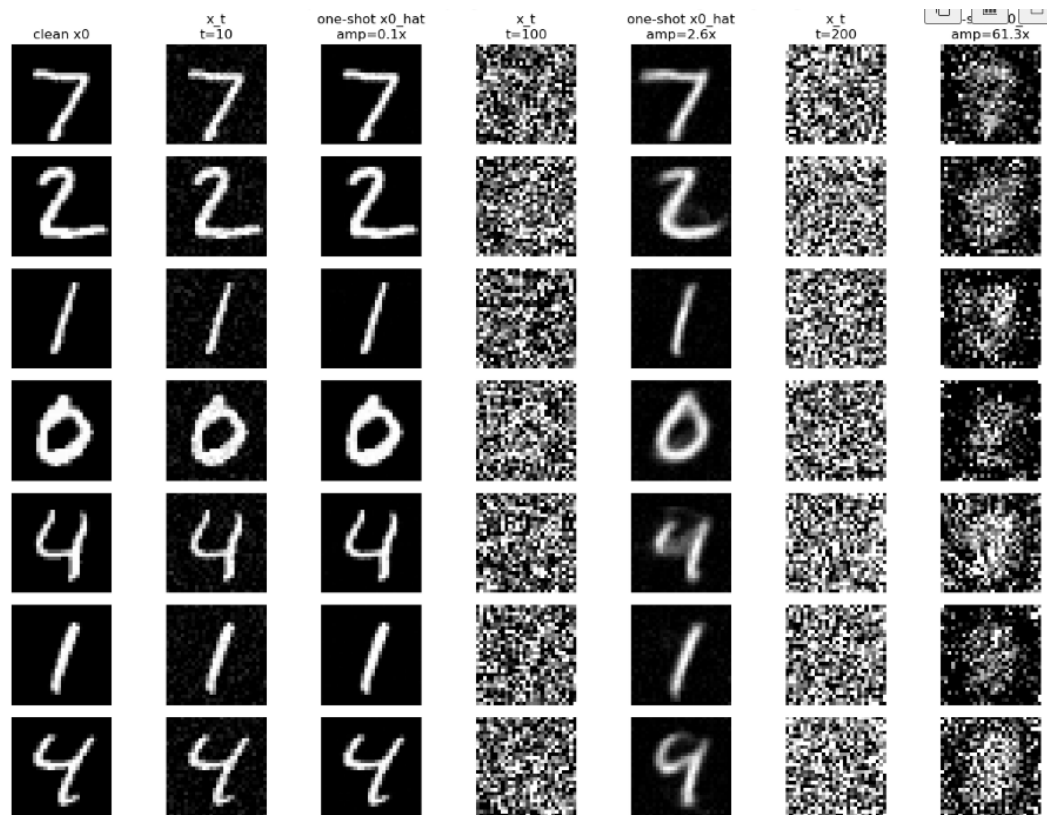
$$\hat{x}_0 = \frac{x_t - \sqrt{1 - \bar{\alpha}_t} \hat{\epsilon}}{\sqrt{\bar{\alpha}_t}}$$



Backward Optimization 3

$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{\text{data}}(x), t \sim U(1, T), \epsilon \sim \mathcal{N}(0, I)} \|\epsilon_\theta(x_t, t) - \epsilon\|_2^2$$

$$x_T \sim \mathcal{N}(0, I) \quad \hat{\epsilon} = \epsilon_{\theta^*}(x_t, t) \quad \rightarrow \quad \hat{x}_0 = \frac{x_t - \sqrt{1 - \bar{\alpha}_t} \hat{\epsilon}}{\sqrt{\bar{\alpha}_t}}$$





Backward Optimization 3

$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{\text{data}}(x), t \sim U(1, T), \epsilon \sim \mathcal{N}(0, I)} \|\epsilon_\theta(x_t, t) - \epsilon\|_2^2$$

$$x_T \sim \mathcal{N}(0, I) \quad \hat{\epsilon} = \epsilon_{\theta^*}(x_t, t)$$

$$\hat{x}_0 = \frac{x_t - \sqrt{1 - \bar{\alpha}_t} \hat{\epsilon}}{\sqrt{\bar{\alpha}_t}}$$

$$x_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \hat{x}_0 + \sqrt{1 - \bar{\alpha}_{t-1}} \hat{\epsilon}$$

Deterministic epsilon rollout snapshots

t=200

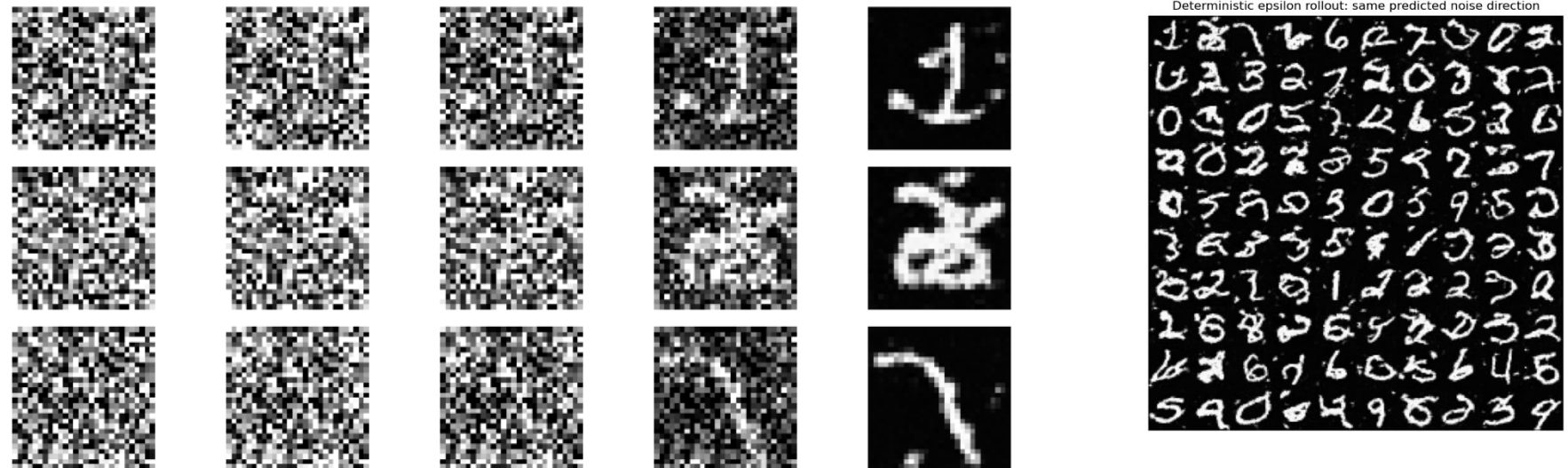
t=150

t=100

t=50

t=0

Deterministic epsilon rollout: same predicted noise direction





Backward Optimization 3

$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{data}(x), t \sim U(1, T), \epsilon \sim \mathcal{N}(0, I)} \|\epsilon_\theta(x_t, t) - \epsilon\|_2^2$$

$$x_T \sim \mathcal{N}(0, I) \quad \hat{\epsilon} = \epsilon_{\theta^*}(x_t, t)$$

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

$$\hat{x}_0 = \frac{x_t - \sqrt{1 - \bar{\alpha}_t} \hat{\epsilon}}{\sqrt{\bar{\alpha}_t}}.$$

$$q(x_{t-1} | x_t, x_0) = \mathcal{N}(\tilde{\mu}_t(x_t, x_0), \tilde{\beta}_t I).$$

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

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

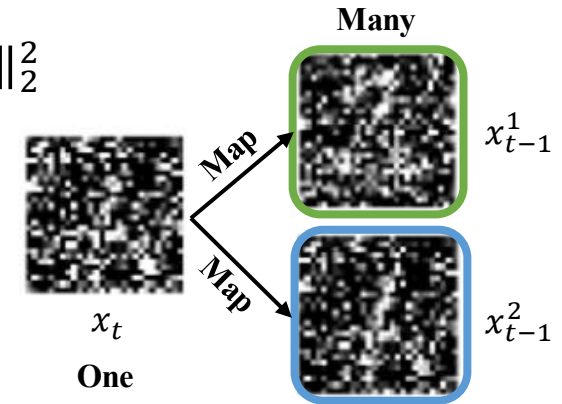
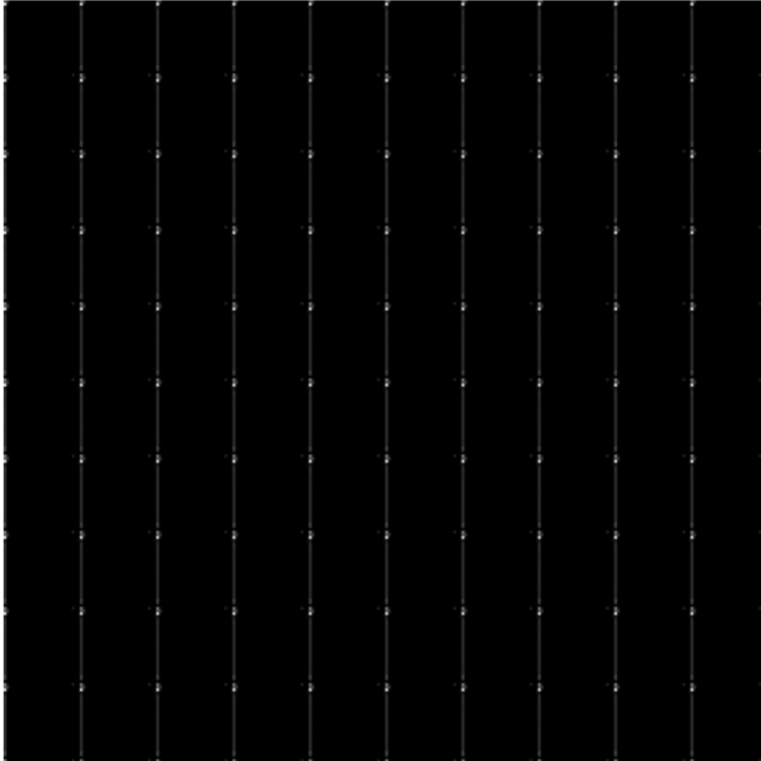
$$x_{t-1} = \mu_\theta(x_t, t) + \sqrt{\tilde{\beta}_t} z, \quad z \sim \mathcal{N}(0, I)$$



Summary of three backward optimization - 1

$$\mathbb{E}_{x_0 \sim p_{data}(x), t \sim U(1, T), x_{t-1} \sim q(x_{t-1} | x_0), x_t \sim q(x_t | x_{t-1})} \| \mathcal{F}(x_t, t) - x_{t-1} \|_2^2$$

Naive generation (deterministic): repeatedly predict x_{t-1}



Forward Process

cat \longrightarrow c#a\$t%

Backward Process

cat \longleftarrow c#a\$t%

What was the exact previous corrupted string?
Was it "c#a\$t%", "ca\$t%", "c#a\$t%", ... ?

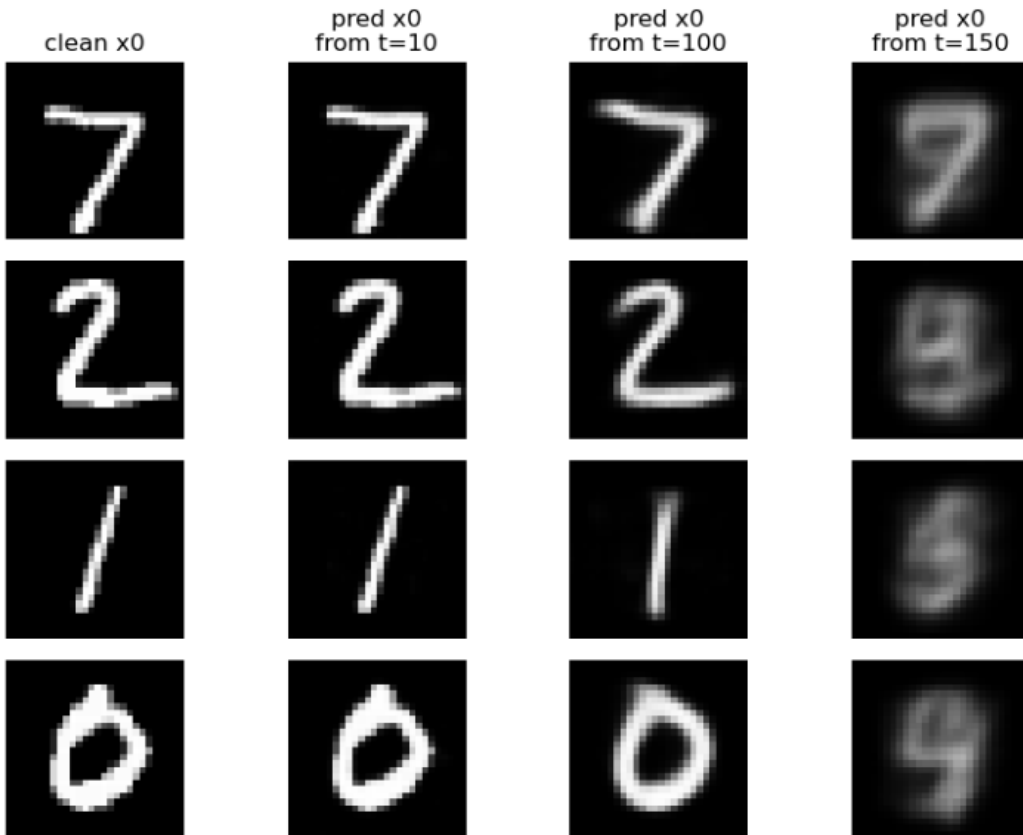
$$x_t \rightarrow x_{t-1}$$



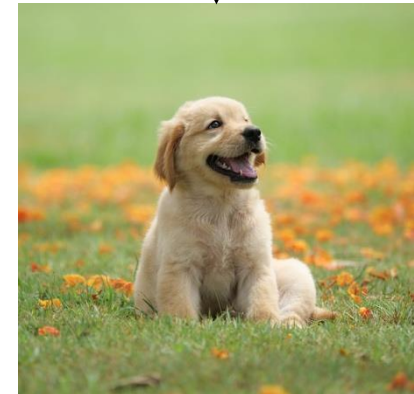
Summary of three backward optimization - 2

$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{data}(x), t \sim U(1, T), x_t \sim q(x_t | x_0)} \|\mathcal{F}_\theta(x_t, t) - x_0\|_2^2$$

$$x_T \sim \mathcal{N}(0, I) \quad \mathcal{F}_{\theta^*}(x_T, t) = \bar{x}_0$$



One-shot
Generation

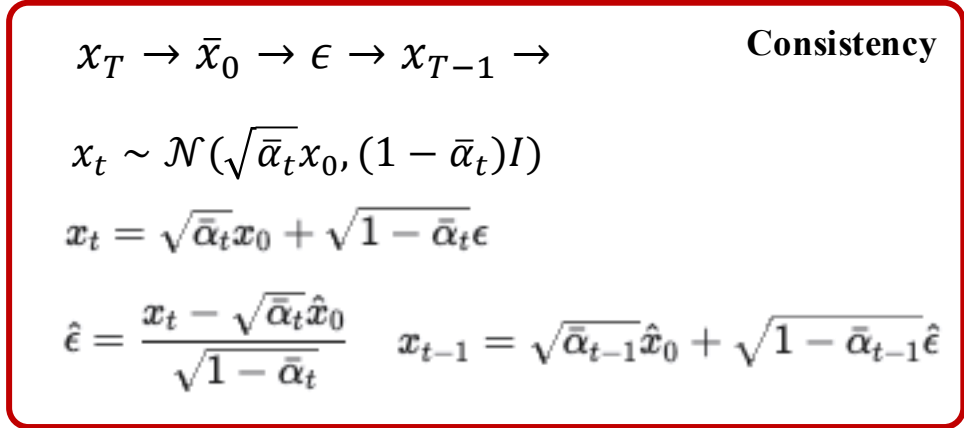
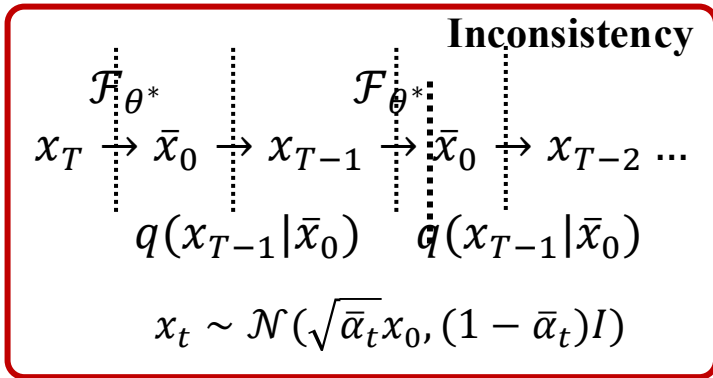




Summary of three backward optimization - 2

$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{\text{data}}(x), t \sim U(1, T), x_t \sim q(x_t | x_0)} \|\mathcal{F}_\theta(x_t, t) - x_0\|_2^2$$

$$x_T \sim \mathcal{N}(0, I) \quad \mathcal{F}_{\theta^*}(x_T, t) = \bar{x}_0$$

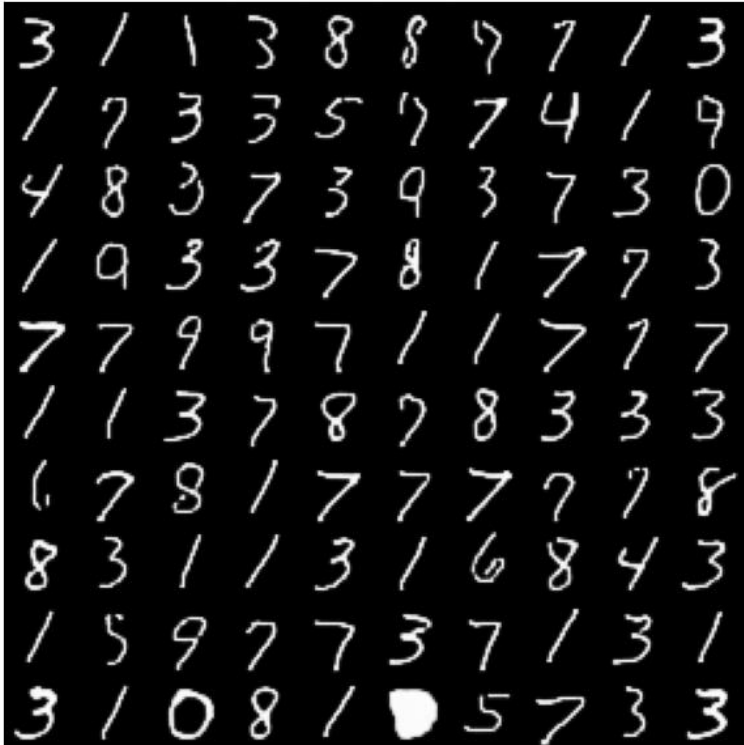




Summary of three backward optimization - 2

$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{\text{data}}(x), t \sim U(1, T), x_t \sim q(x_t | x_0)} \|\mathcal{F}_\theta(x_t, t) - x_0\|_2^2$$

$$x_T \sim \mathcal{N}(0, I) \quad \mathcal{F}_{\theta^*}(x_T, t) = \bar{x}_0$$

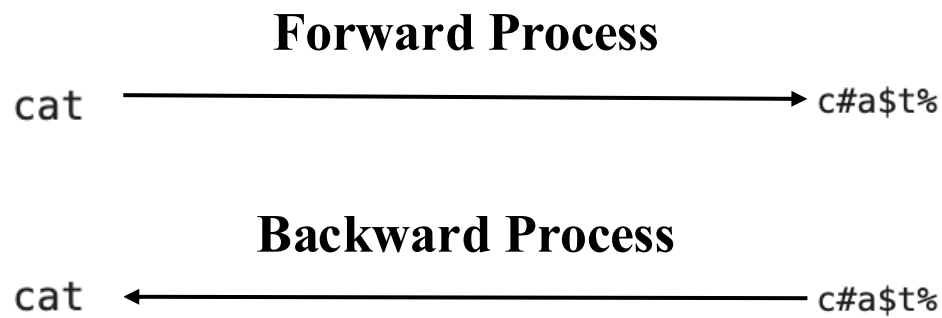




Summary of three backward optimization - 2

$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{\text{data}}(x), t \sim U(1, T), x_t \sim q(x_t | x_0)} \|\mathcal{F}_\theta(x_t, t) - x_0\|_2^2$$

$$x_T \sim \mathcal{N}(0, I) \quad \mathcal{F}_{\theta^*}(x_T, t) = \bar{x}_0$$



What was the clean string?

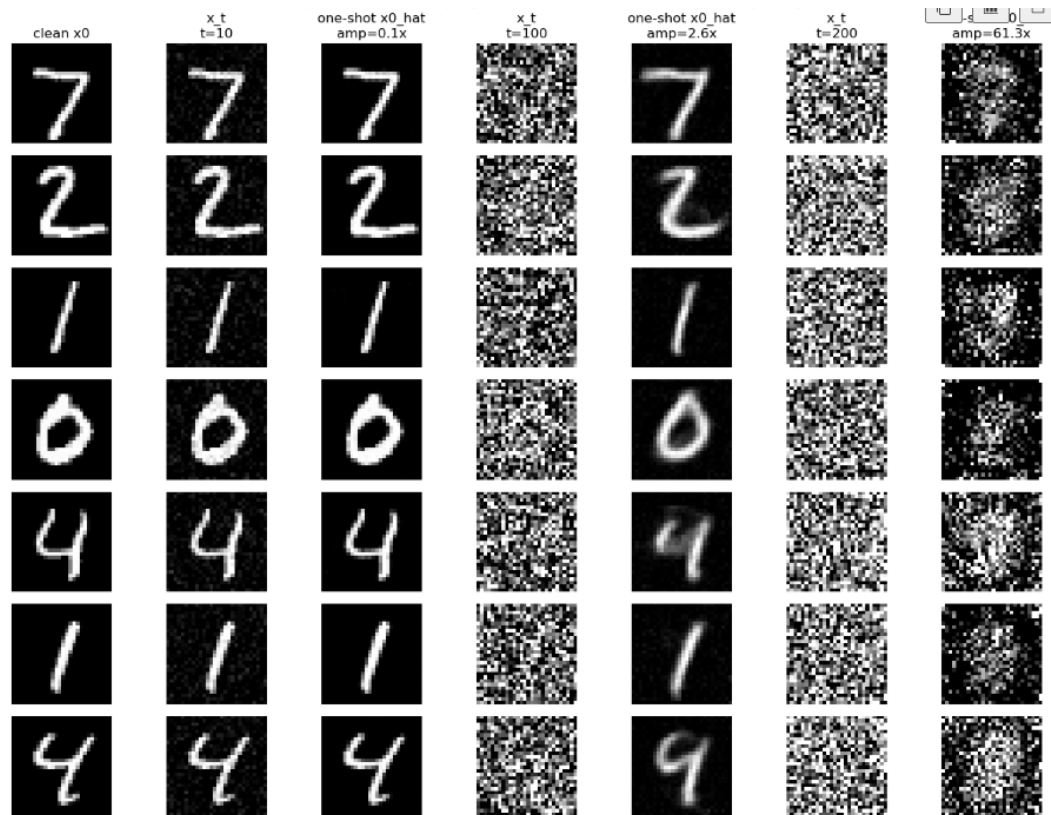
$$x_t \rightarrow x_0$$



Summary of three backward optimization - 3

$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{\text{data}}(x), t \sim U(1, T), \epsilon \sim \mathcal{N}(0, I)} \|\epsilon_{\theta}(x_t, t) - \epsilon\|_2^2$$

$$x_T \sim \mathcal{N}(0, I) \quad \hat{\epsilon} = \epsilon_{\theta^*}(x_t, t) \quad \rightarrow \quad \hat{x}_0 = \frac{x_t - \sqrt{1 - \bar{\alpha}_t} \hat{\epsilon}}{\sqrt{\bar{\alpha}_t}}$$





Summary of three backward optimization - 3

$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{\text{data}}(x), t \sim U(1, T), \epsilon \sim \mathcal{N}(0, I)} \|\epsilon_{\theta}(x_t, t) - \epsilon\|_2^2$$

$$x_T \sim \mathcal{N}(0, I) \quad \hat{\epsilon} = \epsilon_{\theta^*}(x_t, t) \quad \longrightarrow \quad \hat{x}_0 = \frac{x_t - \sqrt{1 - \bar{\alpha}_t} \hat{\epsilon}}{\sqrt{\bar{\alpha}_t}}$$

$$x_T \sim \mathcal{N}(0, I) \quad \hat{\epsilon} = \epsilon_{\theta^*}(x_t, t)$$

$$\hat{x}_0 = \frac{x_t - \sqrt{1 - \bar{\alpha}_t} \hat{\epsilon}}{\sqrt{\bar{\alpha}_t}}$$

Introduce some noise

$$x_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \hat{x}_0 + \sqrt{1 - \bar{\alpha}_{t-1}} \hat{\epsilon}$$

$$q(x_{t-1} | x_t, x_0) = \mathcal{N}(\bar{\mu}_t(x_t, x_0), \bar{\beta}_t I).$$

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

$$x_{t-1} = \mu_{\theta}(x_t, t) + \sqrt{\bar{\beta}_t} z, \quad z \sim \mathcal{N}(0, I)$$



Summary of three backward optimization - 3

$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{\text{data}}(x), t \sim U(1, T), \epsilon \sim \mathcal{N}(0, I)} \|\epsilon_{\theta}(x_t, t) - \epsilon\|_2^2$$

$$x_T \sim \mathcal{N}(0, I) \quad \hat{\epsilon} = \epsilon_{\theta^*}(x_t, t) \quad \longrightarrow \quad \hat{x}_0 = \frac{x_t - \sqrt{1 - \bar{\alpha}_t} \hat{\epsilon}}{\sqrt{\bar{\alpha}_t}}$$

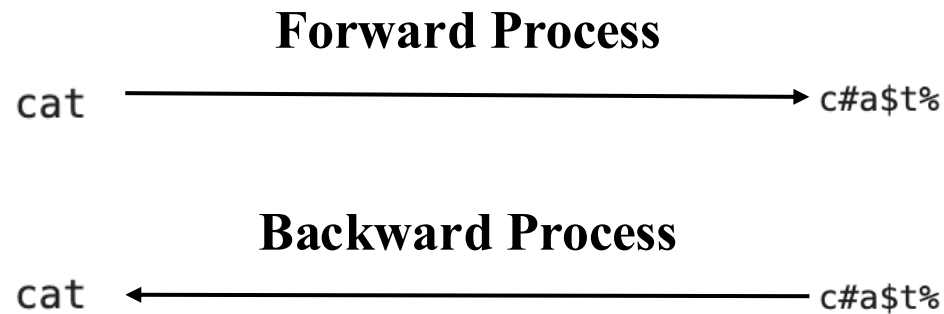




Summary of three backward optimization - 3

$$\theta^* = \operatorname{argmin} \mathbb{E}_{x_0 \sim p_{\text{data}}(x), t \sim U(1, T), \epsilon \sim \mathcal{N}(0, I)} \|\epsilon_{\theta}(x_t, t) - \epsilon\|_2^2$$

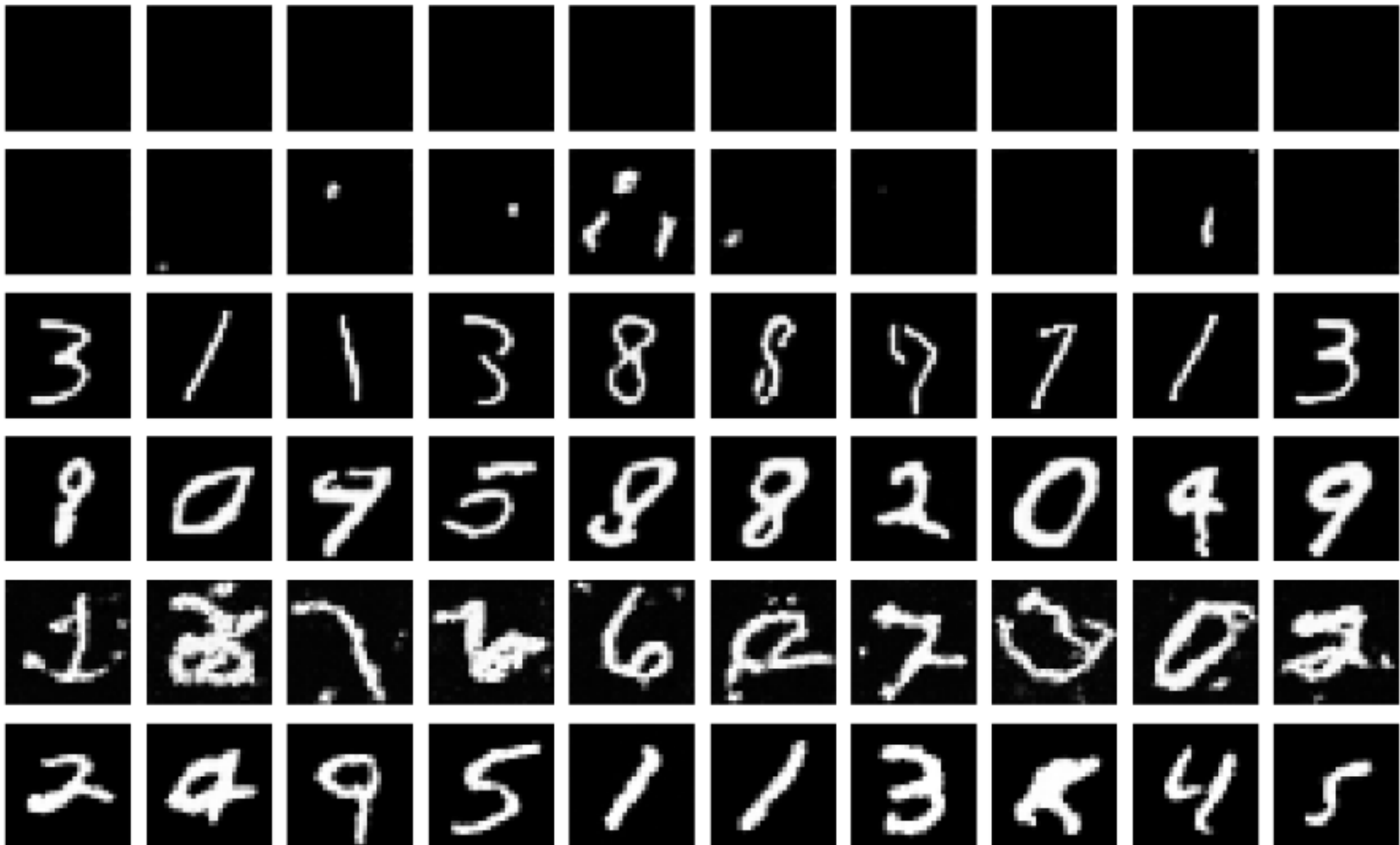
$$x_T \sim \mathcal{N}(0, I) \quad \hat{\epsilon} = \epsilon_{\theta^*}(x_t, t)$$



Which symbols look like noise?

Answer: #, \$, %

$$x_t \rightarrow \sigma$$





Summary

method	NN-MSE	divMSE	div/real	meanGap	stdGap	cover	entropy
A1 xprev det	0.1223	0.0000	0.000	0.3298	0.5671	1	-0.000
A1 xprev +noise	0.1557	0.0874	0.171	0.3088	0.3843	2	0.141
A2 x0 fresh	0.1299	0.2566	0.501	0.1117	0.1677	4	0.530
A2 x0 coherent	0.2295	0.5247	1.025	0.0571	0.0554	6	0.759
A3 eps det	0.3334	0.6679	1.304	0.1542	0.0730	5	0.676
A3 eps DDPM	0.1930	0.5187	1.013	0.0194	0.0262	7	0.819

Nearest-real MSE measures whether generated images lie close to the real MNIST manifold:

$$M_{NN} = \frac{1}{N} \sum_{i=1}^N \min_{1 \leq j \leq M} \frac{1}{D} \|\hat{x}_i - x_j^{\text{real}}\|_2^2.$$



Summary

method	NN-MSE	divMSE	div/real	meanGap	stdGap	cover	entropy
A1 xprev det	0.1223	0.0000	0.000	0.3298	0.5671	1	-0.000
A1 xprev +noise	0.1557	0.0874	0.171	0.3088	0.3843	2	0.141
A2 x0 fresh	0.1299	0.2566	0.501	0.1117	0.1677	4	0.530
A2 x0 coherent	0.2295	0.5247	1.025	0.0571	0.0554	6	0.759
A3 eps det	0.3334	0.6679	1.304	0.1542	0.0730	5	0.676
A3 eps DDPM	0.1930	0.5187	1.013	0.0194	0.0262	7	0.819

Diversity MSE measures spread among generated samples:

$$M_{\text{div}} = \frac{2}{N(N-1)} \sum_{1 \leq i < k \leq N} \frac{1}{D} \|\hat{x}_i - \hat{x}_k\|_2^2.$$

Too low suggests collapse; too high suggests noisy/unstructured samples. The notebook also reports

$$R_{\text{div}} = \frac{M_{\text{div}}}{M_{\text{div}}^{\text{real}}},$$

where values near 1 mean the generated diversity roughly matches the real reference batch.



Summary

method	NN-MSE	divMSE	div/real	meanGap	stdGap	cover	entropy
A1 xprev det	0.1223	0.0000	0.000	0.3298	0.5671	1	-0.000
A1 xprev +noise	0.1557	0.0874	0.171	0.3088	0.3843	2	0.141
A2 x0 fresh	0.1299	0.2566	0.501	0.1117	0.1677	4	0.530
A2 x0 coherent	0.2295	0.5247	1.025	0.0571	0.0554	6	0.759
A3 eps det	0.3334	0.6679	1.304	0.1542	0.0730	5	0.676
A3 eps DDPM	0.1930	0.5187	1.013	0.0194	0.0262	7	0.819

Pixel mean gap compares global brightness:

$$M_{\text{mean}} = \left| \frac{1}{ND} \sum_{i=1}^N \sum_{d=1}^D \hat{x}_{i,d} - \frac{1}{MD} \sum_{j=1}^M \sum_{d=1}^D x_{j,d}^{\text{real}} \right|.$$

Lower means the generated images have a similar average intensity to real MNIST.



Summary

method	NN-MSE	divMSE	div/real	meanGap	stdGap	cover	entropy
A1 xprev det	0.1223	0.0000	0.000	0.3298	0.5671	1	-0.000
A1 xprev +noise	0.1557	0.0874	0.171	0.3088	0.3843	2	0.141
A2 x0 fresh	0.1299	0.2566	0.501	0.1117	0.1677	4	0.530
A2 x0 coherent	0.2295	0.5247	1.025	0.0571	0.0554	6	0.759
A3 eps det	0.3334	0.6679	1.304	0.1542	0.0730	5	0.676
A3 eps DDPM	0.1930	0.5187	1.013	0.0194	0.0262	7	0.819

Centroid coverage and entropy are a lightweight class-spread diagnostic. Let c_k be the centroid of real images with digit label k . Assign each generated sample to the nearest centroid:

$$a_i = \arg \min_{k \in \{0, \dots, 9\}} \frac{1}{D} \|\hat{x}_i - c_k\|_2^2.$$

Coverage is the number of centroids used:

$$C = |\{k : \exists i, a_i = k\}|.$$

Entropy measures how evenly samples spread across these centroid assignments:

$$H = -\frac{1}{\log 10} \sum_{k=0}^9 p_k \log p_k, \quad p_k = \frac{1}{N} \sum_{i=1}^N \mathbf{1}[a_i = k].$$

Higher coverage and entropy suggest better digit-class diversity, but this is only a rough proxy, not a trained classifier evaluation.

