

# 生成模型

Diffusion理论发展与应用

常一帆

# Diffusion

## 理论基础

- DDPM
- DDIM
- SDE
- LDM
- Classifier-guided
- Classifier-free

## 微调

- Dreambooth
- LORA
- Textual inversion

## 细粒度控制

- controlnet

## 应用

- Repaint
- Blended-diffusion
- Diffedit
- P2p
- Pnp
- InstructPix2Pix
- Zero-classifier

CVPR24-oral

- 精准控制
- 复杂场景
- ...

主要调研图像生成（文生图/图片编辑）的研究，  
暂时略过了视频生成等方向

# Generative Models

- Given observed samples  $\textcolor{orange}{x}$  from a distribution of interest,
- The goal of a generative model is to learn to *model* its true data distribution  $p(\textcolor{orange}{x})$

# Generative Models

- ELBO(Evidence Lower Bound)
- VAE
- Hierarchical VAE

# ELBO(Evidence Lower Bound)

- Actually, we try to learn an associated unseen *latent* variable,  $\textcolor{orange}{z}$
- Model  $\textcolor{orange}{x}$  and  $\textcolor{orange}{z}$  by a joint distribution  $p(\textcolor{orange}{x}, \textcolor{orange}{z})$
- Maximize the likelihood  $p(\textcolor{orange}{x})$  of all observed  $\textcolor{orange}{x}$ , termed “likelihood-based”

# ELBO(Evidence Lower Bound)

- Maximize the likelihood  $p(\mathbf{x})$  of all observed  $\mathbf{x}$  ???

$$\begin{aligned}
 \log p(\mathbf{x}) &= \log p(\mathbf{x}) \int q_{\phi}(\mathbf{z}|\mathbf{x}) d\mathbf{z} \\
 &= \int q_{\phi}(\mathbf{z}|\mathbf{x})(\log p(\mathbf{x})) d\mathbf{z} \\
 &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x})] \\
 &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{z}|\mathbf{x})} \right] \\
 &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z}) q_{\phi}(\mathbf{z}|\mathbf{x})}{p(\mathbf{z}|\mathbf{x}) q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \\
 &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] + \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{p(\mathbf{z}|\mathbf{x})} \right] \\
 &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] + D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x})) \\
 &\geq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right]
 \end{aligned}$$

(Multiply by 1 =  $\int q_{\phi}(\mathbf{z}|\mathbf{x}) d\mathbf{z}$ )  
(Bring evidence into integral)  
(Definition of Expectation)  
(Apply Equation 2)  
(Multiply by 1 =  $\frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{q_{\phi}(\mathbf{z}|\mathbf{x})}$ )  
(Split the Expectation)  
(Definition of KL Divergence)  
(KL Divergence always  $\geq 0$ )

$$\mathbb{E}(X) = \int_{-\infty}^{\infty} x f(x) dx.$$

$$\mathbb{E}(g(X)) = \int_{\Omega} g(x) f(x) dx$$

$$p(\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{z}|\mathbf{x})}$$

$$\mathbb{D}_{\text{KL}}(P||Q) = \int_{-\infty}^{\infty} p(\mathbf{x}) \ln\left(\frac{p(\mathbf{x})}{q(\mathbf{x})}\right) d\mathbf{x}$$

# KL散度非负

## 1. 非负性

$\mathbb{D}_{\text{KL}}(P||Q) \geq 0$ ,  $\mathbb{D}_{\text{KL}} = 0$ 当且仅当 $P = Q$ 。

证明 (我们仅对离散情况进行证明, 对于连续随机变量情况, 我们将积分化为求和的极限后可以用相同方式证明) :

我们只需要证明 $\sum_i P(i) \ln\left(\frac{Q(i)}{P(i)}\right) \leq 0$ 。采用不等式 $\ln(x) \leq x - 1, \forall x > 0$ , 则:

$$\sum_i P(i) \ln\left(\frac{Q(i)}{P(i)}\right) \leq \sum_i P(i)\left(\frac{Q(i)}{P(i)} - 1\right) = 0$$

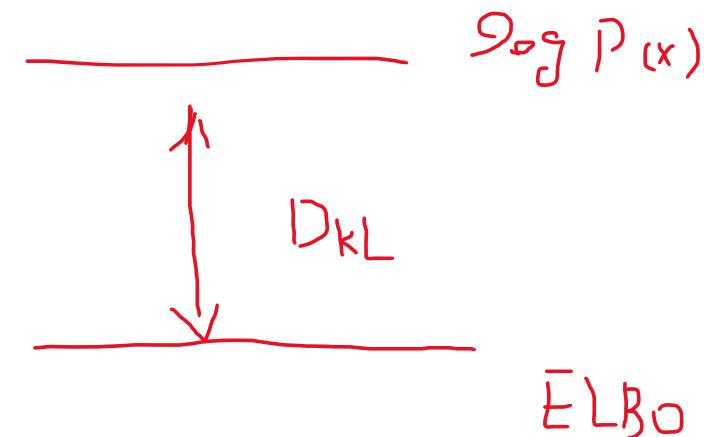
等号当且仅当对于任意的*i*,  $\frac{Q(i)}{P(i)} = 1$ 时取得, 此时有 $P = Q$ 。

# ELBO(Evidence Lower Bound)

- Maximize the likelihood  $p(\mathbf{x})$  of all observed  $\mathbf{x}$  ???

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right]$$

$$\log p(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] + D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x}))$$



# VAE(Variational Autoencoder)

$$\begin{aligned}
 \mathbb{E}_{q_{\phi}(z|x)} \left[ \log \frac{p(x,z)}{q_{\phi}(z|x)} \right] &= \mathbb{E}_{q_{\phi}(z|x)} \left[ \log \frac{p_{\theta}(x|z)p(z)}{q_{\phi}(z|x)} \right] && \text{(Chain Rule of Probability)} \\
 &= \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)] + \mathbb{E}_{q_{\phi}(z|x)} \left[ \log \frac{p(z)}{q_{\phi}(z|x)} \right] && \text{(Split the Expectation)} \\
 &= \underbrace{\mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_{\phi}(z|x) \parallel p(z))}_{\text{prior matching term}} && \text{(Definition of KL Divergence)}
 \end{aligned}$$

- Model

$$q_{\phi}(z|x) = \mathcal{N}(z; \mu_{\phi}(x), \sigma_{\phi}^2(x)\mathbf{I})$$

$$p(z) = \mathcal{N}(z; \mathbf{0}, \mathbf{I})$$

# VAE(Variational Autoencoder)

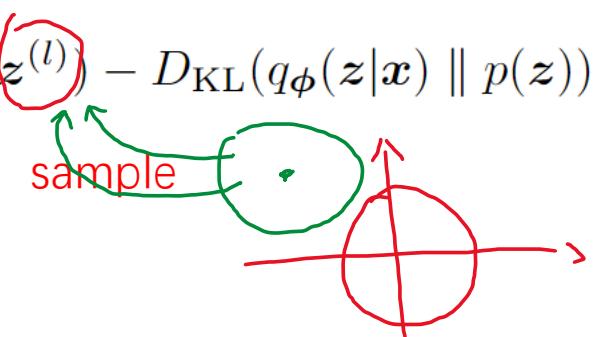
$$\begin{aligned}
 \mathbb{E}_{q_{\phi}(z|x)} \left[ \log \frac{p(x,z)}{q_{\phi}(z|x)} \right] &= \mathbb{E}_{q_{\phi}(z|x)} \left[ \log \frac{p_{\theta}(x|z)p(z)}{q_{\phi}(z|x)} \right] && (\text{Chain Rule of Probability}) \\
 &= \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)] + \mathbb{E}_{q_{\phi}(z|x)} \left[ \log \frac{p(z)}{q_{\phi}(z|x)} \right] && (\text{Split the Expectation}) \\
 &= \underbrace{\mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_{\phi}(z|x) \parallel p(z))}_{\text{prior matching term}} && (\text{Definition of KL Divergence})
 \end{aligned}$$

- optimize

$$\arg \max_{\phi, \theta} \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - D_{\text{KL}}(q_{\phi}(z|x) \parallel p(z)) \approx \arg \max_{\phi, \theta} \sum_{l=1}^L \log p_{\theta}(x|z^{(l)}) - D_{\text{KL}}(q_{\phi}(z|x) \parallel p(z))$$

reparameterization

$$z = \mu_{\phi}(x) + \sigma_{\phi}(x) \odot \epsilon \quad \text{with } \epsilon \sim \mathcal{N}(\epsilon; \mathbf{0}, \mathbf{I})$$



```

1 latent_dim = 2
2 input_dim = 28 * 28
3 inter_dim = 256
4
5 class VAE(nn.Module):
6     def __init__(self, input_dim=input_dim, inter_dim=inter_dim, latent_dim=latent_dim):
7         super(VAE, self).__init__()
8
9         self.encoder = nn.Sequential(
10             nn.Linear(input_dim, inter_dim),
11             nn.ReLU(),
12             nn.Linear(inter_dim, latent_dim * 2),
13         )
14
15         self.decoder = nn.Sequential(
16             nn.Linear(latent_dim, inter_dim),
17             nn.ReLU(),
18             nn.Linear(inter_dim, input_dim),
19             nn.Sigmoid(),
20         )
21
22     def reparameterise(self, mu, logvar):
23         epsilon = torch.randn_like(mu)
24         return mu + epsilon * torch.exp(logvar / 2)
25
26     def forward(self, x):
27         org_size = x.size()
28         batch = org_size[0]
29         x = x.view(batch, -1)
30
31         h = self.encoder(x)
32         mu, logvar = h.chunk(2, dim=1)
33         z = self.reparameterise(mu, logvar)
34         recon_x = self.decoder(z).view(size=org_size)
35
36         return recon_x, mu, logvar

```

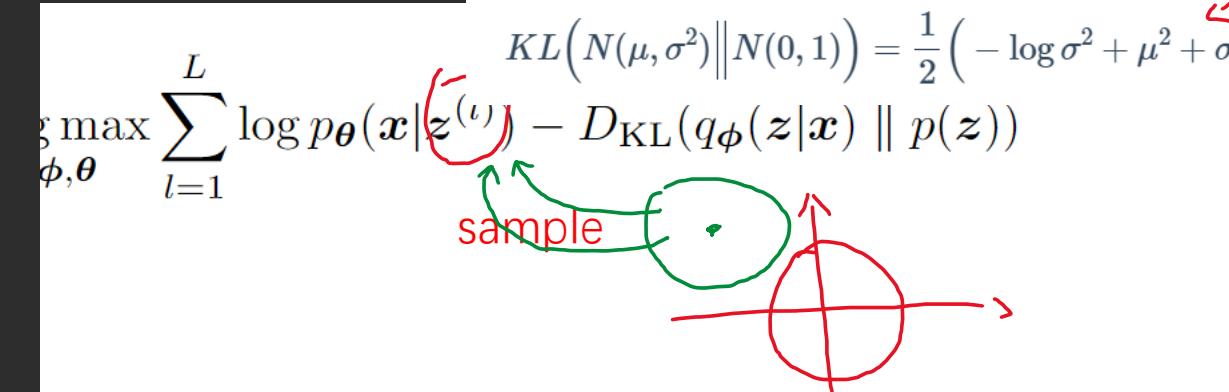
# coder)

$$\begin{aligned}
 &= \mathbb{E}_{q_{\phi}(z|x)} \left[ \log \frac{p_{\theta}(x|z)p(z)}{q_{\phi}(z|x)} \right] \\
 &= \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)] + \mathbb{E}_{q_{\phi}(z|x)} \left[ \log \frac{p(z)}{q_{\phi}(z|x)} \right] \\
 &= \underbrace{\mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_{\phi}(z|x) \parallel p(z))}_{\text{prior matching term}}
 \end{aligned}$$

```

1 kl_loss = lambda mu, logvar: -0.5 * torch.sum(1 + logvar - mu.pow(2) - logvar.exp())
2 recon_loss = lambda recon_x, x: F.binary_cross_entropy(recon_x, x, size_average=False)

```



# Rethink VAE && AE

- AE     $x \rightarrow z \rightarrow x$
- VAE  $x \rightarrow z \rightarrow x \quad z \sim N(0,1)$  放弃点到点的映射，看重分布关系
- 寻找后验分布  $p(z|x_i)$

# Hierarchical VAE

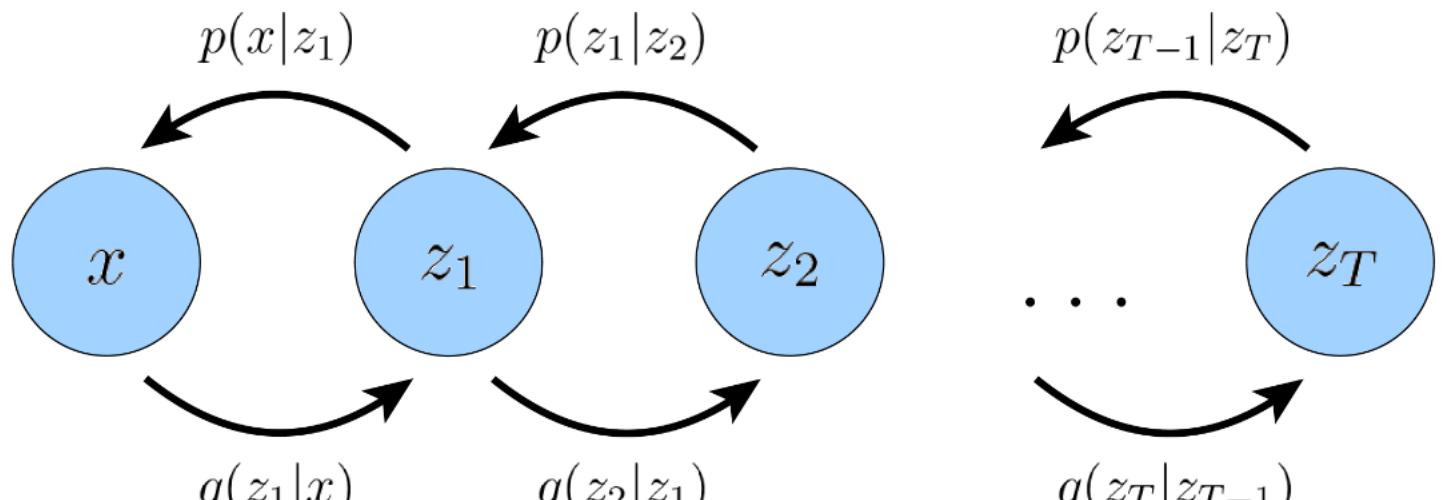
- generative process is a Markov chain

$$p(\mathbf{x}, \mathbf{z}_{1:T}) = p(\mathbf{z}_T)p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}_1) \prod_{t=2}^T p_{\boldsymbol{\theta}}(\mathbf{z}_{t-1}|\mathbf{z}_t)$$

$$q_{\boldsymbol{\phi}}(\mathbf{z}_{1:T}|\mathbf{x}) = q_{\boldsymbol{\phi}}(\mathbf{z}_1|\mathbf{x}) \prod_{t=2}^T q_{\boldsymbol{\phi}}(\mathbf{z}_t|\mathbf{z}_{t-1})$$

$$\begin{aligned} \log p(\mathbf{x}) &= \log \int p(\mathbf{x}, \mathbf{z}_{1:T}) d\mathbf{z}_{1:T} \\ &= \log \int \frac{p(\mathbf{x}, \mathbf{z}_{1:T}) q_{\boldsymbol{\phi}}(\mathbf{z}_{1:T}|\mathbf{x})}{q_{\boldsymbol{\phi}}(\mathbf{z}_{1:T}|\mathbf{x})} d\mathbf{z}_{1:T} \\ &= \log \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{1:T}|\mathbf{x})} \left[ \frac{p(\mathbf{x}, \mathbf{z}_{1:T})}{q_{\boldsymbol{\phi}}(\mathbf{z}_{1:T}|\mathbf{x})} \right] \\ &\geq \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{1:T}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z}_{1:T})}{q_{\boldsymbol{\phi}}(\mathbf{z}_{1:T}|\mathbf{x})} \right] \end{aligned}$$

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



# Diffusion Model

- The latent dimension is exactly equal to the data dimension
- Encoder is pre-defined as a linear Gaussian model

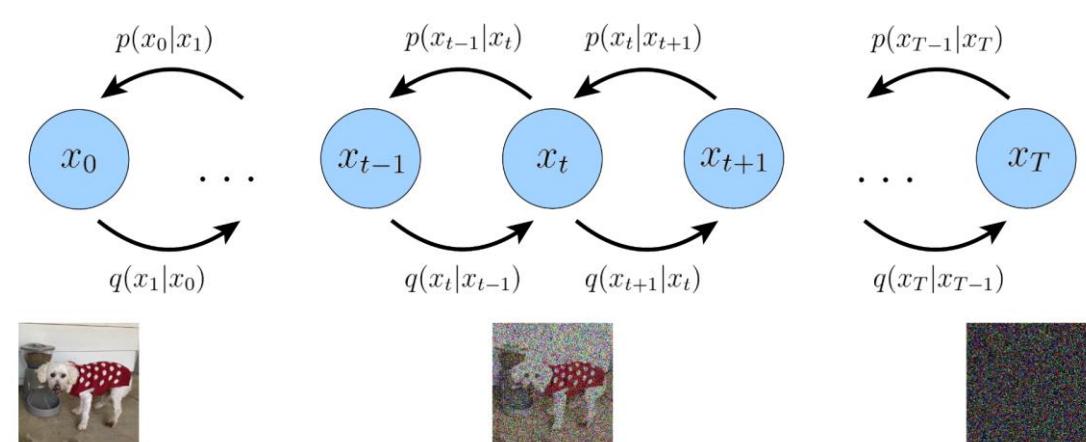
$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t} \mathbf{x}_{t-1}, (1 - \alpha_t) \mathbf{I})$$

- Distribution of the latent at final timestep T is a standard Gaussian

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

where,

$$p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$$



Luo, C. (2022). Understanding Diffusion Models: A Unified Perspective.

Ho, J., Jain, AjayN., & Abbeel, P. (2020). Denoising Diffusion Probabilistic Models. NIPS.

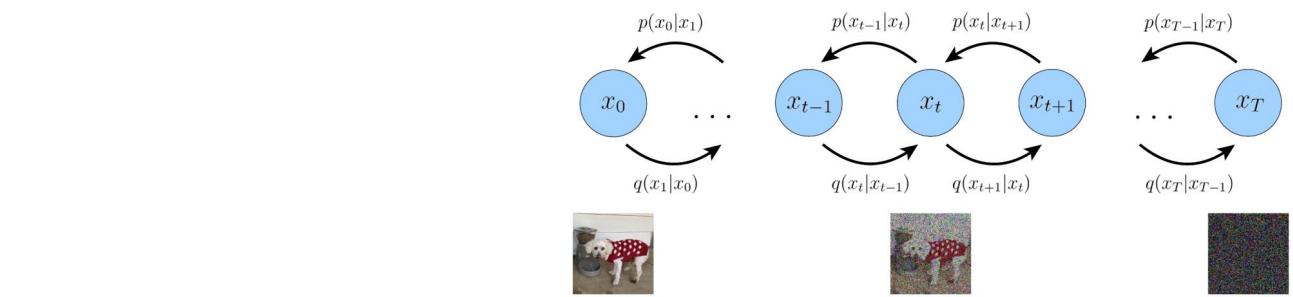
# Diffusion Model

$$\begin{aligned}
 \log p(\mathbf{x}) &\geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1) \prod_{t=2}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_1|\mathbf{x}_0) \prod_{t=2}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \quad \text{提出t=1} \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1) \prod_{t=2}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_1|\mathbf{x}_0) \prod_{t=2}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)} \right] \quad \text{前向马尔可夫性质} \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ \log \frac{p_{\theta}(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)} \right] \\
 &= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_T) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}} \right]
 \end{aligned}$$

$$p(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) \quad q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

where,

$$p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$$



$$\begin{aligned}
 &= \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] \\
 &= \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] \\
 &= \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] \\
 &= \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] \\
 &= \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] \\
 &= \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}}
 \end{aligned}$$

$$\begin{aligned}
 &\cdot \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}} \Bigg] \\
 &\log \prod_{t=2}^T \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \Bigg] \\
 &\log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{q(\mathbf{x}_2|\mathbf{x}_0)} \frac{q(\mathbf{x}_2|\mathbf{x}_0)}{q(\mathbf{x}_3|\mathbf{x}_0)} \dots \frac{q(\mathbf{x}_{T-1}|\mathbf{x}_0)}{q(\mathbf{x}_T|\mathbf{x}_0)} + \log
 \end{aligned}$$

贝叶斯公式

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

$$\begin{aligned} &= \mathbb{E}_{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0 | \mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T | \mathbf{x}_0)} \right] + \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \left[ \log \frac{p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)}{q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)} \right] \\ &= \underbrace{\mathbb{E}_{q(\mathbf{x}_1 | \mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0 | \mathbf{x}_1)]}_{\text{对任意的归一化概率分布 } p(x, y):} + \underbrace{\mathbb{E}_{q(\mathbf{x}_T | \mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T | \mathbf{x}_0)} \right]}_{\text{}} + \sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t, \mathbf{x}_{t-1} | \mathbf{x}_0)} \left[ \log \frac{p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)}{q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)} \right]}_{\text{}} \end{aligned}$$

### 1. 贝叶斯公式:

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

### 2. 联合概率的分解:

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

### 3. 代入联合概率的分解:

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

对任意的归一化概率分布  $p(x, y)$ :

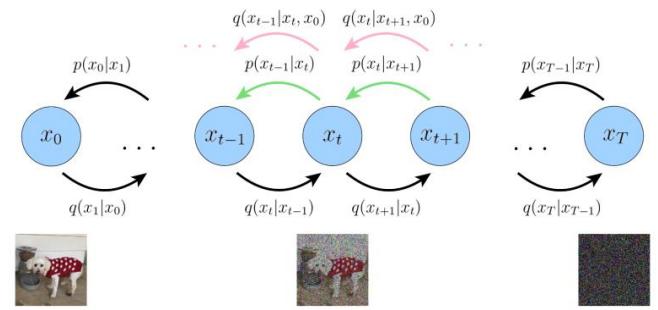
$$\begin{aligned} \mathbb{E}_{p(x,y)}[f(x)] &= \int p(x, y)f(x)dxdy \\ &= \int f(x) \left( \int p(x, y)dy \right) dx \\ &= \int p(x)f(x)dx \end{aligned}$$

其中,  $p(x)$  是边缘分布 (也满足归一化) :

$$p(x) = \int p(x, y)dy$$

就是一个简单的积分变化而已.

# Diffusion Model



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

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

$$\mathbf{x}_t = \sqrt{\alpha_t}\mathbf{x}_{t-1} + \sqrt{1-\alpha_t}\boldsymbol{\epsilon} \quad \text{with } \boldsymbol{\epsilon} \sim \mathcal{N}(\boldsymbol{\epsilon}; \mathbf{0}, \mathbf{I})$$

$$\mathcal{N}(x; \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^k |\Sigma|}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

$$\begin{aligned} \mathbf{x}_t &= \sqrt{\alpha_t}\mathbf{x}_{t-1} + \sqrt{1-\alpha_t}\boldsymbol{\epsilon}_{t-1}^* \\ &= \sqrt{\alpha_t} \left( \sqrt{\alpha_{t-1}}\mathbf{x}_{t-2} + \sqrt{1-\alpha_{t-1}}\boldsymbol{\epsilon}_{t-2}^* \right) + \sqrt{1-\alpha_t}\boldsymbol{\epsilon}_{t-1}^* \\ &= \sqrt{\alpha_t\alpha_{t-1}}\mathbf{x}_{t-2} + \sqrt{\alpha_t - \alpha_t\alpha_{t-1}}\boldsymbol{\epsilon}_{t-2}^* + \sqrt{1-\alpha_t}\boldsymbol{\epsilon}_{t-1}^* \\ &= \sqrt{\alpha_t\alpha_{t-1}}\mathbf{x}_{t-2} + \sqrt{\sqrt{\alpha_t - \alpha_t\alpha_{t-1}}^2 + \sqrt{1-\alpha_t}^2} \boldsymbol{\epsilon}_{t-2} \\ &= \sqrt{\alpha_t\alpha_{t-1}}\mathbf{x}_{t-2} + \sqrt{\alpha_t - \alpha_t\alpha_{t-1} + 1 - \alpha_t} \boldsymbol{\epsilon}_{t-2} \\ &= \sqrt{\alpha_t\alpha_{t-1}}\mathbf{x}_{t-2} + \sqrt{1 - \alpha_t\alpha_{t-1}} \boldsymbol{\epsilon}_{t-2} \\ &= \dots \\ &= \sqrt{\prod_{i=1}^t \alpha_i} \mathbf{x}_0 + \sqrt{1 - \prod_{i=1}^t \alpha_i} \boldsymbol{\epsilon}_0 \\ &= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_0 \\ &\sim \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I}) \end{aligned}$$

最后一个问题，为什么我们选择加噪方式为  $x_t = \sqrt{1 - \beta_t}x_{t-1} + \beta_t\epsilon$ ，换言之，若重写为  $x_t = k_t x_{t-1} + \beta_t\epsilon$ ，为什么要让  $k_t = \sqrt{1 - \beta_t}$  (i.e. 为什么  $k_t^2 + \beta_t^2 = 1$ )？由于我们有式 (3)  $x_t = \sqrt{\bar{\alpha}_t}x_0 + (1 - \bar{\alpha}_t)\epsilon$ ，我们希望  $t$  足够大时， $\sqrt{\bar{\alpha}_t}$  趋近于0，而  $1 - \bar{\alpha}_t$  趋于1，这样我们才能保证  $x_T$  趋于服从标准正态分布。

若  $x_t = k_t x_{t-1} + \beta_t\epsilon$ ，那么我们有

$$\begin{aligned} x_t &= k_t x_{t-1} + \beta_t\epsilon \\ &= k_t(k_{t-1}x_{t-2} + \beta_{t-1}\epsilon) + \beta_t\epsilon \\ &= \bar{k}_t x_0 + (k_t k_{t-1} \dots k_2 \beta_1 \epsilon) + (k_t k_{t-1} \dots k_3 \beta_2 \epsilon) + \dots + (\beta_t \epsilon) \end{aligned}$$

其中  $\bar{k}_t = k_1 k_2 \dots k_t$ 。由于上式括号中每个部分都是正态分布，且均值都是0，因此根据正态分布的可加性，这些括号相加依然是正态分布，且均值为0，方差为

$k_t^2 k_{t-1}^2 \dots k_2^2 \beta_1^2 + k_t^2 k_{t-1}^2 \dots k_3^2 \beta_2^2 + \dots + \beta_t^2$ 。而当  $k_t^2 + \beta_t^2 = 1$  时，我们有  
 $\beta_t^2 = 1 - k_t^2$ ，代入上式可得

$$\begin{aligned} &k_t^2 k_{t-1}^2 \dots k_2^2 (1 - k_1^2) + k_t^2 k_{t-1}^2 \dots k_3^2 (1 - k_2^2) + \dots + 1 - k_t^2 \\ &= k_t^2 k_{t-1}^2 \dots k_2^2 - k_t^2 k_{t-1}^2 \dots k_1^2 + k_t^2 k_{t-1}^2 \dots k_3^2 - k_t^2 k_{t-1}^2 \dots k_2^2 + \dots + 1 - k_t^2 \\ &= 1 - \bar{k}_t^2 \end{aligned}$$

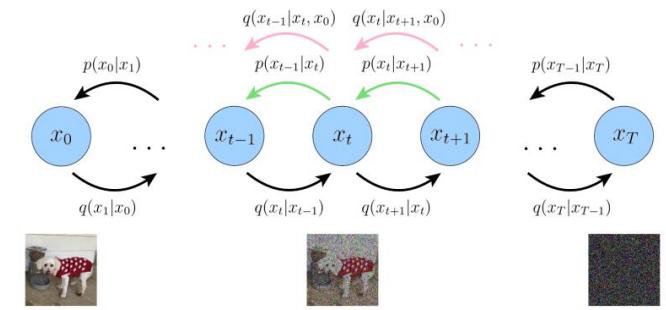
于是满足了我们的需求，即  $t$  足够大时有  $\bar{k}_t$  趋近于0， $1 - \bar{k}_t^2$  趋近于1，另外若令  $k_t = \sqrt{\bar{\alpha}_t}$ ，我们就得到了式 (3)。

yif

$$\begin{aligned}
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)} \\
&= \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})} \\
&\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\} \\
&= \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\} \\
&= \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\} \\
&\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\} \\
&= \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\} \\
&= \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\} \\
&= \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\} \\
&= \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\} \\
&= \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\} \\
&= \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\} \\
&= \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\} \\
&\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} \mathbf{I}}_{\Sigma_q(t)})
\end{aligned}$$

$$\begin{aligned}
\mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \\
&= \sqrt{\alpha_t} \left( \sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* \right) + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \\
&= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^* + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}^* \\
&= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\sqrt{\alpha_t - \alpha_t \alpha_{t-1}}^2 + \sqrt{1 - \alpha_t}^2} \boldsymbol{\epsilon}_{t-2} \\
&= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\alpha_t - \alpha_t \alpha_{t-1} + 1 - \alpha_t} \boldsymbol{\epsilon}_{t-2} \\
&= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2} \\
&= \dots \\
&= \sqrt{\prod_{i=1}^t \alpha_i} \mathbf{x}_0 + \sqrt{1 - \prod_{i=1}^t \alpha_i} \boldsymbol{\epsilon}_0 \\
&= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_0 \\
&\sim \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})
\end{aligned}$$

# Diffusion Model



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

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

- Model  $p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1}|\mathbf{x}_t)$   $\mu_{\boldsymbol{\theta}}(\mathbf{x}_t, t)$

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

$$\mu_{\boldsymbol{\theta}}(\mathbf{x}_t, t) = \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1-\alpha_t)\hat{\mathbf{x}}_{\boldsymbol{\theta}}(\mathbf{x}_t, t)}{1-\bar{\alpha}_t}$$

$$\begin{aligned} & \arg \min_{\boldsymbol{\theta}} D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1}|\mathbf{x}_t)) \\ &= \arg \min_{\boldsymbol{\theta}} D_{\text{KL}}(\mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q(t)) \parallel \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\boldsymbol{\theta}}, \boldsymbol{\Sigma}_{\boldsymbol{\theta}}(t))) \\ &= \arg \min_{\boldsymbol{\theta}} \frac{1}{2} \left[ \log \frac{|\boldsymbol{\Sigma}_q(t)|}{|\boldsymbol{\Sigma}_{\boldsymbol{\theta}}(t)|} - d + \text{tr}(\boldsymbol{\Sigma}_q(t)^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}}(t)) + (\boldsymbol{\mu}_{\boldsymbol{\theta}} - \boldsymbol{\mu}_q)^T \boldsymbol{\Sigma}_q(t)^{-1} (\boldsymbol{\mu}_{\boldsymbol{\theta}} - \boldsymbol{\mu}_q) \right] \\ &= \arg \min_{\boldsymbol{\theta}} \frac{1}{2} \left[ \log 1 - d + d + (\boldsymbol{\mu}_{\boldsymbol{\theta}} - \boldsymbol{\mu}_q)^T \boldsymbol{\Sigma}_q(t)^{-1} (\boldsymbol{\mu}_{\boldsymbol{\theta}} - \boldsymbol{\mu}_q) \right] \\ &= \arg \min_{\boldsymbol{\theta}} \frac{1}{2} \left[ (\boldsymbol{\mu}_{\boldsymbol{\theta}} - \boldsymbol{\mu}_q)^T \boldsymbol{\Sigma}_q(t)^{-1} (\boldsymbol{\mu}_{\boldsymbol{\theta}} - \boldsymbol{\mu}_q) \right] \\ &= \arg \min_{\boldsymbol{\theta}} \frac{1}{2} \left[ (\boldsymbol{\mu}_{\boldsymbol{\theta}} - \boldsymbol{\mu}_q)^T (\sigma_q^2(t)\mathbf{I})^{-1} (\boldsymbol{\mu}_{\boldsymbol{\theta}} - \boldsymbol{\mu}_q) \right] \\ &= \arg \min_{\boldsymbol{\theta}} \frac{1}{2\sigma_q^2(t)} \left[ \|\boldsymbol{\mu}_{\boldsymbol{\theta}} - \boldsymbol{\mu}_q\|_2^2 \right] \\ &= \arg \min_{\boldsymbol{\theta}} \frac{1}{2\sigma_q^2(t)} \left[ \left\| \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1-\alpha_t)\hat{\mathbf{x}}_{\boldsymbol{\theta}}(\mathbf{x}_t, t)}{1-\bar{\alpha}_t} - \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} \right\|_2^2 \right] \\ &= \arg \min_{\boldsymbol{\theta}} \frac{1}{2\sigma_q^2(t)} \left[ \left\| \frac{\sqrt{\bar{\alpha}_{t-1}}(1-\alpha_t)\hat{\mathbf{x}}_{\boldsymbol{\theta}}(\mathbf{x}_t, t)}{1-\bar{\alpha}_t} - \frac{\sqrt{\bar{\alpha}_{t-1}}(1-\alpha_t)\mathbf{x}_0}{1-\bar{\alpha}_t} \right\|_2^2 \right] \\ &= \arg \min_{\boldsymbol{\theta}} \frac{1}{2\sigma_q^2(t)} \left[ \left\| \frac{\sqrt{\bar{\alpha}_{t-1}}(1-\alpha_t)}{1-\bar{\alpha}_t} (\hat{\mathbf{x}}_{\boldsymbol{\theta}}(\mathbf{x}_t, t) - \mathbf{x}_0) \right\|_2^2 \right] \\ &= \arg \min_{\boldsymbol{\theta}} \frac{1}{2\sigma_q^2(t)} \frac{\bar{\alpha}_{t-1}(1-\alpha_t)^2}{(1-\bar{\alpha}_t)^2} \left[ \|\hat{\mathbf{x}}_{\boldsymbol{\theta}}(\mathbf{x}_t, t) - \mathbf{x}_0\|_2^2 \right] \end{aligned}$$

# Diffusion Model

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

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

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

$$\mathbf{x}_0 = \frac{\mathbf{x}_t - \sqrt{1-\bar{\alpha}_t}\epsilon_0}{\sqrt{\bar{\alpha}_t}}$$

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

$$= \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1-\alpha_t)\frac{\mathbf{x}_t - \sqrt{1-\bar{\alpha}_t}\epsilon_0}{\sqrt{\alpha_t}}}{1-\bar{\alpha}_t}$$

$$= \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})\mathbf{x}_t + (1-\alpha_t)\frac{\mathbf{x}_t - \sqrt{1-\bar{\alpha}_t}\epsilon_0}{\sqrt{\alpha_t}}}{1-\bar{\alpha}_t}$$

$$= \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})\mathbf{x}_t + \frac{(1-\alpha_t)\mathbf{x}_t}{(1-\bar{\alpha}_t)\sqrt{\alpha_t}} - \frac{(1-\alpha_t)\sqrt{1-\bar{\alpha}_t}\epsilon_0}{(1-\bar{\alpha}_t)\sqrt{\alpha_t}}}{1-\bar{\alpha}_t}$$

$$= \left( \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} + \frac{1-\alpha_t}{(1-\bar{\alpha}_t)\sqrt{\alpha_t}} \right) \mathbf{x}_t - \frac{(1-\alpha_t)\sqrt{1-\bar{\alpha}_t}\epsilon_0}{(1-\bar{\alpha}_t)\sqrt{\alpha_t}}$$

$$= \left( \frac{\alpha_t(1-\bar{\alpha}_{t-1})}{(1-\bar{\alpha}_t)\sqrt{\alpha_t}} + \frac{1-\alpha_t}{(1-\bar{\alpha}_t)\sqrt{\alpha_t}} \right) \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}\sqrt{\alpha_t}}\epsilon_0$$

$$= \frac{\alpha_t - \bar{\alpha}_t + 1 - \alpha_t}{(1-\bar{\alpha}_t)\sqrt{\alpha_t}} \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}\sqrt{\alpha_t}}\epsilon_0$$

$$= \frac{1-\bar{\alpha}_t}{(1-\bar{\alpha}_t)\sqrt{\alpha_t}} \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}\sqrt{\alpha_t}}\epsilon_0$$

$$= \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}\sqrt{\alpha_t}}\epsilon_0$$

• Model  $\mu_{\boldsymbol{\theta}}(\mathbf{x}_t, t)$

$$\mu_{\boldsymbol{\theta}}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}\sqrt{\alpha_t}} \hat{\epsilon}_{\boldsymbol{\theta}}(\mathbf{x}_t, t)$$

$$\begin{aligned} & \arg \min_{\boldsymbol{\theta}} D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1}|\mathbf{x}_t)) \\ &= \arg \min_{\boldsymbol{\theta}} D_{\text{KL}}(\mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q(t)) \parallel \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\boldsymbol{\theta}}, \boldsymbol{\Sigma}_q(t))) \\ &= \arg \min_{\boldsymbol{\theta}} \frac{1}{2\sigma_q^2(t)} \left[ \left\| \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}\sqrt{\alpha_t}} \hat{\epsilon}_{\boldsymbol{\theta}}(\mathbf{x}_t, t) - \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t + \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}\sqrt{\alpha_t}} \epsilon_0 \right\|_2^2 \right] \\ &= \arg \min_{\boldsymbol{\theta}} \frac{1}{2\sigma_q^2(t)} \left[ \left\| \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}\sqrt{\alpha_t}} \epsilon_0 - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}\sqrt{\alpha_t}} \hat{\epsilon}_{\boldsymbol{\theta}}(\mathbf{x}_t, t) \right\|_2^2 \right] \\ &= \arg \min_{\boldsymbol{\theta}} \frac{1}{2\sigma_q^2(t)} \left[ \left\| \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}\sqrt{\alpha_t}} (\epsilon_0 - \hat{\epsilon}_{\boldsymbol{\theta}}(\mathbf{x}_t, t)) \right\|_2^2 \right] \\ &= \arg \min_{\boldsymbol{\theta}} \frac{1}{2\sigma_q^2(t)} \frac{(1-\alpha_t)^2}{(1-\bar{\alpha}_t)\alpha_t} \left[ \|\epsilon_0 - \hat{\epsilon}_{\boldsymbol{\theta}}(\mathbf{x}_t, t)\|_2^2 \right] \end{aligned}$$

# Diffusion Model

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

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

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

- Model  $\mu_{\theta}(\mathbf{x}_t, t)$

$$\mu_{\theta}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}}\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}\sqrt{\alpha_t}}\hat{\epsilon}_{\theta}(\mathbf{x}_t, t)$$

---

**Algorithm 1** Training

---

```

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

```

---

$$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$$

$$\mu_{\theta}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}}\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}\sqrt{\alpha_t}}\hat{\epsilon}_{\theta}(\mathbf{x}_t, t)$$

$$\Sigma_q(t) = \frac{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}$$

---

**Algorithm 2** Sampling

---

```

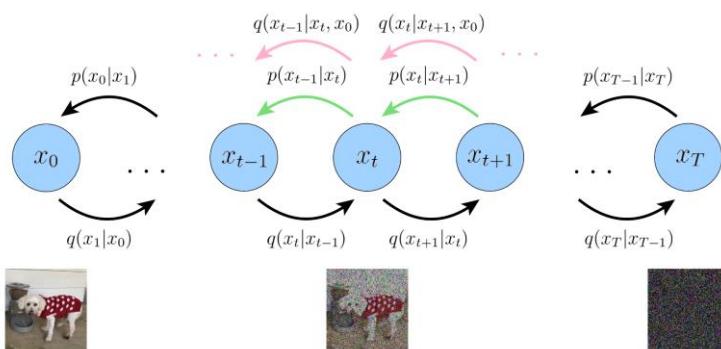
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}}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

```

---

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

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




---

**Algorithm 2** Sampling

---

```

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

```

---

$$q(x_k|x_s, x_0) = \frac{q(x_s|x_k, x_0)q(x_k|x_0)}{q(x_s|x_0)} \quad k \leq s-1$$

$$q(x_k|x_s, x_0) \sim N(kx_0 + mx_s, \sigma^2 I)$$

$$\begin{aligned} x_k &= (kx_0 + mx_s) + \sigma\epsilon \\ &= kx_0 + m(\sqrt{\bar{\alpha}_s}x_0 + \sqrt{1-\bar{\alpha}_s}\epsilon' + \sigma\epsilon) \\ &= (k + m\sqrt{\bar{\alpha}_s})x_0 + \sqrt{m^2(1-\bar{\alpha}_s) + \sigma^2}\epsilon \\ &= \sqrt{\bar{\alpha}_k}x_0 + \sqrt{1-\bar{\alpha}_k}\epsilon \end{aligned}$$

$$m = \frac{\sqrt{1-\bar{\alpha}_k}-\sigma^2}{\sqrt{1-\bar{\alpha}_s}} \quad k = \sqrt{\bar{\alpha}_k} - \frac{\sqrt{1-\bar{\alpha}_k}-\sigma^2}{\sqrt{1-\bar{\alpha}_s}}\sqrt{\bar{\alpha}_s}$$

$$x_k = (\sqrt{\bar{\alpha}_k} - \frac{\sqrt{1-\bar{\alpha}_k}-\sigma^2}{\sqrt{1-\bar{\alpha}_s}}\sqrt{\bar{\alpha}_s})x_0 + \frac{\sqrt{1-\bar{\alpha}_k}-\sigma^2}{\sqrt{1-\bar{\alpha}_s}}x_s + \sigma\epsilon$$

$$q(x_k|x_s, x_0) \sim N((\sqrt{\bar{\alpha}_k} - \frac{\sqrt{1-\bar{\alpha}_k}-\sigma^2}{\sqrt{1-\bar{\alpha}_s}}\sqrt{\bar{\alpha}_s})x_0 + \frac{\sqrt{1-\bar{\alpha}_k}-\sigma^2}{\sqrt{1-\bar{\alpha}_s}}x_s, \sigma^2 I)$$

$$\mu_{\theta} = (\sqrt{\bar{\alpha}_k} - \frac{\sqrt{1-\bar{\alpha}_k}-\sigma^2}{\sqrt{1-\bar{\alpha}_s}}\sqrt{\bar{\alpha}_s})x_0 + \frac{\sqrt{1-\bar{\alpha}_k}-\sigma^2}{\sqrt{1-\bar{\alpha}_s}}x_s$$

$$= \frac{\sqrt{\bar{\alpha}_k}}{\sqrt{\bar{\alpha}_s}}(x_s - (\sqrt{1-\bar{\alpha}_s} - \frac{\sqrt{\bar{\alpha}_s}}{\sqrt{\bar{\alpha}_k}}\sqrt{1-\bar{\alpha}_k-\sigma^2})\epsilon)$$

$$x_0 = \frac{\mathbf{x}_t - \sqrt{1-\bar{\alpha}_t}\epsilon_0}{\sqrt{\bar{\alpha}_t}}$$

# DDIM

Table 1: CIFAR10 and CelebA image generation measured in FID.  $\eta = 1.0$  and  $\hat{\sigma}$  are cases of DDPM (although Ho et al. (2020) only considered  $T = 1000$  steps, and  $S < T$  can be seen as simulating DDPMs trained with  $S$  steps), and  $\eta = 0.0$  indicates DDIM.

$S$	CIFAR10 ( $32 \times 32$ )					CelebA ( $64 \times 64$ )				
	10	20	50	100	1000	10	20	50	100	1000
$\eta$	0.0	<b>13.36</b>	<b>6.84</b>	<b>4.67</b>	<b>4.16</b>	4.04	<b>17.33</b>	<b>13.73</b>	<b>9.17</b>	<b>6.53</b>
	0.2	14.04	7.11	4.77	4.25	4.09	17.66	14.11	9.51	6.79
	0.5	16.66	8.35	5.25	4.46	4.29	19.86	16.06	11.01	8.09
	1.0	41.07	18.36	8.01	5.78	4.73	33.12	26.03	18.48	13.93
$\hat{\sigma}$	367.43	133.37	32.72	9.99	<b>3.17</b>	299.71	183.83	71.71	45.20	<b>3.26</b>

$$q(x_k|x_s, x_0) \sim N\left(\left(\sqrt{1 - \bar{\alpha}_k} - \frac{\sqrt{1 - \bar{\alpha}_k} - \sigma^2}{\sqrt{1 - \bar{\alpha}_s}}\sqrt{\bar{\alpha}_s}\right)x_0 + \frac{\sqrt{1 - \bar{\alpha}_k} - \sigma^2}{\sqrt{1 - \bar{\alpha}_s}}x_s, \sigma^2 I\right)$$

# others

- SDE (随机微分方程)
- ODE (常微分方程)

Song, Y., Sohl-Dickstein, J., Kingma, DiederikP., Kumar, A., Ermon, S., & Poole, B. (2020). Score-Based Generative Modeling through Stochastic Differential Equations. arXiv.

# Guidance—classifier

- 和前向加噪无关 (论文附录有证明)
- 只需推理时对均值做偏移

---

**Algorithm 1** Classifier guided sampling, given a diffusion model  $(\mu_\theta(x_t), \Sigma_\theta(x_t))$ , classifier  $p_\phi(y|x_t)$ , and gradient scale  $s$ .

---

Input: class label  $y$ , gradient scale  $s$

$x_0 \leftarrow$  sample from  $\mathcal{N}(0, \mathbf{I})$

**for all**  $t$  from  $T$  to 1 **do**

$\mu \leftarrow \mu_\theta(x_t)$

$\Sigma \leftarrow \Sigma_\theta(x_t)$

$g \leftarrow s \nabla_{x_t} \log p_\phi(y|x_t)$

$x_{t-1} \leftarrow$  sample from  $\mathcal{N}(\mu + \Sigma g, \Sigma)$

**end for**

**return**  $x_0$

---

# Guidance—classifier

- 只需推理时对均值做偏移

$$\begin{aligned} q(x_{t-1}|x_t, y) &= \frac{q(x_{t-1}|x_t)q(y|x_{t-1}, x_t)}{q(y|x_t)} \\ &= Z q(x_{t-1}|x_t)q(y|x_{t-1}) \quad (\text{由与前向无关推导得出}) \end{aligned}$$

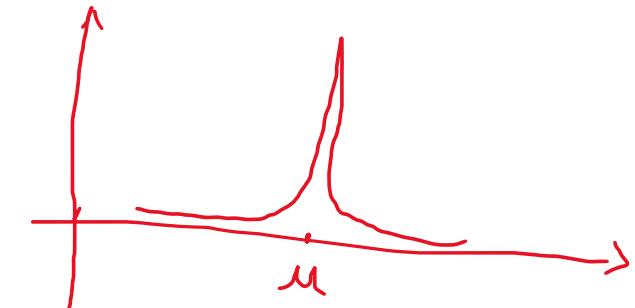
$$p_{\theta, \phi}(x_t|x_{t+1}, y) = Z p_{\theta}(x_t|x_{t+1}) p_{\phi}(y|x_t)$$

$$\log p_{\theta}(x_t|x_{t+1}) = -\frac{1}{2}(x_t - \mu)^T \Sigma^{-1}(x_t - \mu) + C$$

$$\begin{aligned} \log p_{\phi}(y|x_t) &\approx \log p_{\phi}(y|x_t)|_{x_t=\mu} + (x_t - \mu) \nabla_{x_t} \log p_{\phi}(y|x_t)|_{x_t=\mu} \\ &= (x_t - \mu)g + C_1 \end{aligned}$$



$$\begin{aligned} \log(p_{\theta}(x_t|x_{t+1})p_{\phi}(y|x_t)) &\approx -\frac{1}{2}(x_t - \mu)^T \Sigma^{-1}(x_t - \mu) + (x_t - \mu)g + C_2 \\ &= -\frac{1}{2}(x_t - \mu - \Sigma g)^T \Sigma^{-1}(x_t - \mu - \Sigma g) + \frac{1}{2}g^T \Sigma g + C_2 \\ &= -\frac{1}{2}(x_t - \mu - \Sigma g)^T \Sigma^{-1}(x_t - \mu - \Sigma g) + C_3 \\ &= \log p(z) + C_4, z \sim \mathcal{N}(\mu + \Sigma g, \Sigma) \end{aligned}$$



**Algorithm 1** Classifier guided sampling, given a diffusion model  $(\mu_{\theta}(x_t), \Sigma_{\theta}(x_t))$ , classifier  $p_{\phi}(y|x_t)$ , and gradient scale  $s$ .

---

```

Input: class label  $y$ , gradient scale  $s$ 
 $x_0 \leftarrow$  sample from  $\mathcal{N}(0, \mathbf{I})$ 
for all  $t$  from  $T$  to 1 do
     $\mu \leftarrow \mu_{\theta}(x_t)$ 
     $\Sigma \leftarrow \Sigma_{\theta}(x_t)$ 
     $g \leftarrow s \nabla_{x_t} \log p_{\phi}(y|x_t)$ 
     $x_{t-1} \leftarrow$  sample from  $\mathcal{N}(\mu + \Sigma g, \Sigma)$ 
end for
return  $x_0$ 

```

---

<https://kexue.fm/archives/9257>

Dhariwal, P., & Nichol, A. (2021). Diffusion Models Beat GANs on Image Synthesis. NIPS.

[https://www.bilibili.com/video/BV1s8411i7cU/?spm\\_id\\_from=333.788&vd\\_source=26cdccafc07aa98082eae4058dbfcf75](https://www.bilibili.com/video/BV1s8411i7cU/?spm_id_from=333.788&vd_source=26cdccafc07aa98082eae4058dbfcf75)

# Guidance—classifier in DDIM

$$\nabla_{x_t} \log p_\theta(x_t) = -\frac{1}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(x_t)$$

$$\begin{aligned} \nabla_{x_t} \log(p_\theta(x_t)p_\phi(y|x_t)) &= \nabla_{x_t} \log p_\theta(x_t) + \nabla_{x_t} \log p_\phi(y|x_t) \\ &\quad - \frac{\hat{\epsilon}}{\sqrt{1-\bar{\alpha}_t}} = -\frac{1}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(x_t) + \nabla_{x_t} \log p_\phi(y|x_t) \end{aligned}$$

$$\hat{\epsilon}(x_t) := \epsilon_\theta(x_t) - \sqrt{1-\bar{\alpha}_t} \nabla_{x_t} \log p_\phi(y|x_t)$$

**Algorithm 1** Classifier guided sampling, given a diffusion model  $(\mu_\theta(x_t), \Sigma_\theta(x_t))$ , classifier  $p_\phi(y|x_t)$ , and gradient scale  $s$ .

---

```

Input: class label  $y$ , gradient scale  $s$ 
 $x_0 \leftarrow$  sample from  $\mathcal{N}(0, \mathbf{I})$ 
for all  $t$  from  $T$  to 1 do
     $\mu \leftarrow \mu_\theta(x_t)$ 
     $\Sigma \leftarrow \Sigma_\theta(x_t)$ 
     $g \leftarrow s \nabla_{x_t} \log p_\phi(y|x_t)$ 
     $x_{t-1} \leftarrow$  sample from  $\mathcal{N}(\mu + \Sigma g, \Sigma)$ 
end for
return  $x_0$ 
```

---

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

$$\nabla \log p(x_t|x_0) = -\frac{(x_t - \sqrt{\bar{\alpha}_t}x_0)}{1 - \bar{\alpha}_t} = -\frac{\epsilon}{\sqrt{1 - \bar{\alpha}_t}}$$

# Guidance—classifier-free

---

**Algorithm 1** Joint training a diffusion model with classifier-free guidance

**Require:**  $p_{\text{uncond}}$ : probability of unconditional training

- ```

1: repeat
2:    $(\mathbf{x}, \mathbf{c}) \sim p(\mathbf{x}, \mathbf{c})$                                  $\triangleright$  Sample data with conditioning from the dataset
3:    $\mathbf{c} \leftarrow \emptyset$  with probability  $p_{\text{uncond}}$   $\triangleright$  Randomly discard conditioning to train unconditionally
4:    $\lambda \sim p(\lambda)$  $\triangleright$  Sample log SNR value
5:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
6:    $\mathbf{z}_\lambda = \alpha_\lambda \mathbf{x} + \sigma_\lambda \epsilon$                  $\triangleright$  Corrupt data to the sampled log SNR value
7:   Take gradient step on  $\nabla_\theta \|\epsilon_\theta(\mathbf{z}_\lambda, \mathbf{c}) - \epsilon\|^2$      $\triangleright$  Optimization of denoising model
8: until converged

```

---

**Algorithm 2** Conditional sampling with classifier-free guidance

**Require:**  $w$ : guidance strength

**Require:**  $c$ : conditioning information for conditional sampling

**Require:**  $\lambda_1, \dots, \lambda_T$ : increasing log SNR sequence with  $\lambda_1 = \lambda_{\min}$ ,  $\lambda_T = \lambda_{\max}$

- ```

1:  $\mathbf{z}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = 1, \dots, T$  do
   ▷ Form the classifier-free guided score at log SNR  $\lambda_t$ 
3:    $\tilde{\epsilon}_t = (1 + w)\epsilon_\theta(\mathbf{z}_t, \mathbf{c}) - w\epsilon_\theta(\mathbf{z}_t)$ 
   ▷ Sampling step (could be replaced by another sampler, e.g. DDIM)
4:    $\tilde{\mathbf{x}}_t = (\mathbf{z}_t - \sigma_{\lambda_t} \tilde{\epsilon}_t) / \alpha_{\lambda_t}$ 
5:    $\mathbf{z}_{t+1} \sim \mathcal{N}(\tilde{\mu}_{\lambda_{t+1} | \lambda_t}(\mathbf{z}_t, \tilde{\mathbf{x}}_t), (\tilde{\sigma}_{\lambda_{t+1} | \lambda_t}^2)^{1-v} (\sigma_{\lambda_t | \lambda_{t+1}}^2)^v)$  if  $t < T$  else  $\mathbf{z}_{t+1} = \tilde{\mathbf{x}}_t$ 
6: end for
7: return  $\mathbf{z}_{T+1}$ 

```

# LDM—Stable Diffusion

- Latent Diffusion Model

- 为什么单独VAE效果不好

1. unkown->N(0,1)

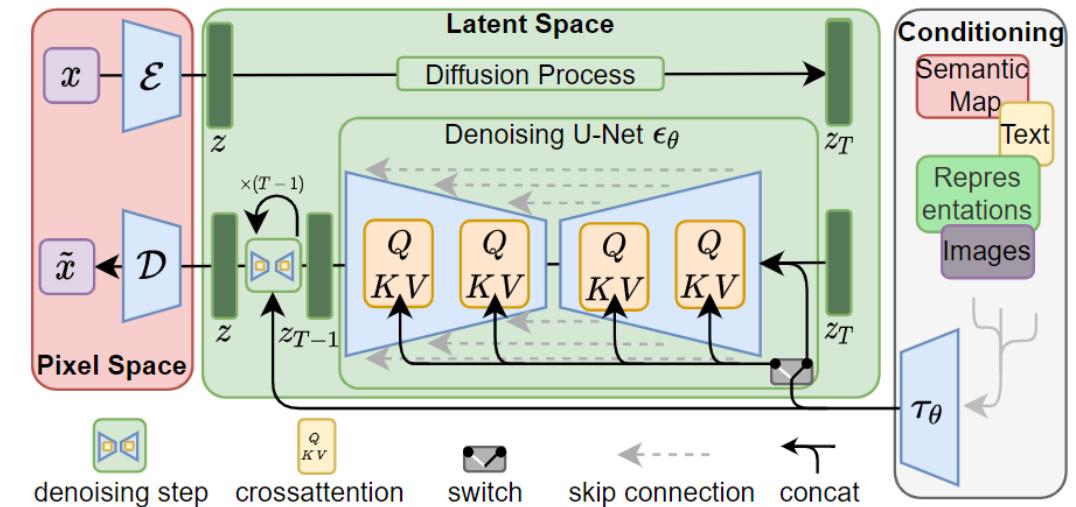


Figure 3. We condition LDMs either via concatenation or by a more general cross-attention mechanism. See Sec. 3.3

$$\begin{aligned}
 \mathbb{E}_{q_\phi(z|x)} \left[ \log \frac{p(x, z)}{q_\phi(z|x)} \right] &= \mathbb{E}_{q_\phi(z|x)} \left[ \log \frac{p_\theta(x|z)p(z)}{q_\phi(z|x)} \right] && \text{(Chain Rule of Probability)} \\
 &= \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z)] + \mathbb{E}_{q_\phi(z|x)} \left[ \log \frac{p(z)}{q_\phi(z|x)} \right] && \text{(Split the Expectation)} \\
 &= \underbrace{\mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_\phi(z|x) \parallel p(z))}_{\text{prior matching term}} && \text{(Definition of KL Divergence)}
 \end{aligned}$$

# Application

## 微调

- Dreambooth
- LORA
- Textual inversion

## 细粒度控制

- controlnet

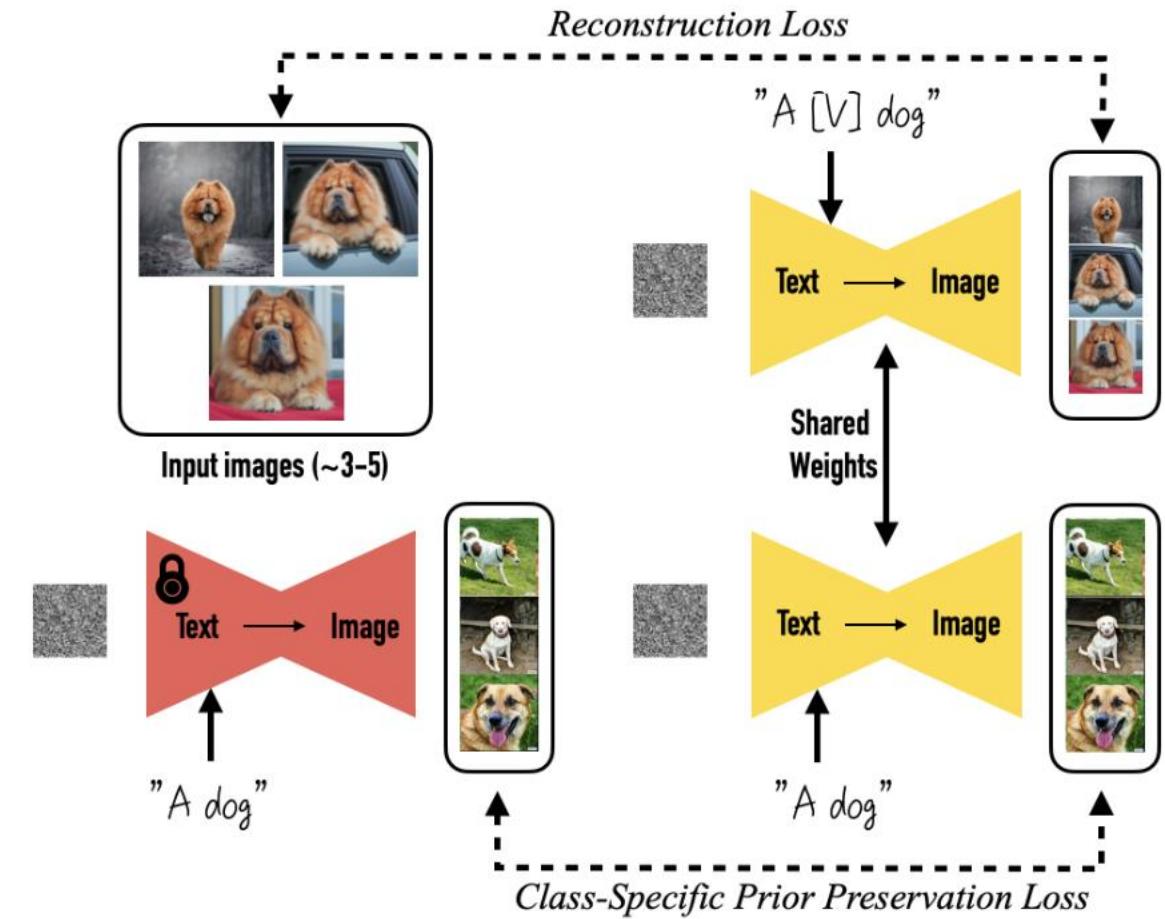
## 应用

- Repaint
- Blended-diffusion
- Diffedit
- P2p
- Pnp
- InstructPix2Pix
- Zero-classifier

# Dreambooth

- Personalization
- Rare-token $\leftarrow\rightarrow$ subject

$$\mathbb{E}_{\mathbf{x}, \mathbf{c}, \boldsymbol{\epsilon}, \boldsymbol{\epsilon}', t} [w_t \|\hat{\mathbf{x}}_\theta(\alpha_t \mathbf{x} + \sigma_t \boldsymbol{\epsilon}, \mathbf{c}) - \mathbf{x}\|_2^2 + \lambda w_{t'} \|\hat{\mathbf{x}}_\theta(\alpha_{t'} \mathbf{x}_{\text{pr}} + \sigma_{t'} \boldsymbol{\epsilon}', \mathbf{c}_{\text{pr}}) - \mathbf{x}_{\text{pr}}\|_2^2],$$



# Dreambooth



Input images



A [V] backpack in the Grand Canyon



A wet [V] backpack in water



A [V] backpack in Boston



A [V] backpack at night



DreamBooth (Imagen)



Input images



A [V] teapot floating in milk



A transparent [V] teapot with milk inside



A [V] teapot pouring tea



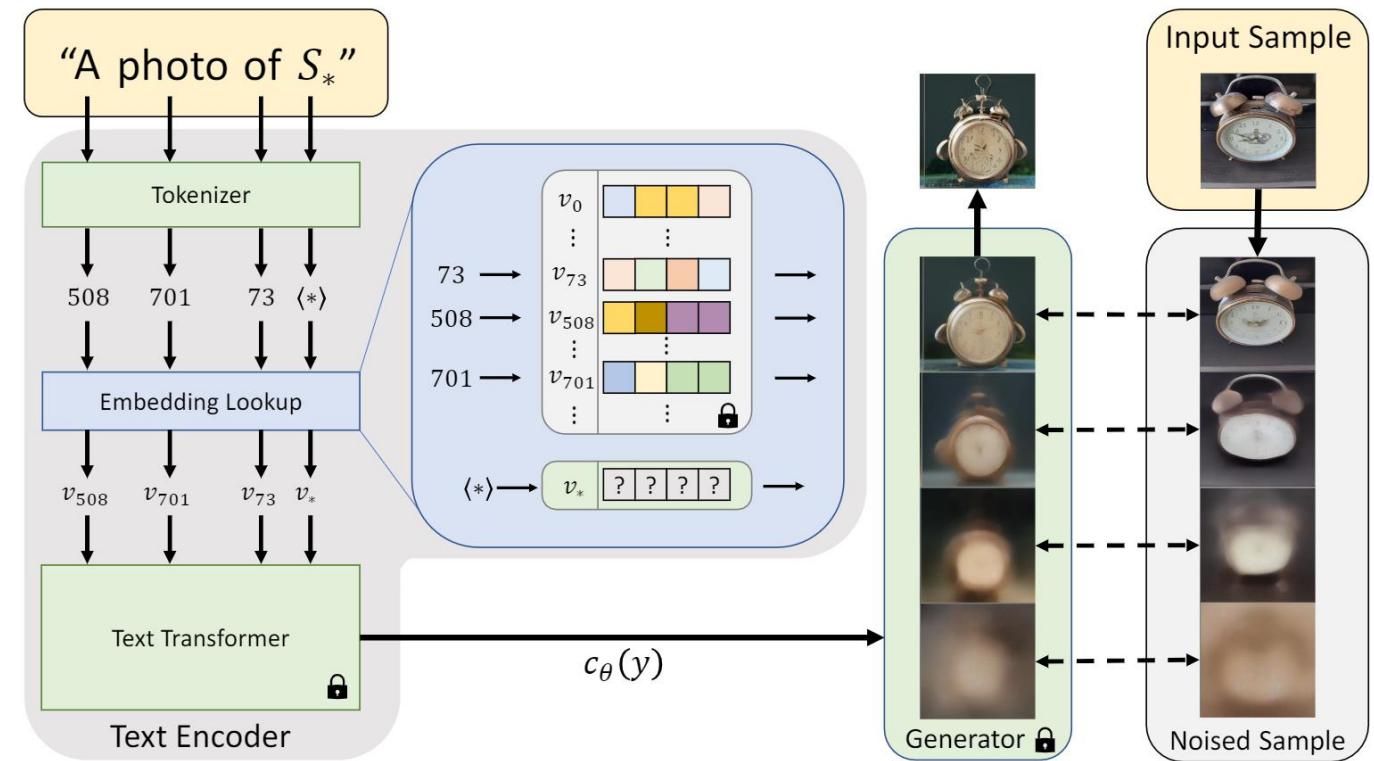
A [V] teapot in the sun



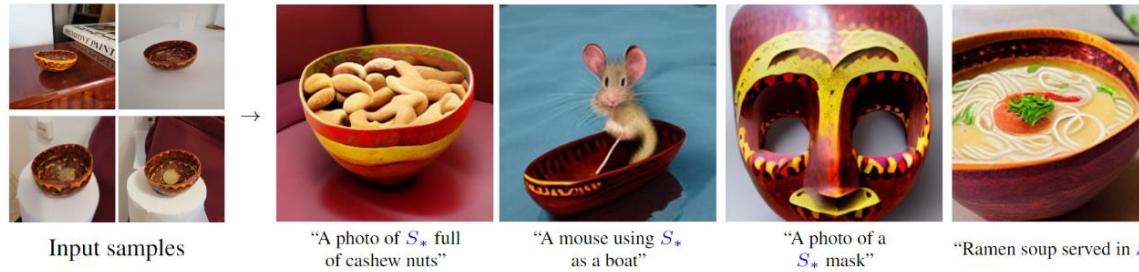
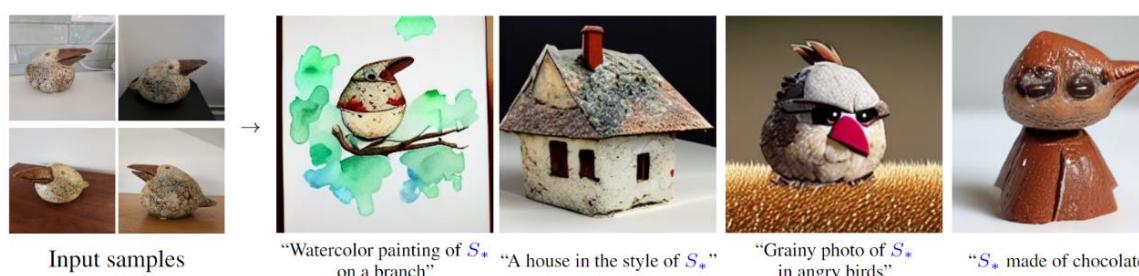
Textual Inversion (Stable Diffusion)



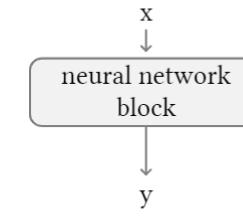
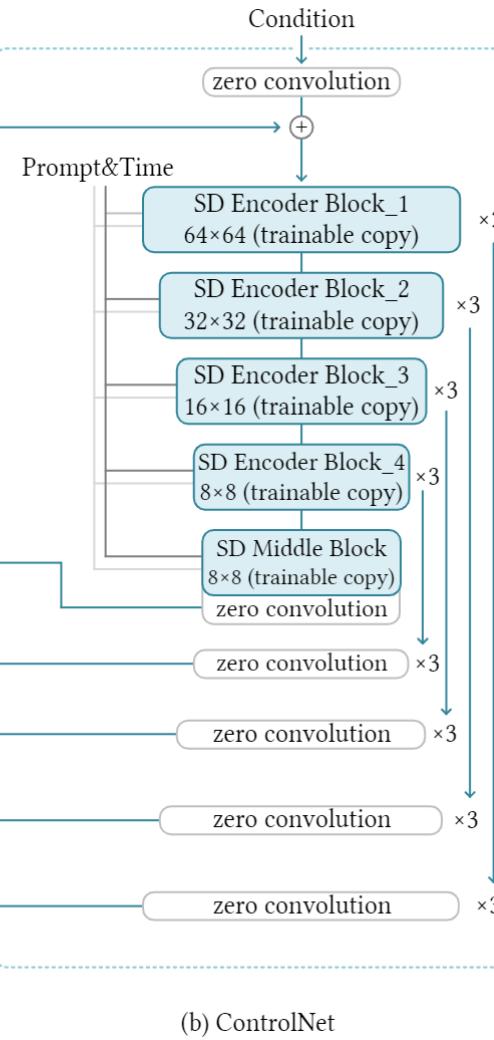
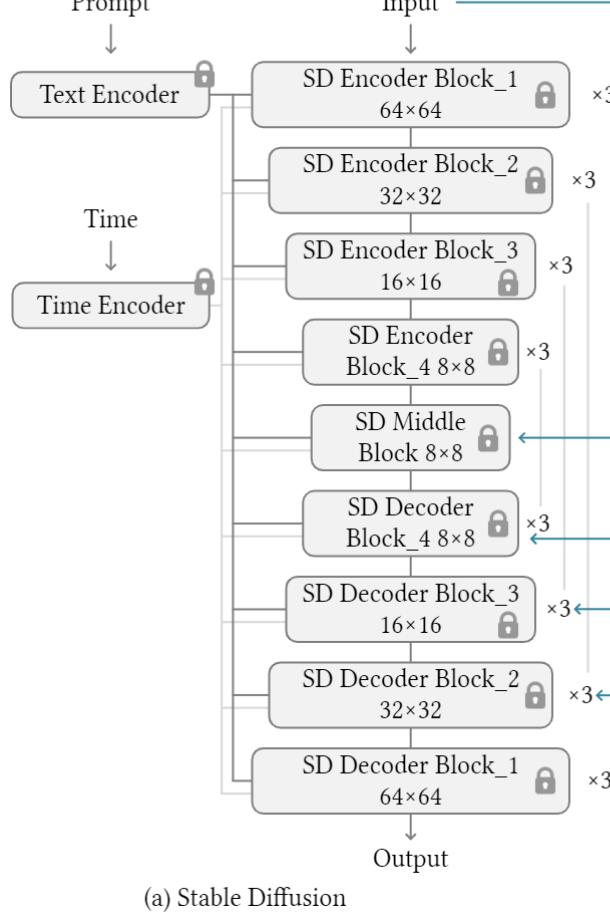
# Textual inversion



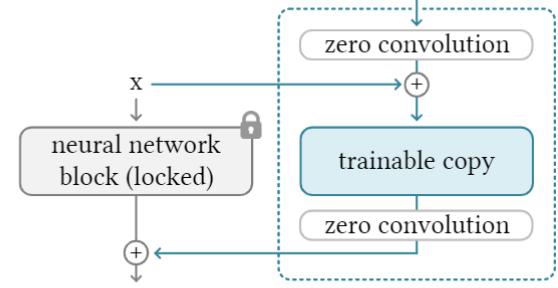
Gal, R., Alaluf, Y., Atzmon, Y., Patashnik, O., Bermano, A., Chechik, G., & Cohen-Or, D. (2022). An Image is Worth One Word: Personalizing Text-to-Image Generation using Textual Inversion.



# Controlnet



(a) Before



<https://github.com/Ilyasviel/ControlNet/discussions/188>

# Controlnet



Input Canny edge



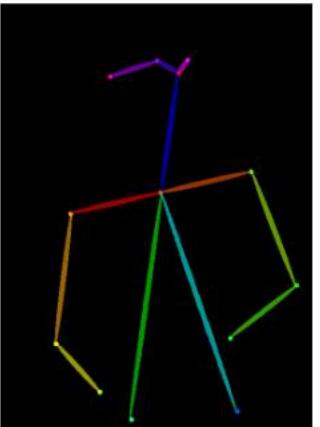
Default



"masterpiece of fairy tale, giant deer, golden antlers"



"..., quaint city Galic"



Input human pose



Default



"chef in kitchen"



"Lincoln statue"

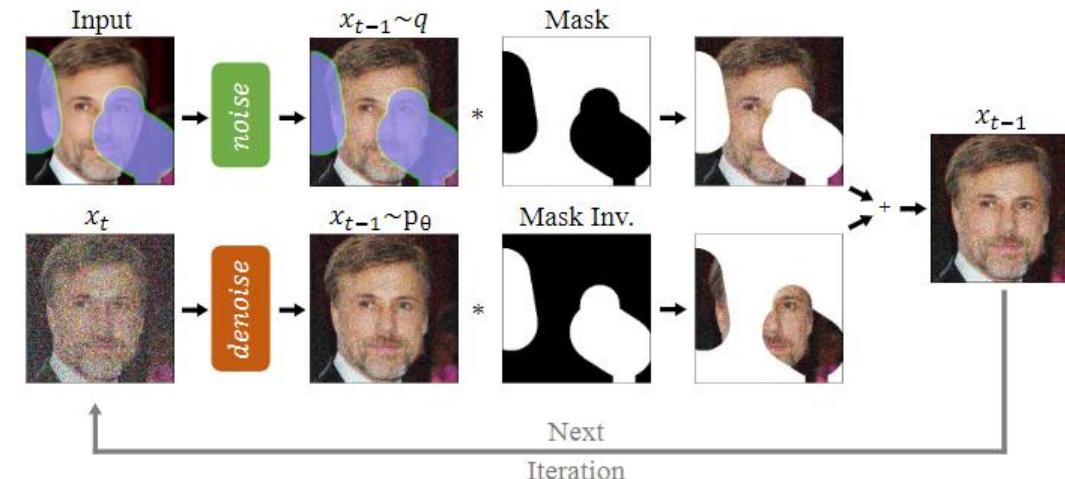
# Repaint

## Algorithm 1 Inpainting using our RePaint approach.

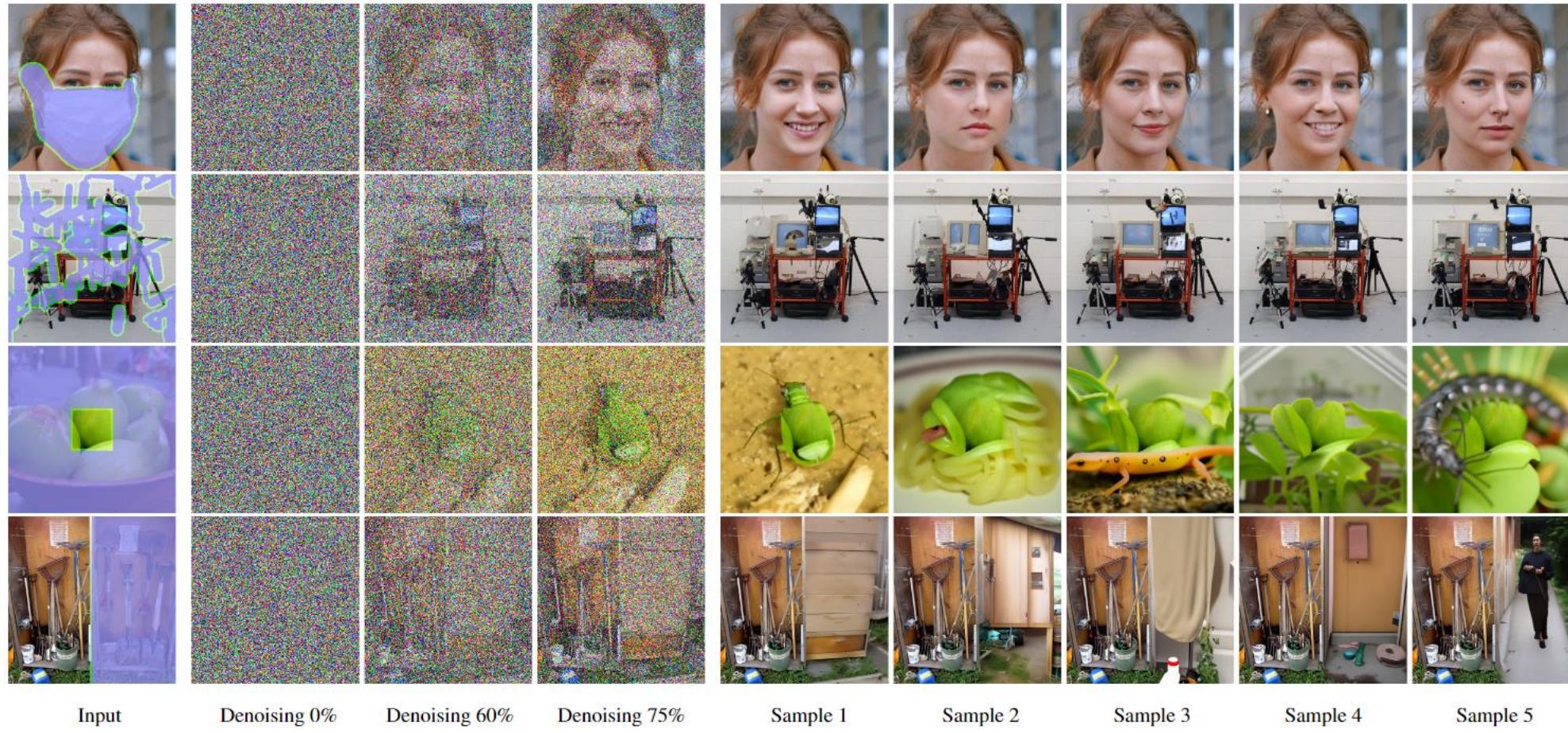
```

1:  $x_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:   for  $u = 1, \dots, U$  do
4:      $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\epsilon = \mathbf{0}$ 
5:      $x_{t-1}^{\text{known}} = \sqrt{\alpha_t}x_0 + (1 - \bar{\alpha}_t)\epsilon$ 
6:      $z \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $z = \mathbf{0}$ 
7:      $x_{t-1}^{\text{unknown}} = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(x_t, t) \right) + \sigma_t z$ 
8:      $x_{t-1} = m \odot x_{t-1}^{\text{known}} + (1 - m) \odot x_{t-1}^{\text{unknown}}$ 
9:   if  $u < U$  and  $t > 1$  then
10:     $x_t \sim \mathcal{N}(\sqrt{1 - \beta_{t-1}}x_{t-1}, \beta_{t-1}\mathbf{I})$ 
11:   end if
12: end for
13: end for
14: return  $x_0$ 

```



# Repaint

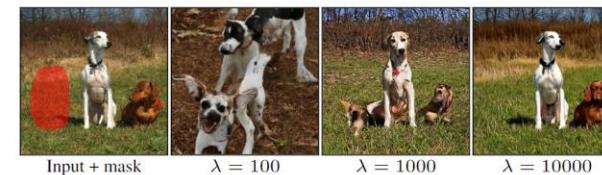


# Blended-diffusion



**Algorithm 1** Local CLIP-guided diffusion, given a diffusion model  $(\mu_\theta(x_t), \Sigma_\theta(x_t))$  and CLIP model

**Input:** source image  $x$ , target text description  $d$ , input mask  $m$ , diffusion steps  $k$ , background preservation coefficient  $\lambda$   
**Output:** edited image  $\widehat{x}$  that differs from input image  $x$  inside area  $m$  according to text description  $d$   
 $x_k \sim \mathcal{N}(\sqrt{\bar{\alpha}_k}x_0, (1 - \bar{\alpha}_k)\mathbf{I})$   
**for all**  $t$  from  $k$  to 1 **do**  
 $\mu, \Sigma \leftarrow \mu_\theta(x_t), \Sigma_\theta(x_t)$   
 $\widehat{x}_0 \leftarrow \frac{x_t}{\sqrt{\bar{\alpha}_t}} - \frac{\sqrt{1 - \bar{\alpha}_t}\epsilon_\theta(x_t, t)}{\sqrt{\bar{\alpha}_t}}$   
 $\widehat{x}_{0,aug} \leftarrow \text{ExtendingAugmentations}(\widehat{x}_0, N)$   
 $\mathcal{L} \leftarrow \mathcal{D}_{CLIP}(\widehat{x}_{0,aug}, d, m) + \lambda \mathcal{D}_{bg}(x, \widehat{x}_0, aug, m)$   
 $x_{t-1} \sim \mathcal{N}(\mu + \Sigma \nabla \widehat{x}_0 \mathcal{L}, \Sigma)$   
**end for**  
**return**  $x_0$

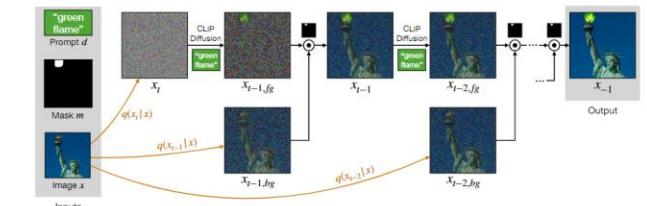


**Figure 3. Effect of  $\lambda$  in local CLIP-guided diffusion.** Given an input image with a mask, and the prompt “a dog”: with  $\lambda$  set too low ( $\lambda = 100$ ), the entire image changes completely, while if  $\lambda$  is too high ( $\lambda = 10000$ ), the model fails to change the foreground (and the background preservation is not perfect). Using an intermediate value ( $\lambda = 1000$ ) the model changes the foreground while resembling the original background (zoom for more details).

The above images are from an Intrinsic Comparison

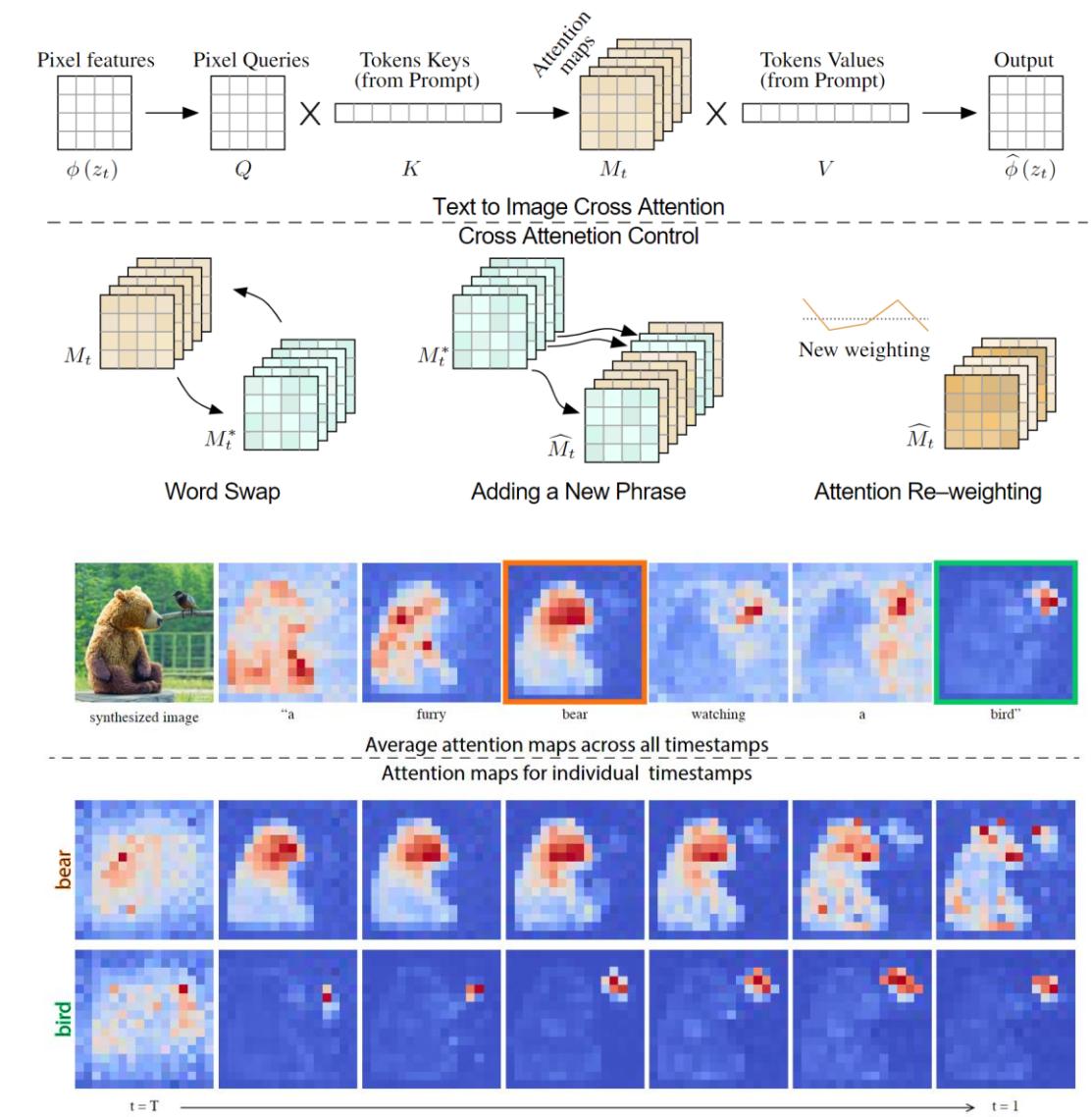
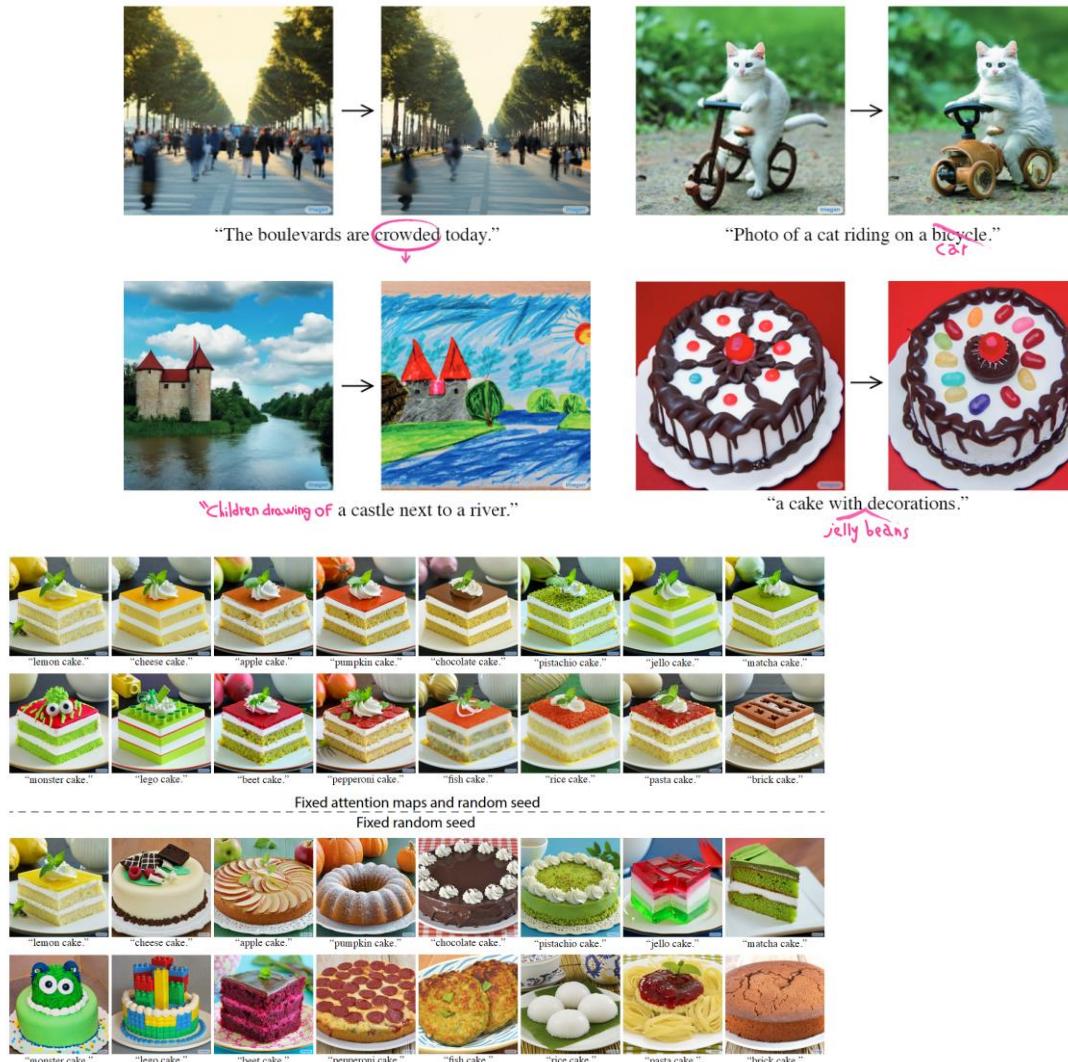
**Algorithm 2** Text-driven blended diffusion: given a diffusion model  $(\mu_\theta(x_t), \Sigma_\theta(x_t))$ , and CLIP model

**Input:** source image  $x$ , target text description  $d$ , input mask  $m$ , diffusion steps  $k$ , number of extending augmentations  $N$   
**Output:** edited image  $\widehat{x}$  that differs from input image  $x$  inside area  $m$  according to text description  $d$   
 $x_k \sim \mathcal{N}(\sqrt{\bar{\alpha}_k}x_0, (1 - \bar{\alpha}_k)\mathbf{I})$   
**for all**  $t$  from  $k$  to 0 **do**  
 $\mu, \Sigma \leftarrow \mu_\theta(x_t), \Sigma_\theta(x_t)$   
 $\widehat{x}_0 \leftarrow \frac{x_t}{\sqrt{\bar{\alpha}_t}} - \frac{\sqrt{1 - \bar{\alpha}_t}\epsilon_\theta(x_t, t)}{\sqrt{\bar{\alpha}_t}}$   
 $\widehat{x}_{0,aug} \leftarrow \text{ExtendingAugmentations}(\widehat{x}_0, N)$   
 $\nabla_{text} \leftarrow \frac{1}{N} \sum_{i=1}^N \nabla_{\widehat{x}_{0,aug}} \mathcal{D}_{CLIP}(\widehat{x}_{0,aug}, d, m)$   
 $x_{fg} \sim \mathcal{N}(\mu + \Sigma \nabla_{text}, \Sigma)$   
 $x_{bg} \sim \mathcal{N}(\sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I})$   
 $x_{t-1} \leftarrow x_{fg} \odot m + x_{bg} \odot (1 - m)$   
**end for**  
**return**  $x_{-1}$

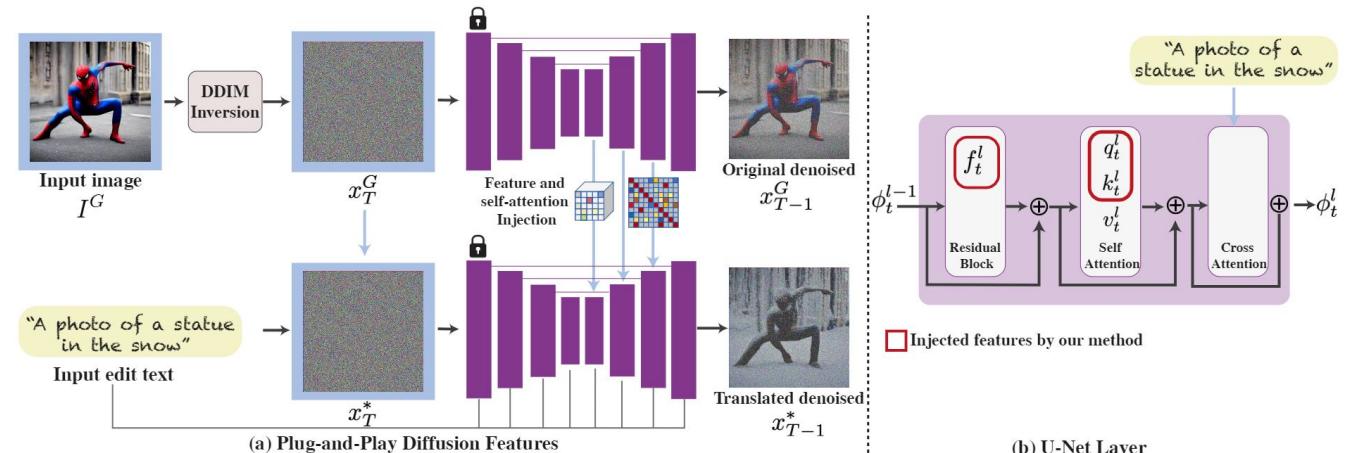


**Figure 4. Text-driven blended diffusion.** Given input image  $x$ , input mask  $m$ , and a text prompt  $d$ , we leverage the diffusion process to edit the image locally and coherently. We denote with  $\odot$  the element-wise blending of two images using the input mask  $m$ .

P2p



# Pnp




---

**Algorithm 1** Plug-and-Play Diffusion Features
 

---

**Inputs:**

$I^G$  ▷ real guidance image

$P$  ▷ target text prompt

$\tau_f, \tau_A$  ▷ injection thresholds

```

 $x_T^G \leftarrow \text{DDIM-inv}(I^G)$ 
 $x_T^* \leftarrow x_T^G$                                 ▷ Starting from same seed
 $\text{for } t \leftarrow T \dots 1 \text{ do}$ 
     $z_{t-1}^G, f_t^4, \{A_t^l\} \leftarrow \epsilon_\theta(x_t^G, \emptyset, t)$ 
     $x_{t-1}^G \leftarrow \text{DDIM-samp}(x_t^G, z_{t-1}^G)$ 
     $\text{if } t > \tau_f \text{ then } f_t^{*4} \leftarrow f_t^4 \text{ else } f_t^{*4} \leftarrow \emptyset$ 
     $\text{if } t > \tau_A \text{ then } A_t^{*l} \leftarrow A_t^l \text{ else } A_t^{*l} \leftarrow \emptyset$ 
     $z_{t-1}^* \leftarrow \hat{\epsilon}_\theta(x_t^*, P, t; f_t^{*4}, \{A_t^{*l}\})$ 
     $x_{t-1}^* \leftarrow \text{DDIM-samp}(x_t^*, z_{t-1}^*)$ 
 $\text{end for}$ 
Output:  $I^* \leftarrow x_0^*$ 
  
```

---

# InstructPix2Pix



Figure 5. Mona Lisa transformed into various artistic mediums.



Figure 6. The Creation of Adam with new context and subjects (generated at 768 resolution).



Figure 7. The iconic Beatles Abbey Road album cover transformed in a variety of ways.



$$\begin{aligned}\tilde{e}_\theta(z_t, c_I, c_T) &= e_\theta(z_t, \emptyset, \emptyset) \\ &+ s_I \cdot (e_\theta(z_t, c_I, \emptyset) - e_\theta(z_t, \emptyset, \emptyset)) \\ &+ s_T \cdot (e_\theta(z_t, c_I, c_T) - e_\theta(z_t, c_I, \emptyset))\end{aligned}$$

# InstructPix2Pix



“Zoom into the image”

“Move it to Mars”

“Color the tie blue”

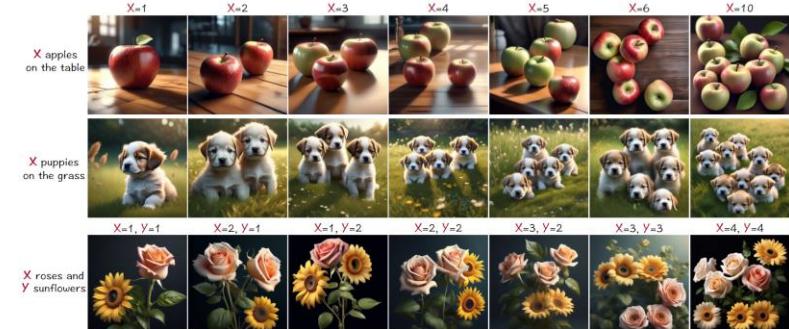
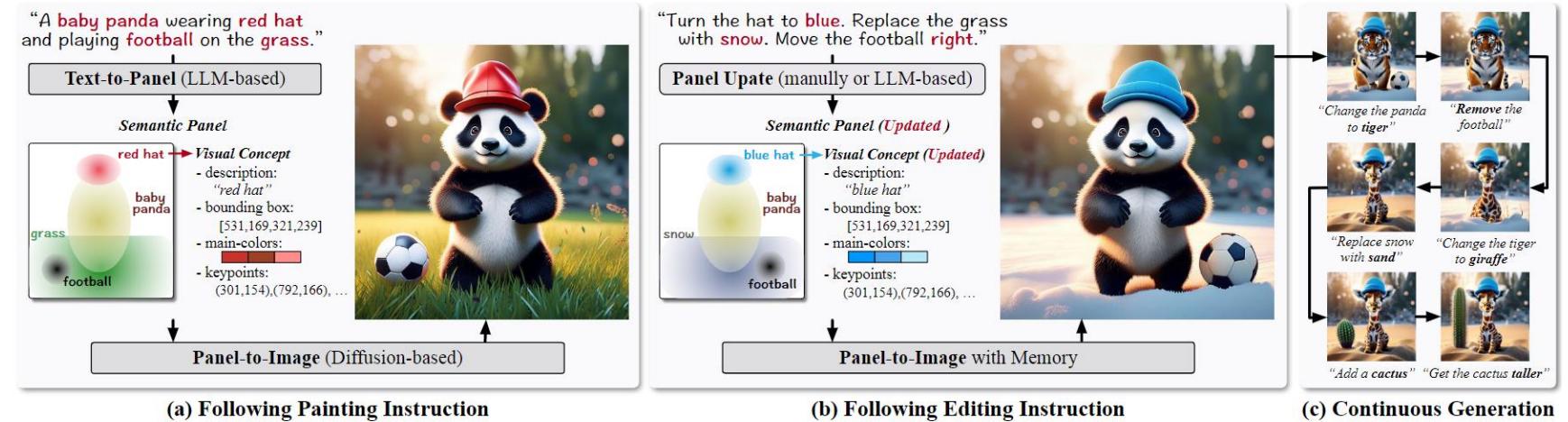
“Have the people swap places”

Figure 13. Failure cases. Left to right: our model is not capable of performing viewpoint changes, can make undesired excessive changes to the image, can sometimes fail to isolate the specified object, and has difficulty reorganizing or swapping objects with each other.

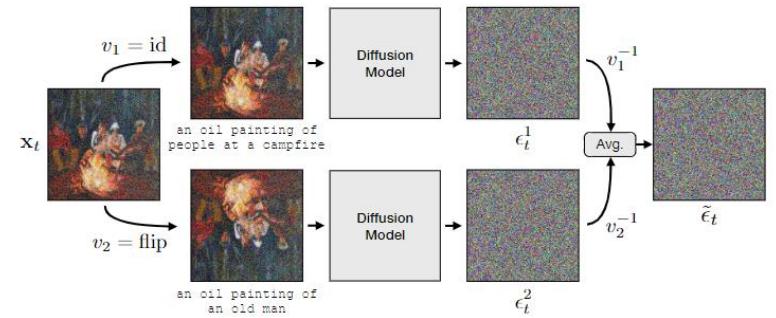
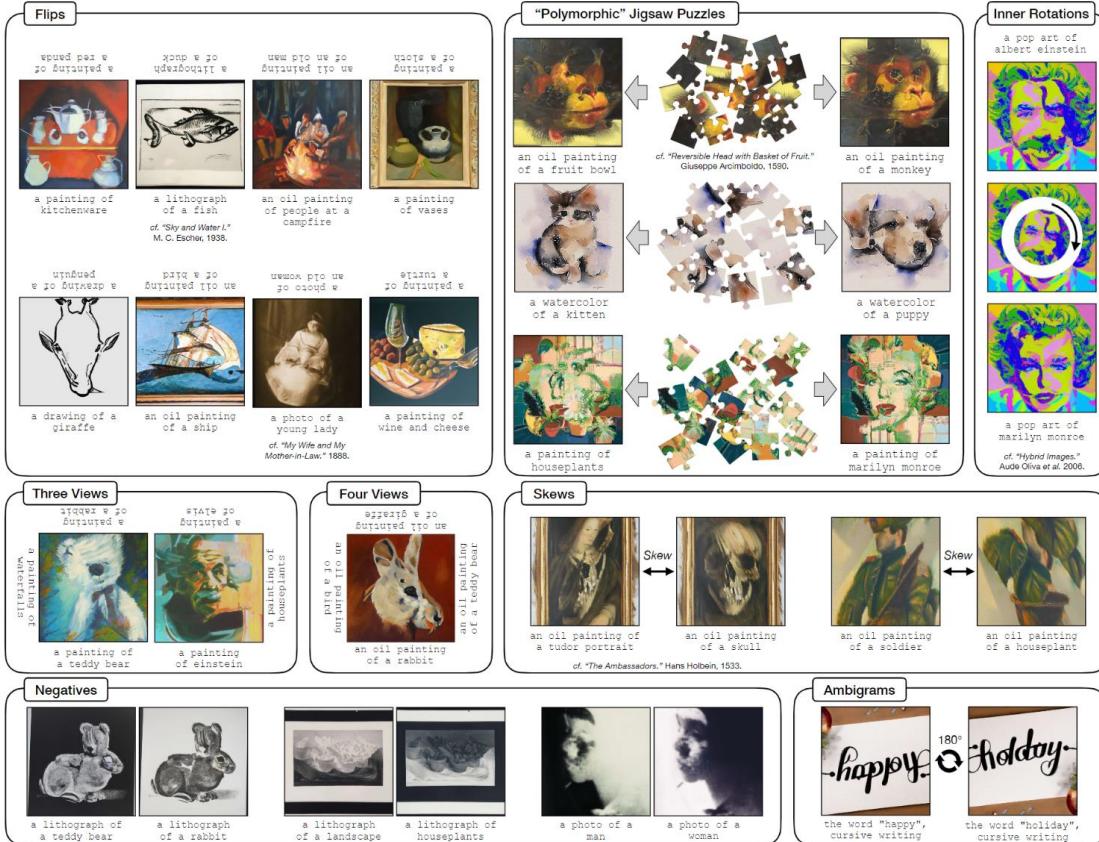
## CVPR24-oral

- 精准控制
- 复杂场景
- ...

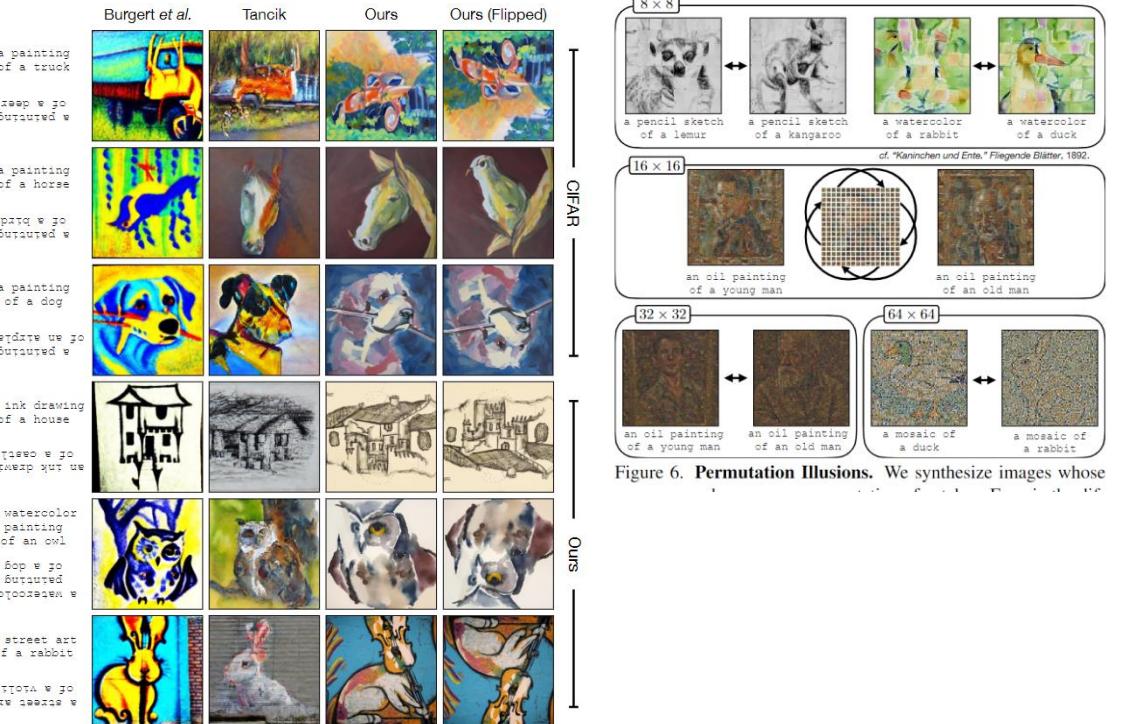
# Ranni

Figure 5. Samples generated by Ranni on **quantity-awareness** prompts.Figure 6. Samples generated by Ranni on **spatial relationship** prompts.

# Visual Anagrams

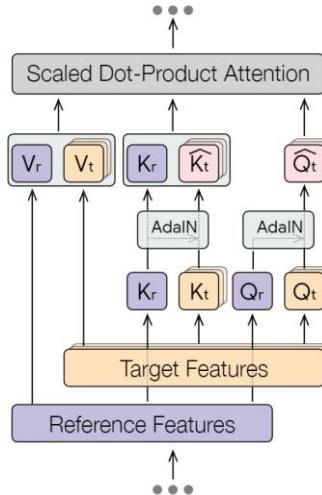
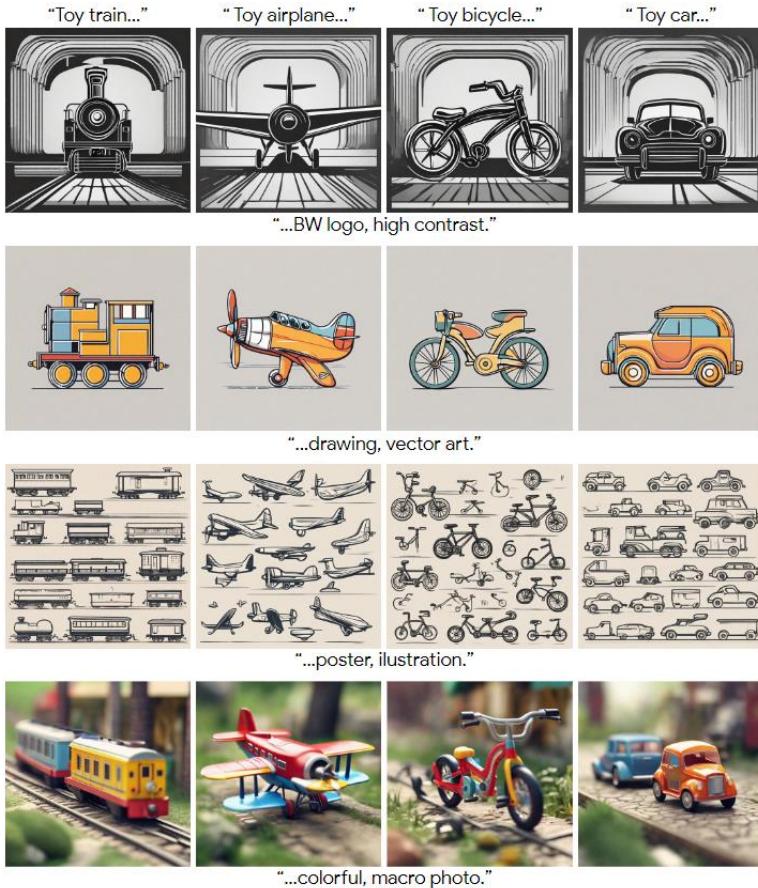


**Figure 2. Algorithm Overview.** Our method works by simultaneously denoising multiple views of an image. Given a noisy image  $\mathbf{x}_t$ , we compute noise estimates,  $\epsilon_t^i$ , conditioned on different prompts, after applying views  $v_i$ . We then apply the inverse view  $v_i^{-1}$  to align estimates, average the estimates, and perform a reverse diffusion step. The final output is an optical illusion.



**Figure 6. Permutation Illusions.** We synthesize images whose

# StyleAligned



**Figure 4. Shared attention layer.** The target images attend to the reference image by applying AdaIN over their queries and keys using the reference queries and keys respectively. Then, we apply shared attention where the target features are updated by both the target values  $V_t$  and the reference values  $V_r$ .

As illustrated in Fig. 4, to enable balanced attention reference, we normalize the queries  $Q_t$  and keys  $K_t$  of the target image using the queries  $Q_r$  and keys  $K_r$  of the reference image using the adaptive normalization operation (AdaIN) [26]:

$$\hat{Q}_t = \text{AdaIN}(Q_t, Q_r) \quad \hat{K}_t = \text{AdaIN}(K_t, K_r),$$

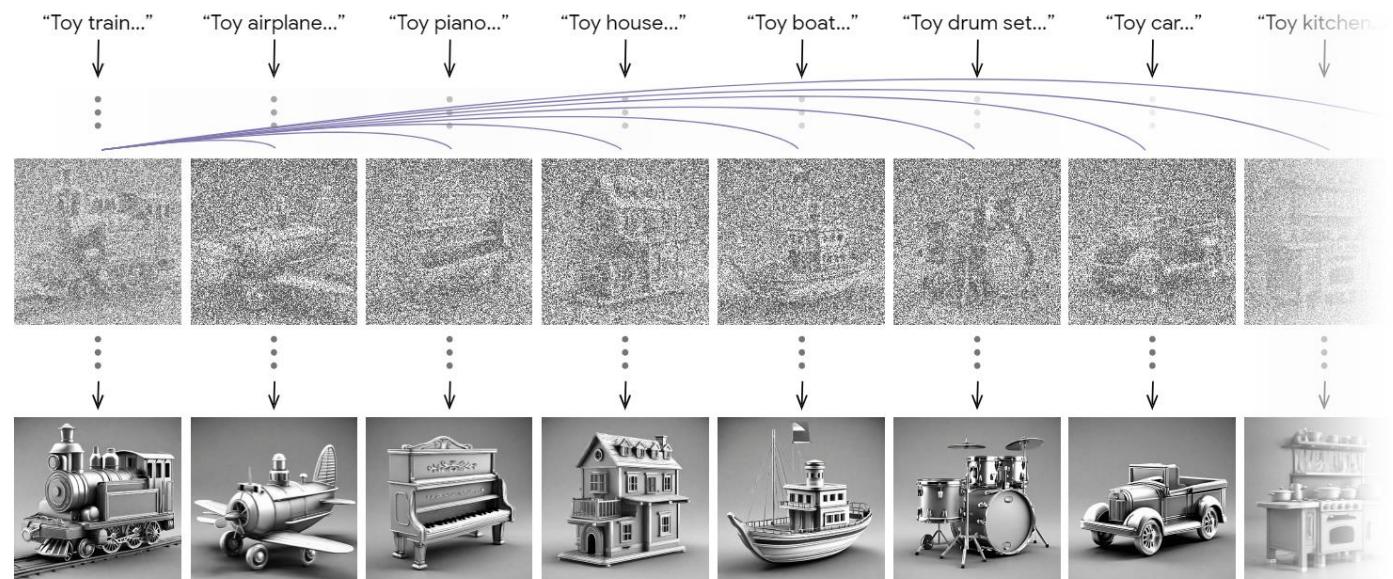
where the AdaIn operation is given by:

$$\text{AdaIN}(x, y) = \sigma(y) \left( \frac{x - \mu(x)}{\sigma(x)} \right) + \mu_y,$$

and  $\mu(x), \sigma(x) \in \mathbb{R}^{d_k}$  are the mean and the standard deviation of queries and keys across different pixels. Finally, our shared attention is given by

$$\text{Attention}(\hat{Q}_t, K_{rt}^T, V_{rt}),$$

$$\text{where } K_{rt} = \begin{bmatrix} K_r \\ \hat{K}_t \end{bmatrix} \text{ and } V_{rt} = \begin{bmatrix} V_r \\ V_t \end{bmatrix}.$$



**Figure 3. Style Aligned Diffusion.** Generation of images with a style aligned to the reference image on the left. In each diffusion denoising step all the images, except the reference, perform a shared self-attention with the reference image.

# Instruct-Imagen

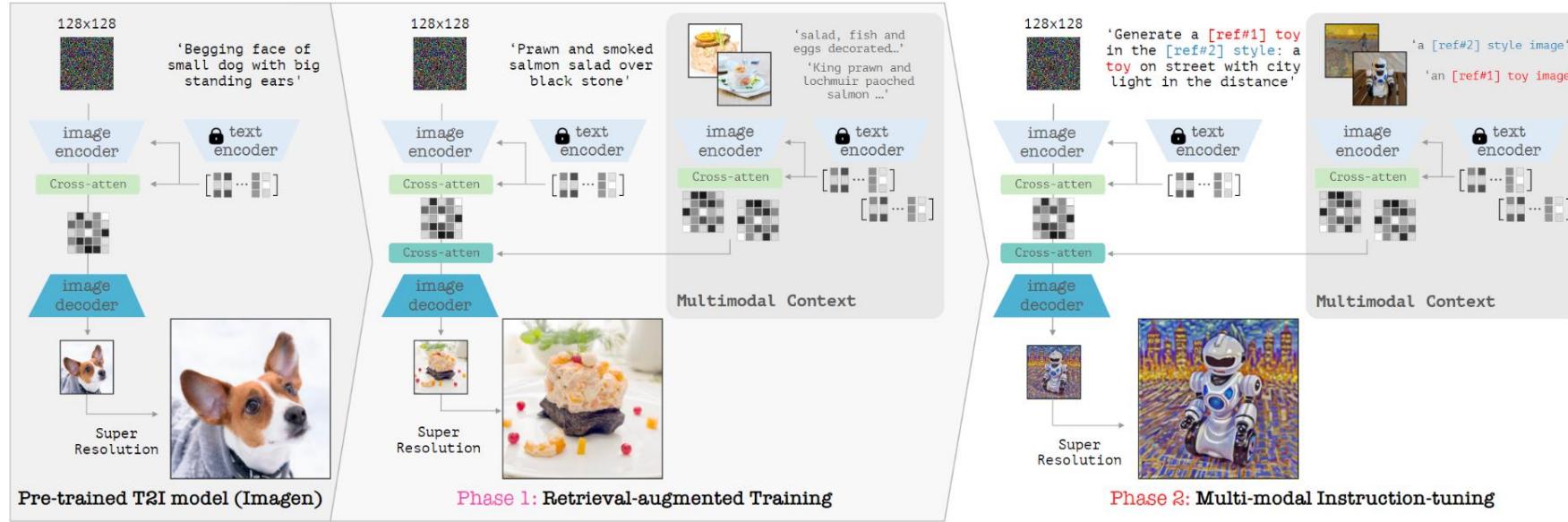
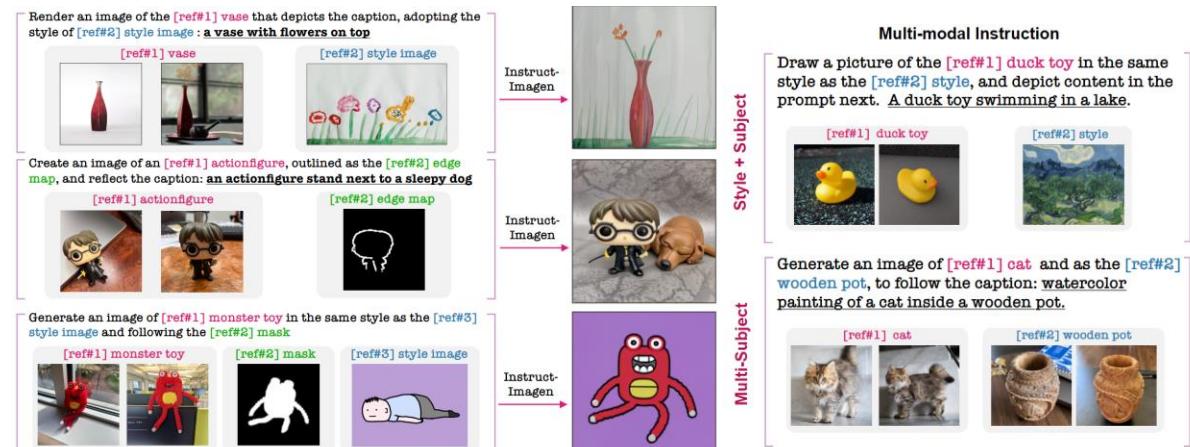


Figure 3. Overview of the two-staged training pipeline for the proposed Instruct-Imagen model.



Instruct-Imagen : Image Generation with Multi-modal Instruction. (n.d.).

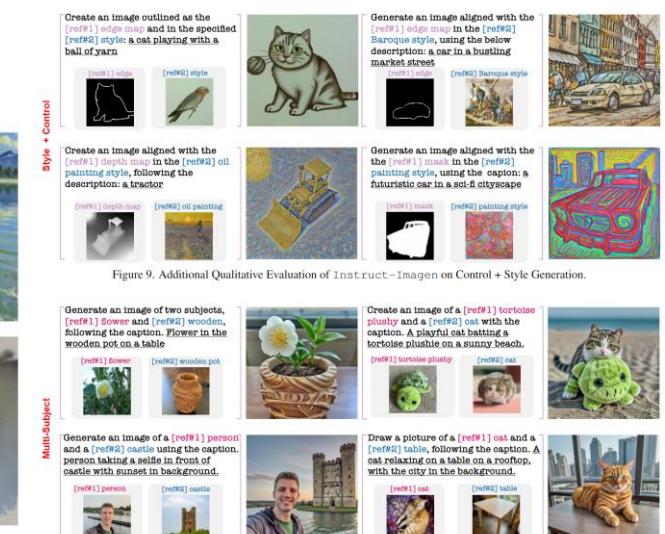
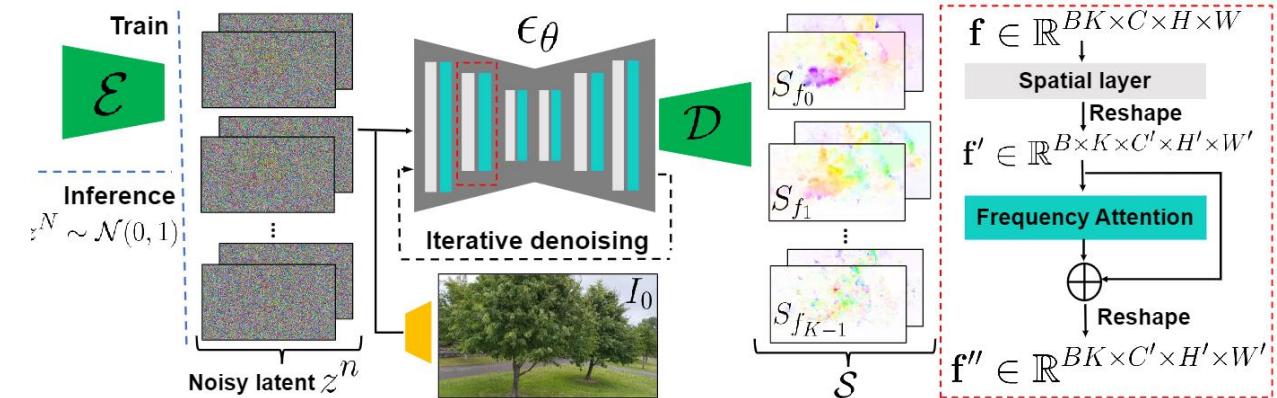


Figure 9. Additional Qualitative Evaluation of Instruct-Imagen on Control + Style Generation.



Figure 10. Additional Qualitative Evaluation of Instruct-Imagen on Multi-Subject Generation.

# Best paper |



# Best paper ||

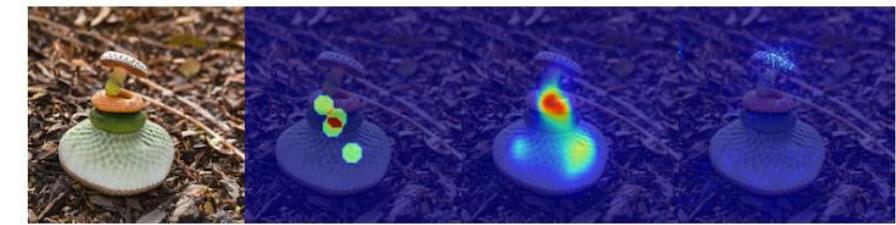


**Figure 1. An illustration of our annotation UI.** Annotators mark points on the image to indicate artifact/implausibility regions (red points) or misaligned regions (blue points) w.r.t the text prompt. Then, they click on the words to mark the misaligned keywords (underlined and shaded) and choose the scores for plausibility, text-image alignment, aesthetics, and overall quality (underlined).



(a) Image      (b) GT      (c) Our model      (d) ResNet-50

Figure 5. Examples of implausibility heatmaps. Prompt: *photo of a slim asian little girl ballerina with long hair wearing white tights running on a beach from behind nikon D5*



(a) Image      (b) GT      (c) Our model      (d) CLIP gradient

Figure 6. Examples of misalignment heatmaps. Prompt: *A snake on a mushroom.*

# 总结

- 研究广泛但不深入
  - 广泛：多种新奇任务、结合LLM
  - 不深入：设计新模型困难
- 机会仍然很多
  - 精准图像生成/编辑
  - Image-video、text-video
  - Text-3D
  - Text-4D