# Diffusion Model

A brief introduction & Applications

Huawei Fan    Cheng Wan    Anqi Zeng    Jiasheng Xu
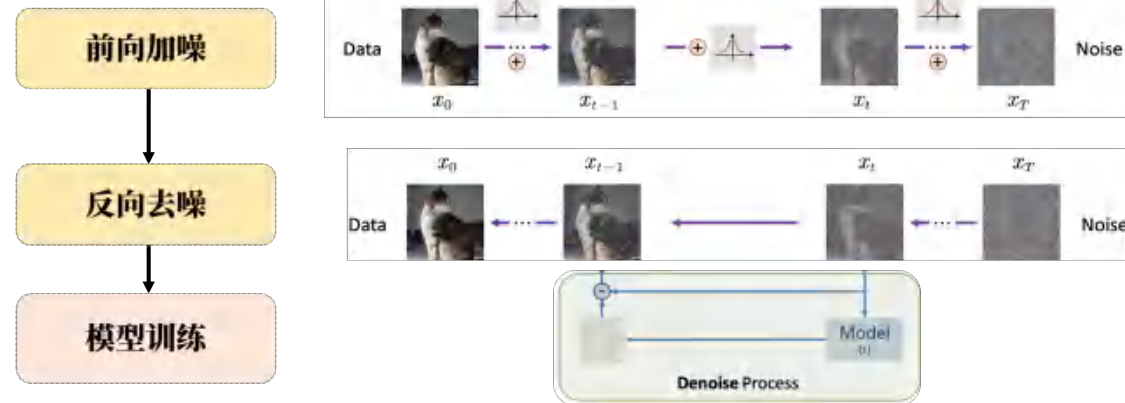
# Inspiration

# Structure

## Introduction



### What is Diffusion Model?

前向加噪

反向去噪

模型训练

## Application

### How Diffusion Model can be applied in different fields?

Stable Diffusio(LDM)

DALLE 2 By OpenAI

**TXT2IMG**

**LLM**

**Robotics**

# Brief Introduction

前向加噪

反向去噪

模型训练

# What is Diffusion Model?

## Denoising Diffusion Probabilistic Models
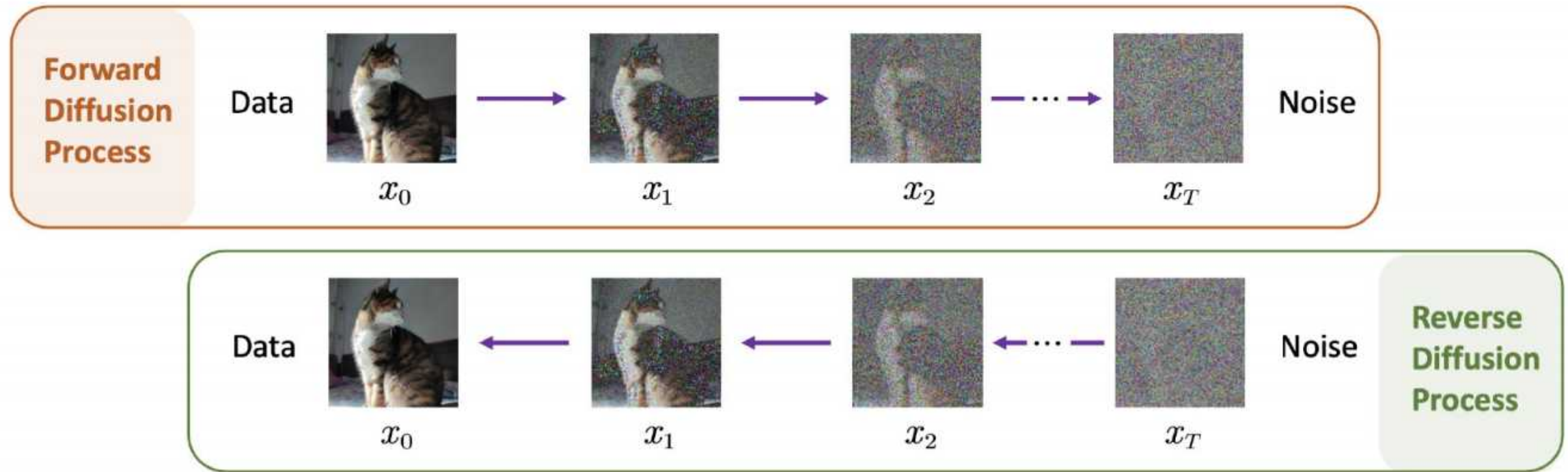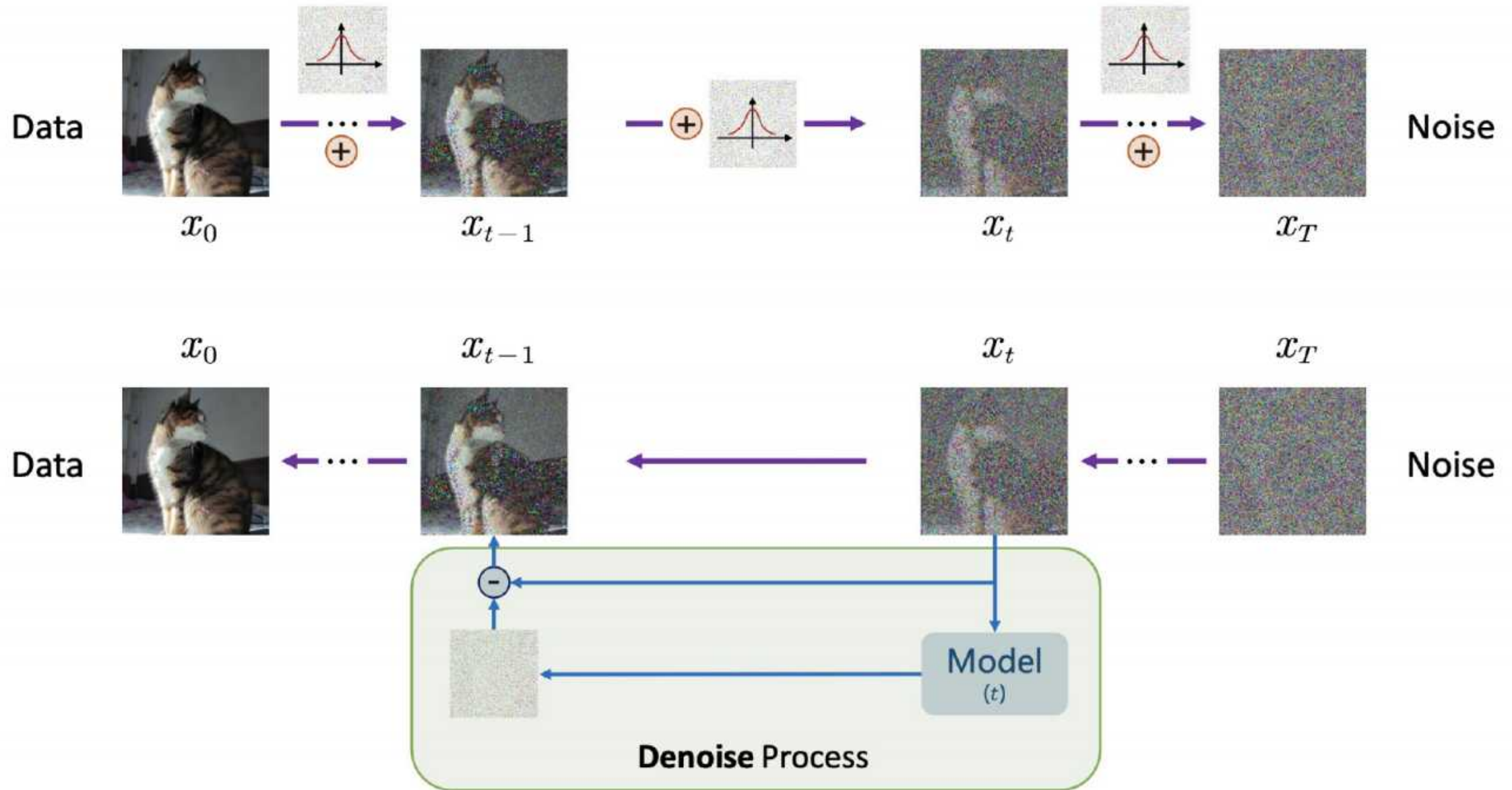
# What is Diffusion Model?



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

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

Data     Noise

Model $(t)$

**Denoise** Process

# Forward Diffusion Process



Data ⋯ → $x_0$ → $x_{t-1}$ → $q(\boldsymbol{x}_t | \boldsymbol{x}_{t-1})$ **Gauss distribution** → $x_t$ → ⋯ → $x_T$ Noise

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

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

$$x_t = \sqrt{\alpha_t}x_{t-1} + \sqrt{1-\alpha_t}\varepsilon_{t-1}, \quad \varepsilon \sim \mathcal{N}(0, I)$$

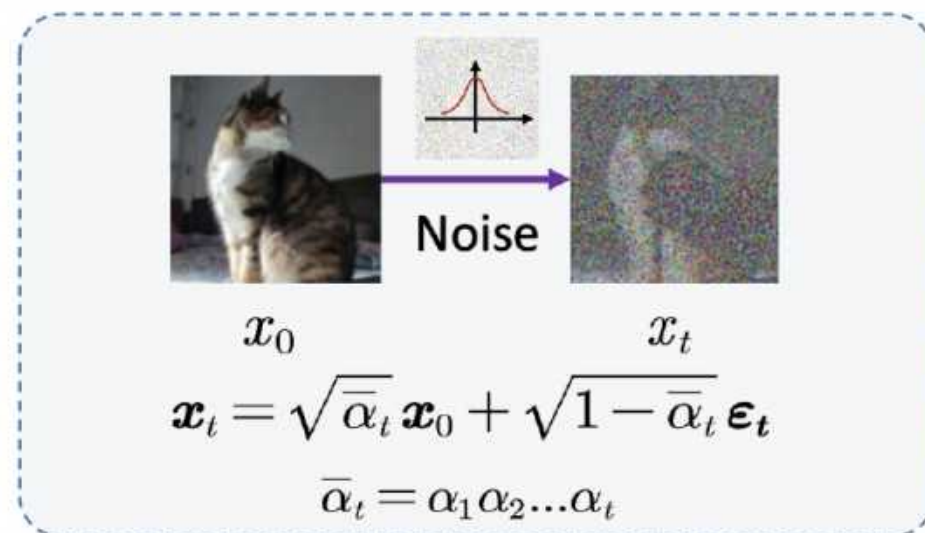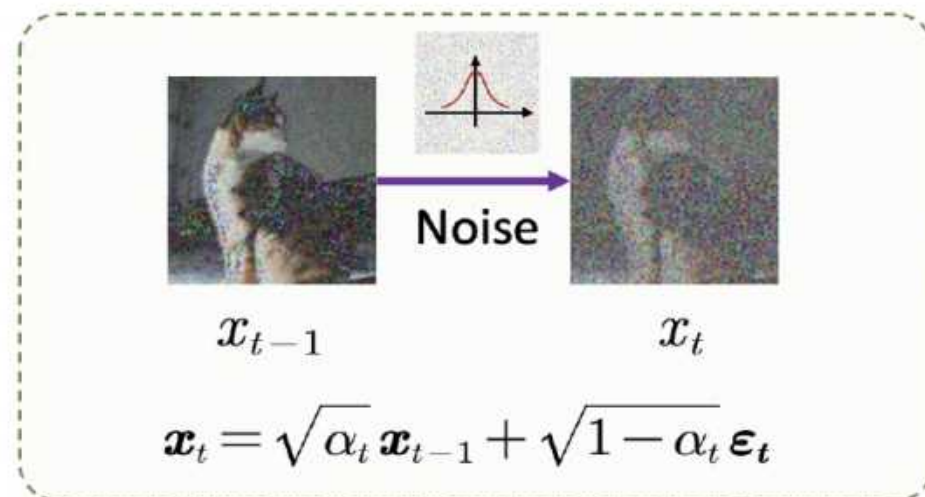$$q(x_t \mid x_{t-1}) = \mathcal{N}(x_t; \sqrt{\alpha_t}x_{t-1}, (1-\alpha_t)I)$$

# Forward Diffusion Process

## ① Forward (closed-form)
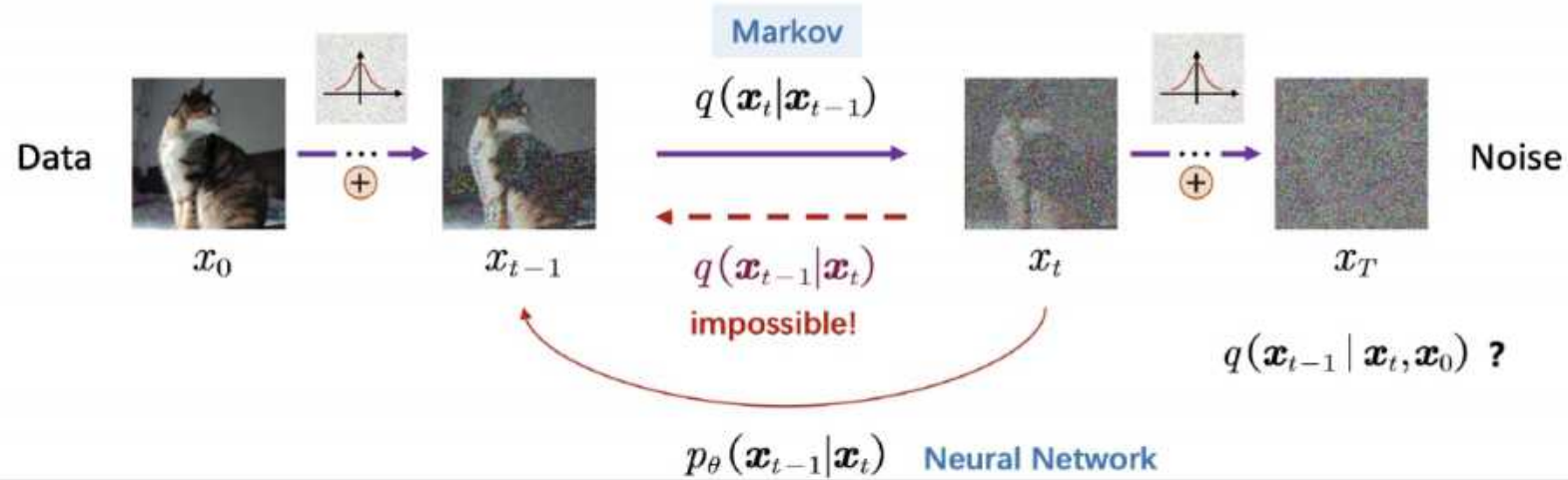
$$\mathbf{x}_t = \sqrt{\alpha_t}\mathbf{x}_{t-1} + \sqrt{1-\alpha_t}\boldsymbol{\varepsilon}_{t-1}$$

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

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

$$\vdots$$

$$= \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t}\boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, I)$$

Where $\bar{\alpha}_t = \prod_{s=1}^{t} \alpha_s$



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



$$x_0 \qquad x_t$$
$$\boldsymbol{x}_t = \sqrt{\bar{\alpha}_t}\,\boldsymbol{x}_0 + \sqrt{1-\bar{\alpha}_t}\,\boldsymbol{\varepsilon}_t$$
$$\bar{\alpha}_t = \alpha_1\alpha_2...\alpha_t$$

# Reverse Diffusion Process

## **Reverse** Diffusion Process



**Assume:** the output is gaussian

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

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

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



**If we know** $x_0$ **and** $x_t$

$q(x_{t-1} \mid x_t, x_0)$ **is deterministic**

**Assume: Markov**

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

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

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

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

③**Reverse**    If we know $x_0$

$$q(x_{t-1}|x_t, x_0) = \mathcal{N}\left(x_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1-\alpha_t)x_0}{1-\bar{\alpha}_t}}_{\mu_q(\mathbf{x}_t, t)}, \underbrace{\frac{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}}_{\Sigma_q(t)} I\right)$$
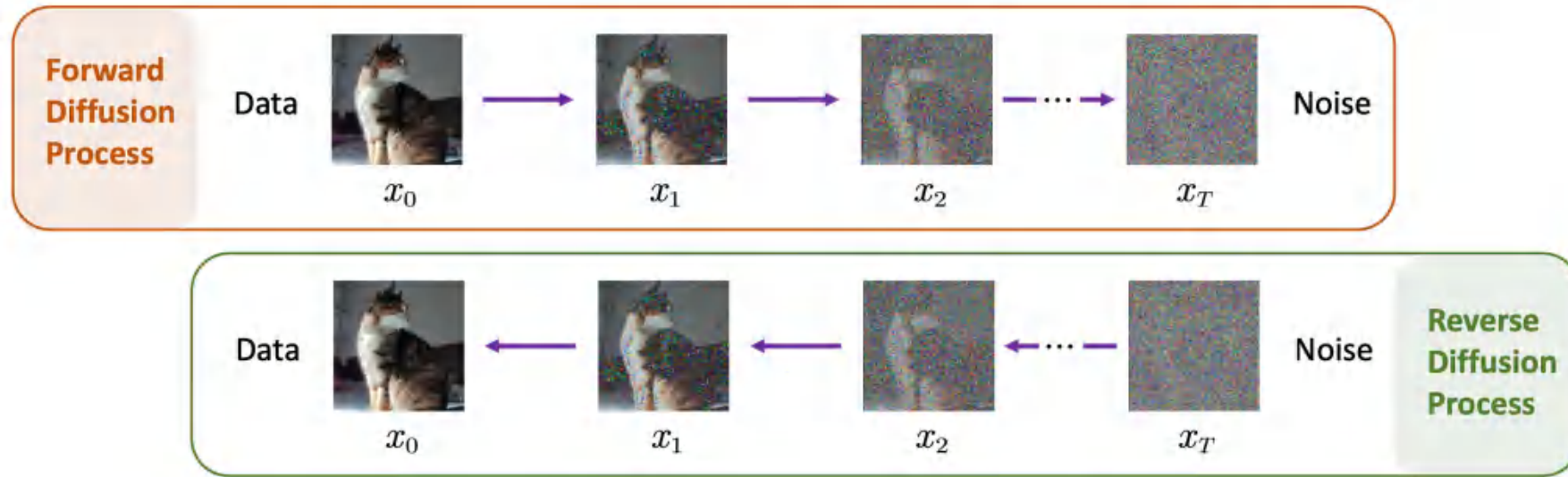
① **Forward (close-form)**

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1-\bar{\alpha}_t}\varepsilon_t \Rightarrow x_0 = \frac{x_t - \sqrt{1-\bar{\alpha}_t}\varepsilon_t}{\sqrt{\bar{\alpha}_t}}$$

$$\frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1-\alpha_t)x_0}{1-\bar{\alpha}_t} \Rightarrow \frac{1}{\sqrt{\bar{\alpha}_t}}\left(x_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}}\varepsilon_t\right) \qquad \text{noise predictor } \varepsilon_t(x_0 \to x_t)$$

Why not predict $x_0$ directly?
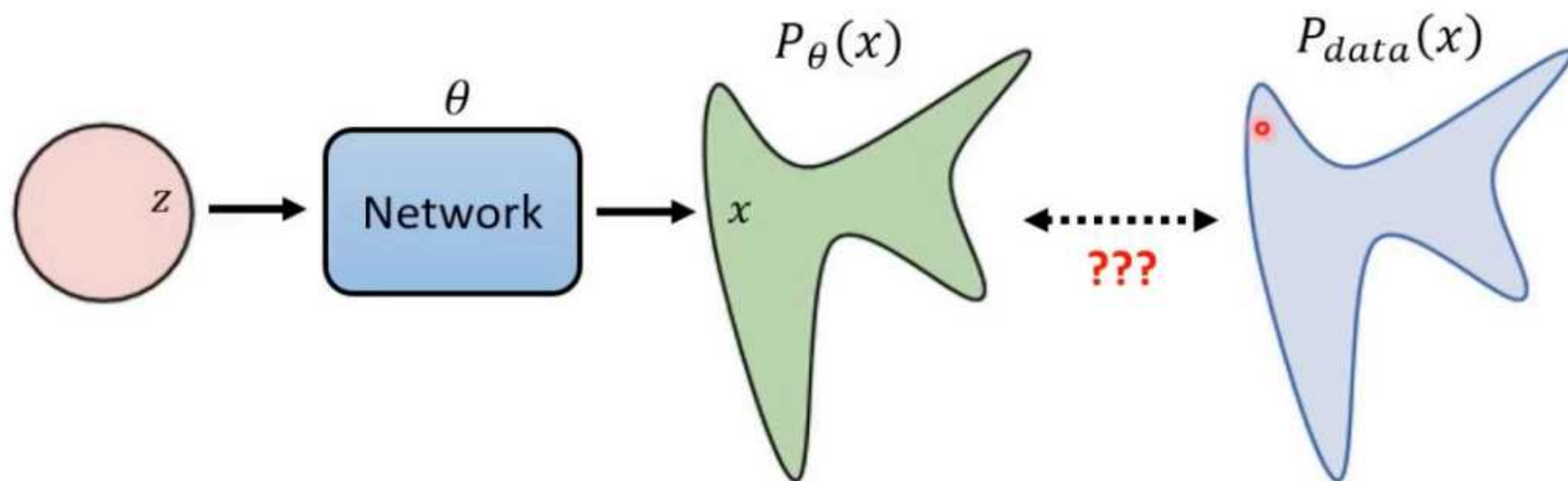
# Summary



$$Forward\ Process: \ q(x_t|x_{t-1}) = N(x_t; \sqrt{\alpha_t}x_{t-1}, (1-\alpha_t)I)$$

$$Reverse\ Process: \ p(x_{t-1}|x_t) = N\left(x_{t-1}; \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{1-\alpha_t}{\sqrt{1-\alpha_t}}\epsilon_t\right), \frac{1-\alpha_{t-1}}{1-\alpha_t}\beta_t I\right)$$

# Reverse Diffusion Process



Maximum Likelihood Estimation

$P_\theta(x)$

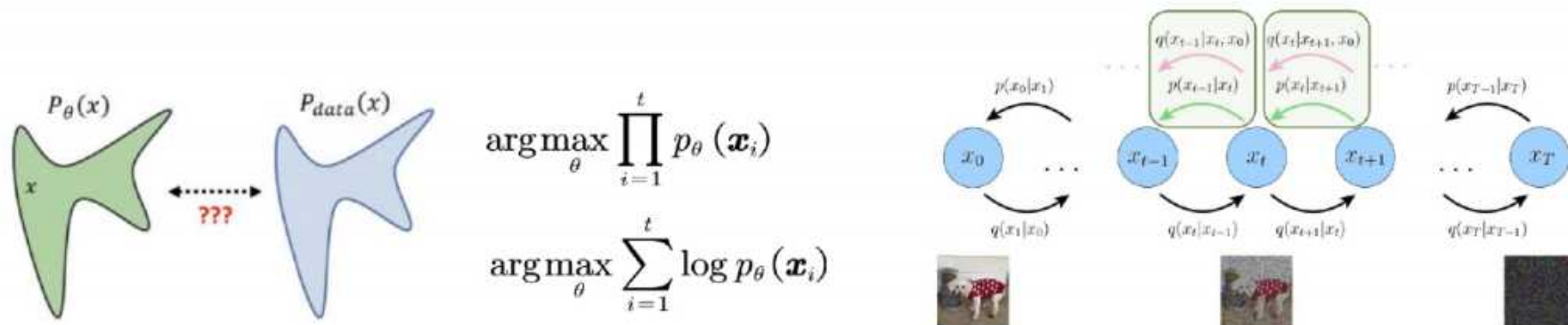$P_{data}(x)$

$\theta$

Network

$z$

$x$

???

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

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

???

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

# Maximum Likelihood Estimation



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

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

## ②Optimization (view 1)

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

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

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

# Derivation Process

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

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

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

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

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

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

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

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

$$= \mathbb{E}_q [ \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p_\theta(\mathbf{x}_T))}_{L_T} + \sum_{t=2}^{T} \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))}_{L_{t-1}} - \underbrace{\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)}_{L_0} ]$$

# Derivation Process

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

$$L_{t-1} - C = \mathbb{E}_{\mathbf{x}_0, \boldsymbol{\epsilon}} \left[ \frac{1}{2\sigma_t^2} \left\| \tilde{\boldsymbol{\mu}}_t \left( \mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\epsilon}), \frac{1}{\sqrt{\bar{\alpha}_t}} (\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\epsilon}) - \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}) \right) - \boldsymbol{\mu}_\theta(\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\epsilon}), t) \right\|^2 \right]$$

$$= \mathbb{E}_{\mathbf{x}_0, \boldsymbol{\epsilon}} \left[ \frac{1}{2\sigma_t^2} \left\| \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\epsilon}) - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon} \right) - \boldsymbol{\mu}_\theta(\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\epsilon}), t) \right\|^2 \right]$$

$$\boldsymbol{\mu}_\theta(\mathbf{x}_t, t) = \tilde{\boldsymbol{\mu}}_t \left( \mathbf{x}_t, \frac{1}{\sqrt{\bar{\alpha}_t}} (\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t)) \right) = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right)$$

$$\mathbb{E}_{\mathbf{x}_0, \boldsymbol{\epsilon}} \left[ \frac{\beta_t^2}{2\sigma_t^2 \alpha_t (1 - \bar{\alpha}_t)} \| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t) \|^2 \right]$$

$$L_{\text{simple}}(\theta) := \mathbb{E}_{t, \mathbf{x}_0, \boldsymbol{\epsilon}} \left[ \| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t) \|^2 \right]$$

# Diffusion Model Test by Pytorch

# Applications

| TXT2IMG | LLM | Robotics |
|---------|-----|----------|
| LDM&DALLE | LLaDA | Diffusion Policy |

什么是txt2img?

通过给定文本提示词（text prompt）输出一张匹配提示词的图片。



Stable Diffusio(LDM)

High-Resolution Image Synthesis with Latent Diffusion Models



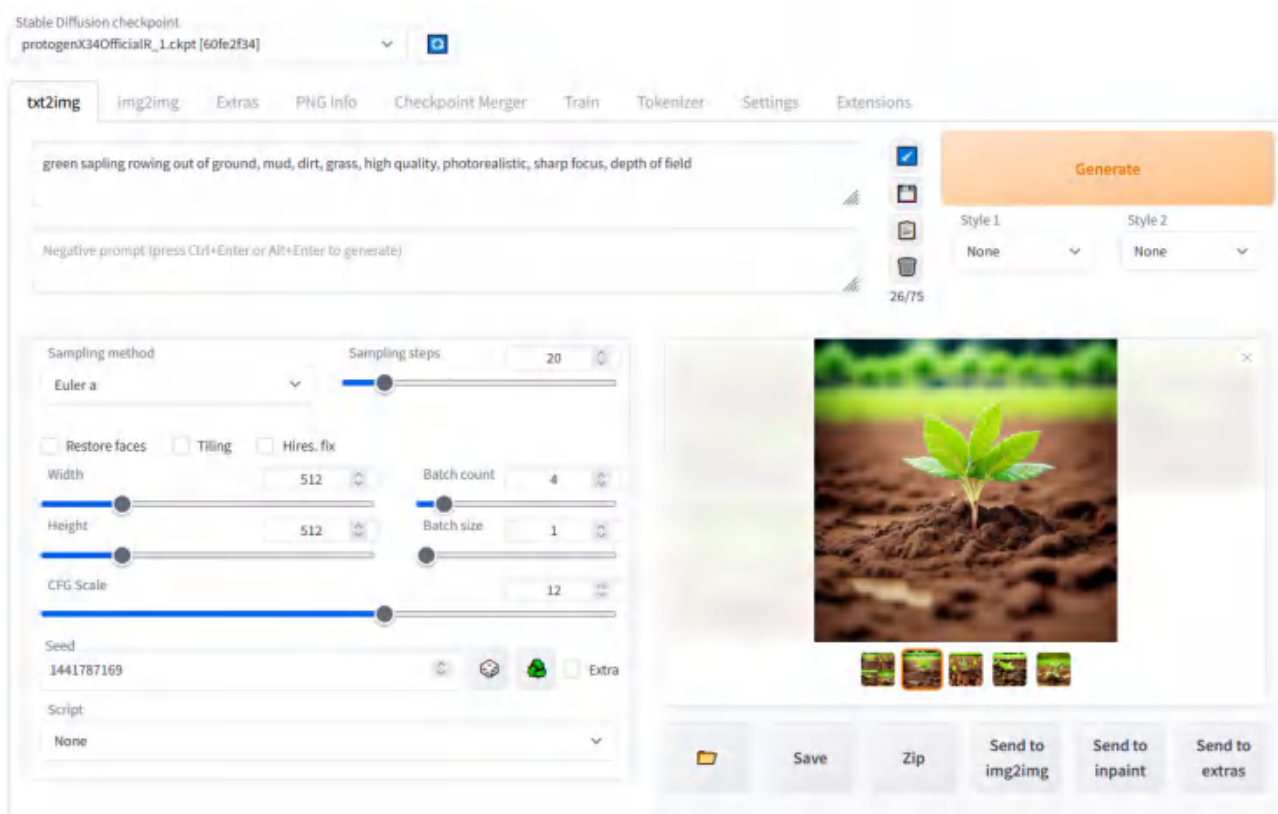DALLE 2 By OpenAI

Hierarchical Text-Conditional Image Generation with CLIP Latents
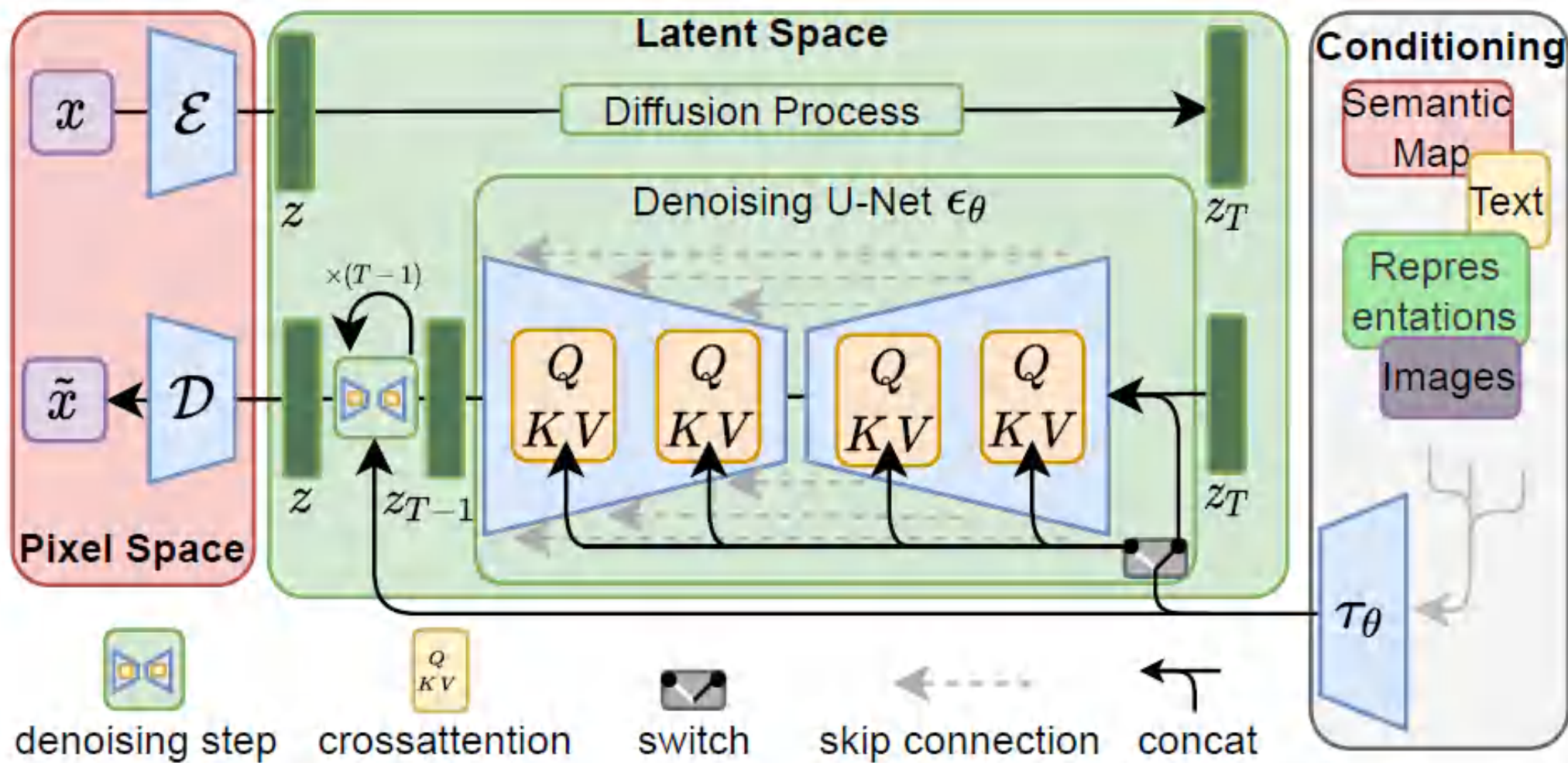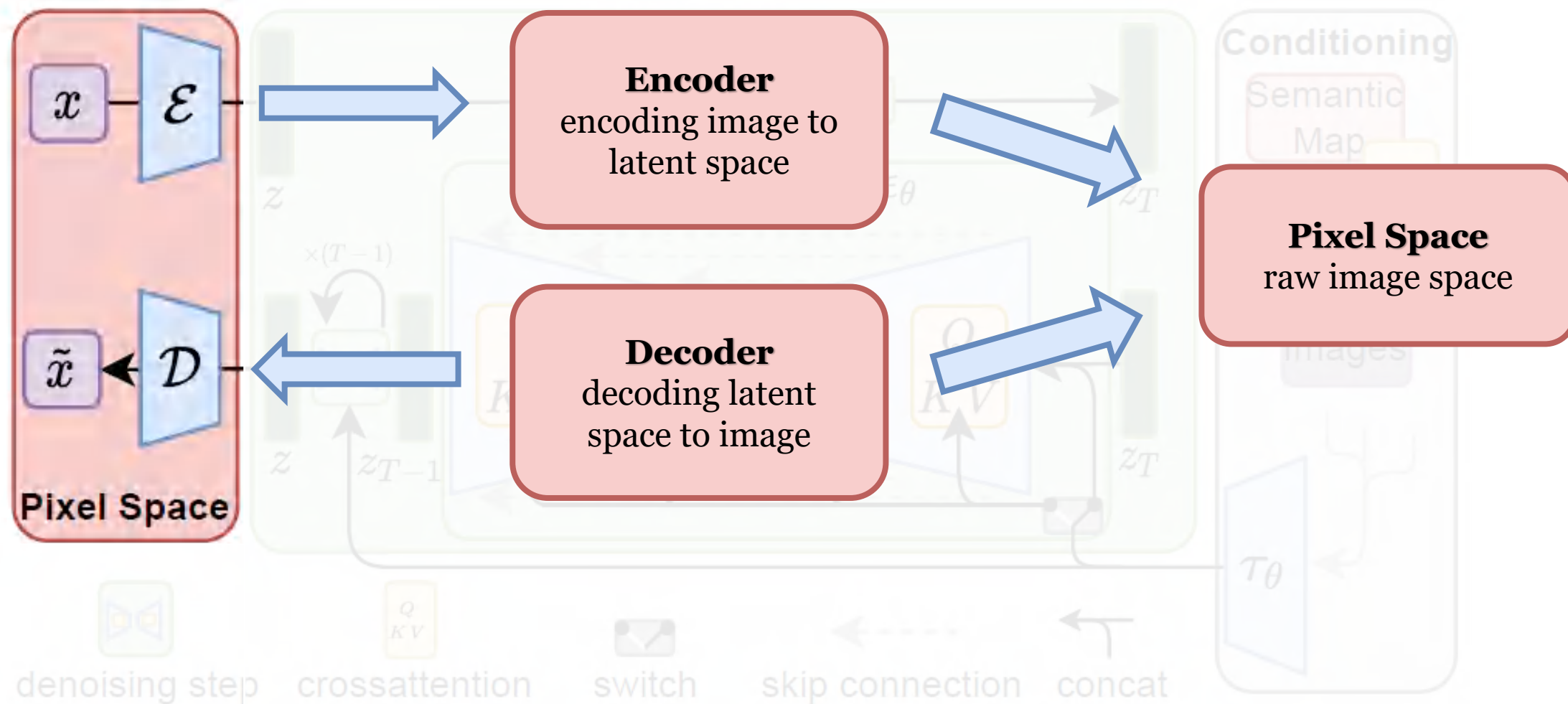
## Have a try! Open Source



**本地部署**

**stability.ai**

**MidJourney**

**网站访问**

High-Resolution Image Synthesis with Latent Diffusion Models

**Encoder**
encoding image to latent space

**Decoder**
decoding latent space to image

**Pixel Space**
raw image space

Pixel Space

denoising step    crossattention    switch    skip connection    concat

Conditioning
Influencing model output

Semantic Map
Visual representations

Text & Images
Textual and visual inputs

Representation:
Contextual conditioning

# Application 01 —— LDM 原理



Step 1
生成随机向量

Random tensor
in latent space

controlled by seed

Step 2 根据prompt预测噪声

Image In
Latent Space

Noise
Predictor
(UNet)

Text
Prompt

Predicted Noise
in Latent Space

Step 3 去除噪声

New Image = Latent Image - Predicted Noise

Step 4 Decode 生成图像

After N Samping

# Application 01 —— 图像生成 DALLE2

Have a try! API supported https://platform.openai.com/docs/overview



A photo of a white fur monster standing in a purple room

panda mad scientist mixing sparkling chemicals, artstation

vibrant portrait painting of Salvador Dalí with a robotic half face

扩散模型也能玩转大语言模型？ 　　Large Language Diffusion Models

What is LLaDA？ **L**arge **La**nguage **D**iffusion with m**A**sking

A text generation method different from the traditional **left-to-right** approach

## 传统路径：Auto Regression

每次生成一个token，新生成的token会拼到序列末尾
每个token的生成依赖于之前所有已生成的tokens

**Talk is Cheap，Show me the Code!**

## LLaDA

Mask的过程体现了diffusion加噪去噪的思想

# Performance

## Diffusion Policy

Image Diffusion

Action Diffusion

**Training Stability**

Compared to Energy Based Models

Diffusion Policy

IBC [1]

0.95

0.65

Checkpoint selection requires expensive realworld evaluation

# **Approach**

Step1 训练阶段：学习"去噪"过程　　经过多次迭代去噪，模型生成出符合当前观察条件的"干净"动作x0

$$x_{k-1} = \alpha(x_k - \gamma\epsilon_\theta(x_k, k) + N(0, \sigma^2 I))$$

$$L = MSE(\epsilon_k, \epsilon_\theta(x_0 + \epsilon_k, k))$$

去噪步长的缩放系数　　　噪声预测函数　　　　　　　　　　噪声预测函数

Step2 得分匹配来优化动作的能量

$$\nabla_a \log p(a|o) = -\nabla_a E_\theta(a, o)$$

**沿着降低能量的方向调整动作**

Step3 InfoNCE损失：区分"好"动作与"坏"动作

"好"动作

$$L_{infoNCE} = -\log\left(\frac{e^{-E_\theta(o,a)}}{e^{-E_\theta(o,a)} + \sum_{j=1}^{N_{neg}} e^{-E_\theta(o,\tilde{a}_j)}}\right)$$

"坏"动作

Step4 推理阶段：实时生成动作序列

# Reference



If you are interested in this topic, check these papers uploaded in DingTalk Group

# Thanks