

Diffusion and Consistency models in Latent Spaces

Jian Li
IIS, Tsinghua
2023



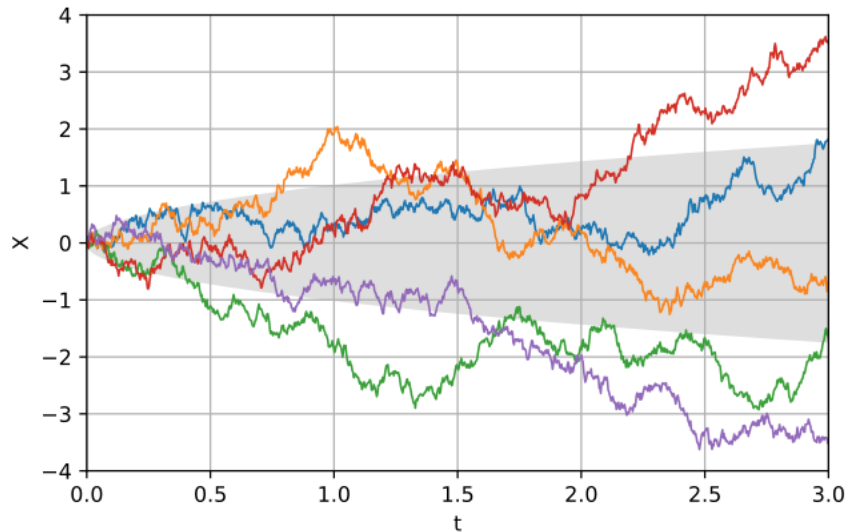
- Intro to Diffusion Process
- Intro to Generative Diffusion Models
- Probability Flow ODE
- Latent Diffusion Models (Stable Diffusion)
- Consistency Models
- Latent Consistency Models

Introduction to diffusion process

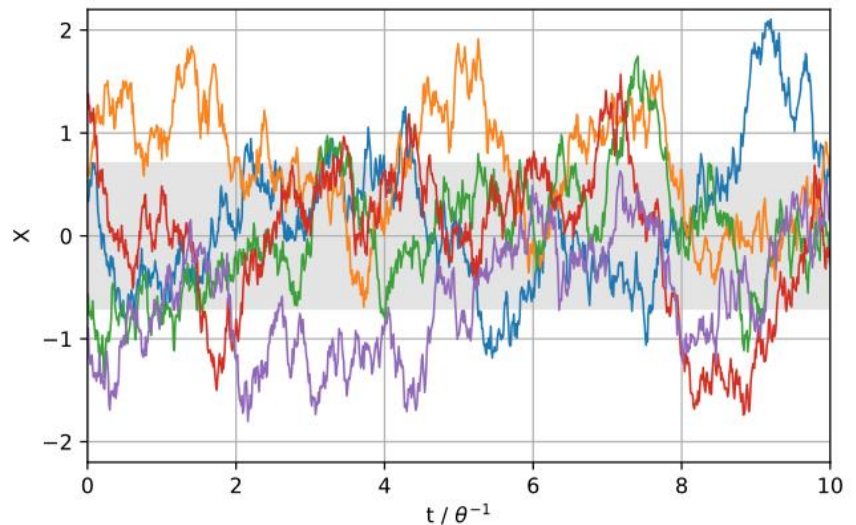
(Forward) Diffusion process (adding noise)

$$\text{SDE: } d\mathbf{x}_t = \underbrace{f(t)\mathbf{x}_t dt}_{\text{Drift term}} + \underbrace{g(t)d\mathbf{w}_t}_{\text{Diffusion term}}, \quad \mathbf{x}_0 \sim p_{\text{data}}(\mathbf{x}_0)$$

Brownian motion (Wiener process)



Standard 1-d Brownian motion
 $f=0, g=1$



Ornstein-Uhlenbeck Process
 $f=-1, g=1$

Introduction to diffusion process

Forward process (adding noise)

$$\text{SDE: } d\mathbf{x}_t = f(t)\mathbf{x}_t dt + g(t)d\mathbf{w}_t, \quad \mathbf{x}_0 \sim p_{data}(\mathbf{x}_0)$$

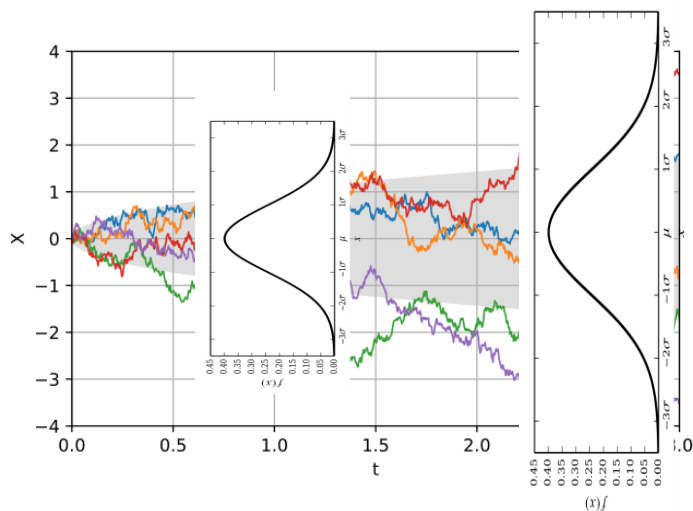
Reverse process (denoising)

Anderson's theorem

Reverse-time Brownian motion

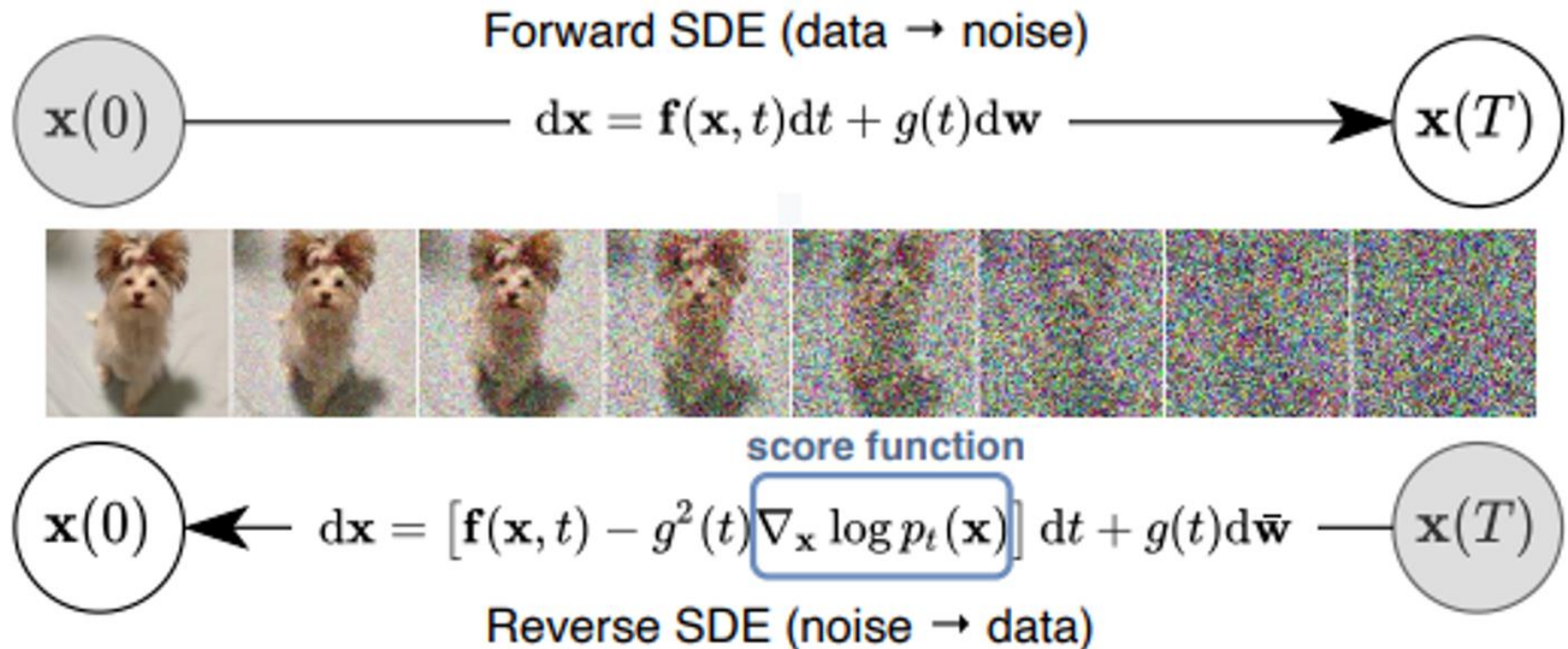
$$d\mathbf{x}_t = [f(t)\mathbf{x}_t - g^2(t)\nabla_{\mathbf{x}} \log q_t(\mathbf{x}_t)] dt + g(t)d\bar{\mathbf{w}}_t, \quad \mathbf{x}_T \sim q_T(\mathbf{x}_T)$$

Marginal probability distribution of x_t at time t

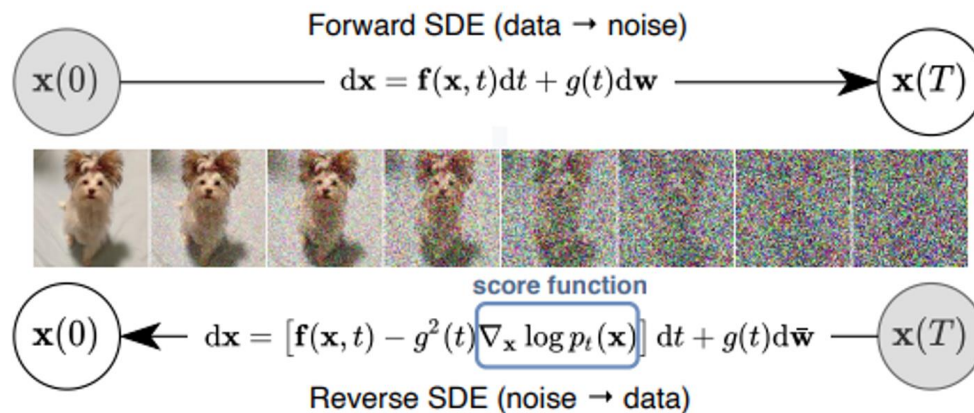


- Intro to Diffusion Process
- Intro to Generative Diffusion Models
- Probability Flow ODE
- Latent Diffusion Models (Stable Diffusion)
- Consistency Models
- Latent Consistency Models

Introduction to diffusion models



Score-based diffusion models



However, we do not know the score function. We only have training images.

Learn the score function from the training data!

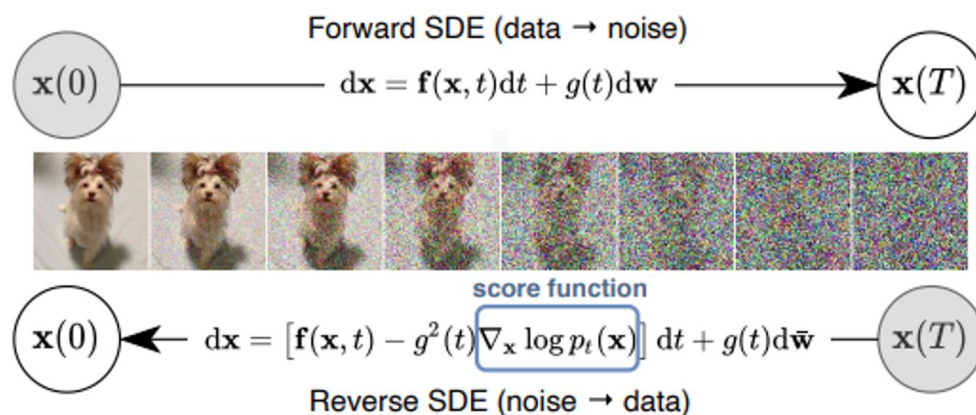
Score matching objective [Song & Ermon (2019)]:

$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N \sigma_i^2 \mathbb{E}_{p_{\text{data}}(\mathbf{x})} \mathbb{E}_{p_{\sigma_i}(\tilde{\mathbf{x}}|\mathbf{x})} \left[\|\mathbf{s}_{\theta}(\tilde{\mathbf{x}}, \sigma_i) - \nabla_{\tilde{\mathbf{x}}} \log p_{\sigma_i}(\tilde{\mathbf{x}} | \mathbf{x})\|_2^2 \right].$$

A deep neural network

Forward probability
we know how to compute

Score-based diffusion models



However, we do not know the score function. We only have training images.

Learn the score function from the training data!

Score-matching objective [Song & Ermon (2019)]:

$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N \sigma_i^2 \mathbb{E}_{p_{\text{data}}(\mathbf{x})} \mathbb{E}_{p_{\sigma_i}(\tilde{\mathbf{x}}|\mathbf{x})} \left[\left\| \mathbf{s}_{\theta}(\tilde{\mathbf{x}}, \sigma_i) - \nabla_{\tilde{\mathbf{x}}} \log p_{\sigma_i}(\tilde{\mathbf{x}} | \mathbf{x}) \right\|_2^2 \right].$$

A deep neural network

Forward probability
we know how to compute

Score-based diffusion models

- It is possible to further simplify the loss:

$$\begin{aligned}\mathcal{L}(\theta) &= \mathbb{E}_{t \in [0, T], \mathbf{x}_t \sim q_t} [w(t) || \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) + \sigma(t) \nabla \log q_t(\mathbf{x}_t) ||^2] \\ &= \mathbb{E}_{t \in [0, T], \mathbf{x}_0 \sim q_0, \boldsymbol{\epsilon}} [w(t) || \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) - \boldsymbol{\epsilon} ||^2]\end{aligned}$$

A deep neural network

where $\boldsymbol{\epsilon} \sim N(0, I)$ and $\mathbf{x}_t = \alpha(t)\mathbf{x}_0 + \sigma(t)\boldsymbol{\epsilon}$

So $\boldsymbol{\epsilon}_\theta(.,.)$ is also called the **noise prediction model**

- Intro to Diffusion Process
- Intro to Generative Diffusion Models
- Probability Flow ODE
- Latent Diffusion Models (Stable Diffusion)
- Consistency Models
- Latent Consistency Models

Prob-Flow ODE

Reverse process

Anderson's theorem

Reverse-time Brownian motion

$$d\mathbf{x}_t = [f(t)\mathbf{x}_t - g^2(t)\nabla_{\mathbf{x}} \log q_t(\mathbf{x}_t)] dt + g(t)d\overline{\mathbf{w}}_t, \quad \mathbf{x}_T \sim q_T(\mathbf{x}_T)$$

The above SDE defines the same marginal distribution $q_t(\mathbf{x}_t)$ as the following **Probability-Flow ODE** [Song et al. 2020]

$$\frac{d\mathbf{x}_t}{dt} = f(t)\mathbf{x}_t - \frac{1}{2}g^2(t)\nabla_{\mathbf{x}} \log q_t(\mathbf{x}_t), \quad \mathbf{x}_T \sim q_T(\mathbf{x}_T)$$

Since the noise prediction model $\epsilon_{\theta}(\cdot, \cdot)$ tries to fit the score function, the following **empirical PF-ODE** is an approximation of the above PF-ODE

$$\frac{d\mathbf{x}_t}{dt} = f(t)\mathbf{x}_t + \frac{g^2(t)}{2\sigma_t}\epsilon_{\theta}(\mathbf{x}_t, t), \quad \mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \tilde{\sigma}^2 \mathbf{I})$$

Probability Flow ODE

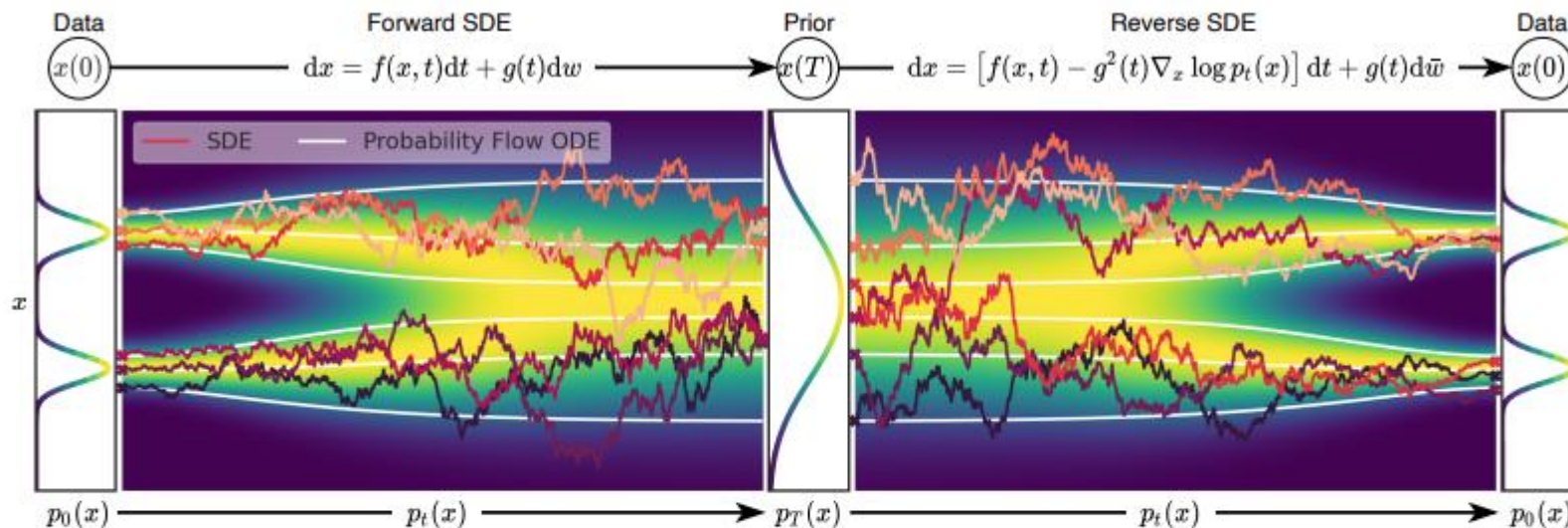


Figure 2: **Overview of score-based generative modeling through SDEs.** We can map data to a noise distribution (the prior) with an SDE (Section 3.1), and reverse this SDE for generative modeling (Section 3.2). We can also reverse the associated probability flow ODE (Section 4.3), which yields a deterministic process that samples from the same distribution as the SDE. Both the reverse-time SDE and probability flow ODE can be obtained by estimating the score $\nabla_x \log p_t(x)$ (Section 3.3).

$$\text{prob flow ODE: } dx = \left[f(x, t) - \frac{1}{2}g(t)^2 \nabla_x \log p_t(x) \right] dt$$

Sample from PF-ODE

- Empirical PF-ODE

$$\frac{d\mathbf{x}_t}{dt} = f(t)\mathbf{x}_t + \frac{g^2(t)}{2\sigma_t}\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t), \quad \mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \tilde{\sigma}^2 \mathbf{I})$$

- ODE trajectories are smoother than SDE trajectories
- To accelerate sampling, need to build better ODE solvers

These ODE solvers can
approximate the integration

$$\int_{t_{n+1}}^{t_n} \left(f(t)\mathbf{z}_t + \frac{g^2(t)}{2\sigma_t}\boldsymbol{\epsilon}_\theta(\mathbf{z}_t, \mathbf{c}, t) \right) dt \approx \Psi(\mathbf{z}_{t_{n+1}}, t_{n+1}, t_n, \mathbf{c})$$

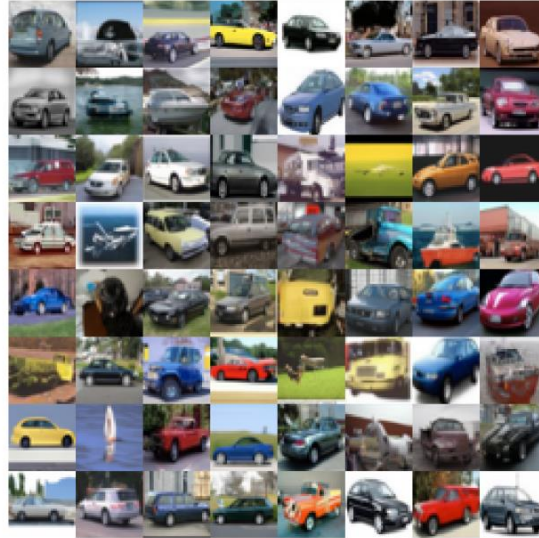
- Several diffusion models are based on the above idea: DDPM, DDIM, DPM-Solver, DPM-Solver++
- Current best ODE solvers (DPM-Solver, DPM-Solver++) requires 10-20 steps [Lu et al. 2022]

Generated Images

class: airplane



class: automobile



class: bird



Cifar 10



CelebA-HQ

Song et al. ICLR 2021

- Intro to Diffusion Process
- Intro to Generative Diffusion Models
- Probability Flow ODE
- Latent Diffusion Models (Stable Diffusion)
- Consistency Models
- Latent Consistency Models

Stable Diffusion

'A painting of the last supper by Picasso.'



High-Resolution Image Synthesis with Latent Diffusion Models

Stable diffusion

Idea: the reconstructions are confined to the **image manifold** (enforcing local realism), rather than on pixel-space (**Diffusion in latent space**)

Use an autoencoder in VQGAN (trained by combination of a perceptual loss and a patch-based adversarial objective)

Encoder E encodes image x into a latent representation $z = E(x)$

The encoder downsamples the image by a factor f

Decoder D reconstructs the image from the latent, giving $\tilde{x} = D(z) = D(E(x))$

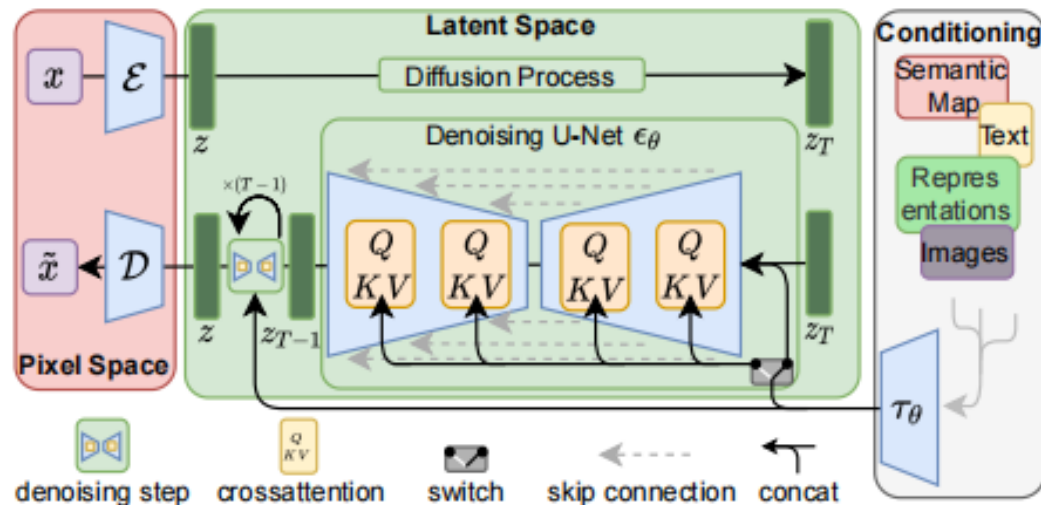


Figure 3. We condition LDMs either via concatenation or by a more general cross-attention mechanism. See Sec. 3.3

Stable diffusion

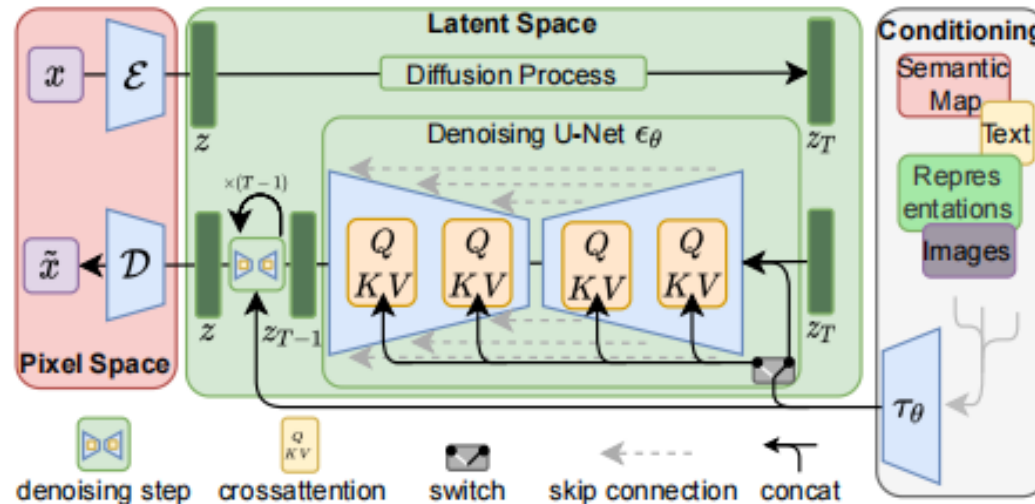


Figure 3. We condition LDMs either via concatenation or by a more general cross-attention mechanism. See Sec. 3.3

Loss from variational lower bound on $p(x)$, which mirrors denoising score-matching

$$\text{Loss of Diffusion model: } L_{DM} = \mathbb{E}_{x, \epsilon \sim \mathcal{N}(0,1), t} \left[\|\epsilon - \epsilon_\theta(x_t, t)\|_2^2 \right]$$

denoising autoencoders $\epsilon_\theta(x_t, t)$

Loss of Diffusion model in latent space: $L_{LDM} := \mathbb{E}_{\mathcal{E}(x), \epsilon \sim \mathcal{N}(0,1), t} \left[\|\epsilon - \epsilon_\theta(z_t, t)\|_2^2 \right]$.

The neural backbone of the model $\epsilon_\theta(z_t, t)$ is realized as a time-conditional UNet

Stable Diffusion: Text-to-Image

Text-to-Image Synthesis on LAION. 1.45B Model.

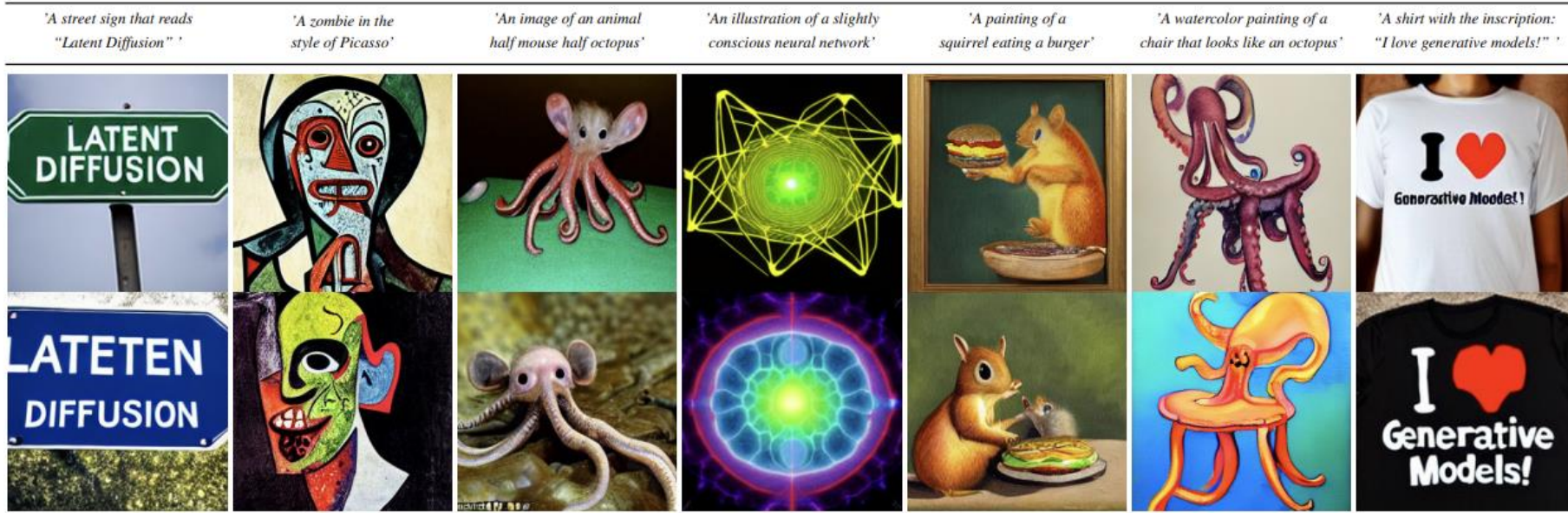


Figure 5. Samples for user-defined text prompts from our model for text-to-image synthesis, *LDM-8 (KL)*, which was trained on the LAION [78] database. Samples generated with 200 DDIM steps and $\eta = 1.0$. We use unconditional guidance [32] with $s = 10.0$.

- Intro to Diffusion Process
- Intro to Generative Diffusion Models
- Probability Flow ODE
- Latent Diffusion Models (Stable Diffusion)
- Consistency Models
- Latent Consistency Models

Consistency Models

Goal: Accelerate inference

How: Learn a **consistency function** $f : (\mathbf{x}_t, t) \mapsto \mathbf{x}_\epsilon$

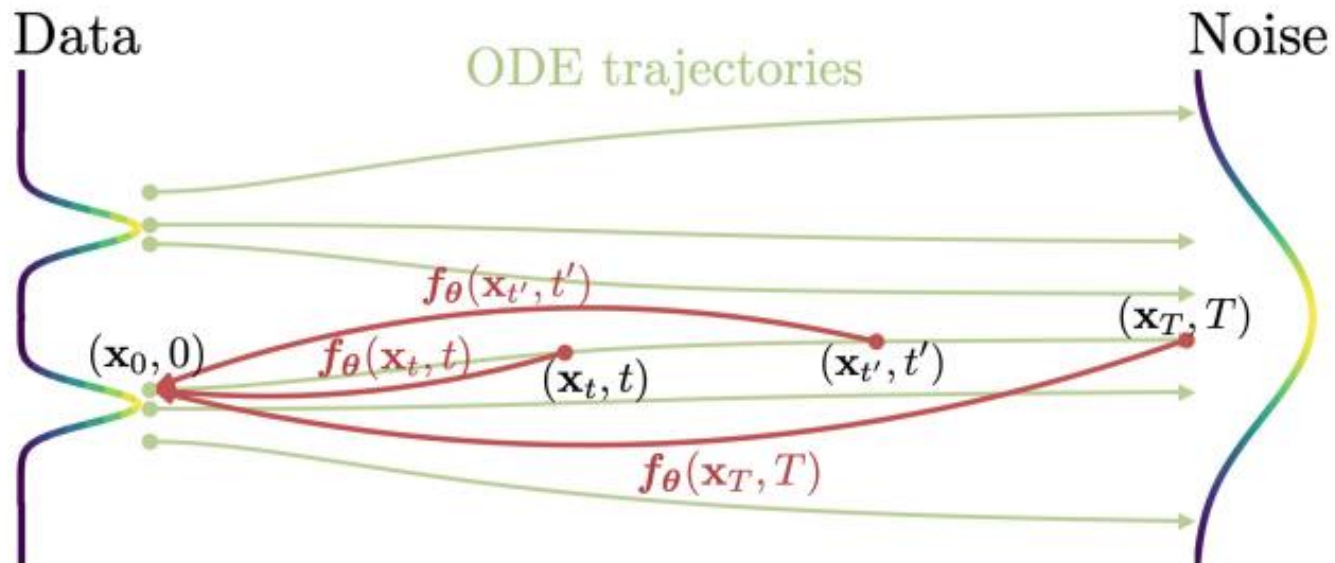


Figure 2: **Consistency models** are trained to map points on any trajectory of the **PF ODE** to the trajectory's origin.

Consistency Models

Consistency Property:

$$\mathbf{f}(\mathbf{x}_t, t) = \mathbf{f}(\mathbf{x}_{t'}, t') \text{ for all } t, t' \in [\epsilon, T]$$

How to ensure such property?

Define the following consistency loss:

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^-; \Phi) = \mathbb{E}_{\mathbf{x}, t} \left[d \left(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}_{t_{n+1}}, t_{n+1}), \mathbf{f}_{\boldsymbol{\theta}^-}(\hat{\mathbf{x}}_{t_n}^{\phi}, t_n) \right) \right]$$

Here: $\boldsymbol{\theta}^- \leftarrow \mu \bar{\boldsymbol{\theta}}^- + (1 - \mu) \boldsymbol{\theta}$

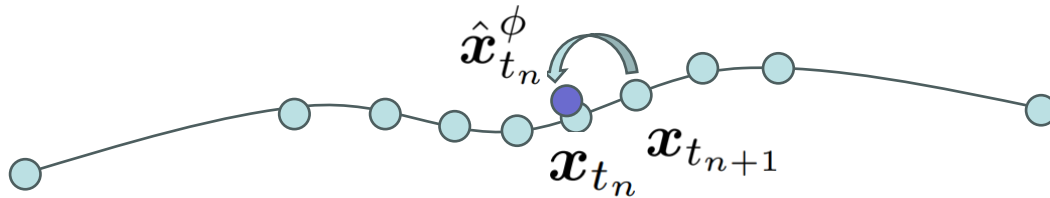
$$\hat{\mathbf{x}}_{t_n}^{\phi} \leftarrow \mathbf{x}_{t_{n+1}} + (t_n - t_{n+1}) \Phi(\mathbf{x}_{t_{n+1}}, t_{n+1}; \phi)$$

Explain in next page

Consistency Models

$$\hat{\mathbf{x}}_{t_n}^\phi \leftarrow \mathbf{x}_{t_{n+1}} + (t_n - t_{n+1})\Phi(\mathbf{x}_{t_{n+1}}, t_{n+1}; \phi)$$

- How do we estimate \mathbf{x}_{t_n} from $\mathbf{x}_{t_{n+1}}$



Previous ODE solvers requires many steps to generate high quality images

Consistency Models

We can distill a consistency model from an existing diffusion models

Algorithm 2 Consistency Distillation (CD)

Input: dataset \mathcal{D} , initial model parameter θ , learning rate η , ODE solver $\Phi(\cdot, \cdot; \phi)$, $d(\cdot, \cdot)$, $\lambda(\cdot)$, and μ

$\theta^- \leftarrow \theta$

repeat

 Sample $\mathbf{x} \sim \mathcal{D}$ and $n \sim \mathcal{U}[[1, N - 1]]$

 Sample $\mathbf{x}_{t_{n+1}} \sim \mathcal{N}(\mathbf{x}; t_{n+1}^2 \mathbf{I})$

$\hat{\mathbf{x}}_{t_n}^\phi \leftarrow \mathbf{x}_{t_{n+1}} + (t_n - t_{n+1})\Phi(\mathbf{x}_{t_{n+1}}, t_{n+1}; \phi)$

$\mathcal{L}(\theta, \theta^-; \phi) \leftarrow$

$\lambda(t_n)d(\mathbf{f}_\theta(\mathbf{x}_{t_{n+1}}, t_{n+1}), \mathbf{f}_{\theta^-}(\hat{\mathbf{x}}_{t_n}^\phi, t_n))$

$\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}(\theta, \theta^-; \phi)$

$\theta^- \leftarrow \text{stopgrad}(\mu \theta^- + (1 - \mu)\theta)$

until convergence

Image generated by CM

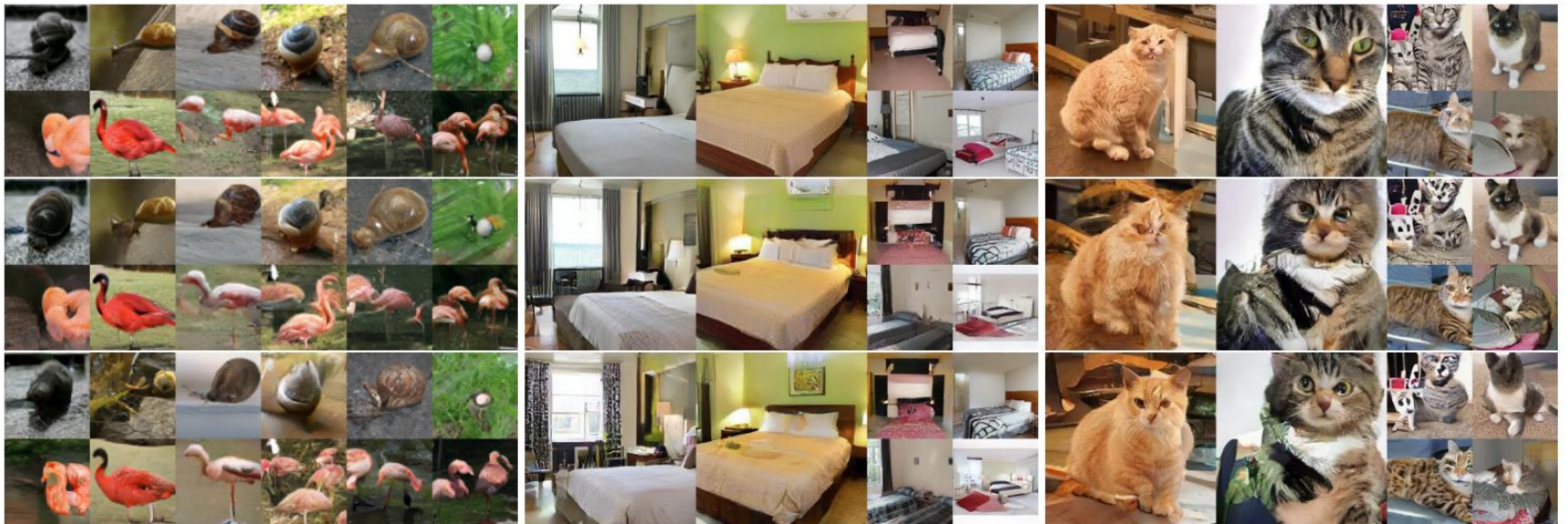


Figure 5: Samples generated by EDM (*top*), CT + single-step generation (*middle*), and CT + 2-step generation (*Bottom*). All corresponding images are generated from the same initial noise.

- Intro to Diffusion Process
- Intro to Generative Diffusion Models
- Probability Flow ODE
- Latent Diffusion Models (Stable Diffusion)
- Consistency Models
- Latent Consistency Models

Latent Consistency Models

- Goal:
 - Synthesizing high quality images, like Stable Diffusion
 - Conditional generation: text-to-image
 - Incorporate classifier free guidance
 - Very fast inference: 2-4 steps, even 1 step
 - Allow easy finetuning on downstream customized dataset

Incorporate Classifier-Free Guidance

- Classifier-Free Guidance (CFG) [Ho & Salimans, 2022] is an important technique for generating high quality images
- Given a CFG scale ω , the original noise prediction is replaced by a linear combination of conditional and unconditional noise prediction

$$\tilde{\epsilon}_{\theta}(\mathbf{z}_t, \omega, \mathbf{c}, t) = (1 + \omega)\epsilon_{\theta}(\mathbf{z}_t, \mathbf{c}, t) - \omega\epsilon_{\theta}(\mathbf{z}, \emptyset, t)$$

Latent Consistency Distillation

Algorithm 1 Latent Consistency Distillation (LCD)

Input: dataset \mathcal{D} , initial model parameter θ , learning rate η , ODE solver $\Psi(\cdot, \cdot, \cdot, \cdot)$, distance metric $d(\cdot, \cdot)$, EMA rate μ , noise schedule $\alpha(t), \sigma(t)$, guidance scale $[w_{\min}, w_{\max}]$, skipping interval k , and encoder $E(\cdot)$

Encoding training data into latent space: $\mathcal{D}_z = \{(z, c) | z = E(x), (x, c) \in \mathcal{D}\}$

$\theta^- \leftarrow \theta$

repeat

Sample $(z, c) \sim \mathcal{D}_z, n \sim \mathcal{U}[1, N - k]$ and $\omega \sim [\omega_{\min}, \omega_{\max}]$

Sample $z_{t_{n+k}} \sim \mathcal{N}(\alpha(t_{n+k})z; \sigma^2(t_{n+k})\mathbf{I})$

$\hat{z}_{t_n}^{\Psi, \omega} \leftarrow z_{t_{n+k}} + (1 + \omega)\Psi(z_{t_{n+k}}, t_{n+k}, t_n, c) - \omega\Psi(z_{t_{n+k}}, t_{n+k}, t_n, \emptyset)$

$\mathcal{L}(\theta, \theta^-; \Psi) \leftarrow d(f_\theta(z_{t_{n+k}}, \omega, c, t_{n+k}), f_\theta(\hat{z}_{t_n}^{\Psi, \omega}, \omega, c, t_n))$

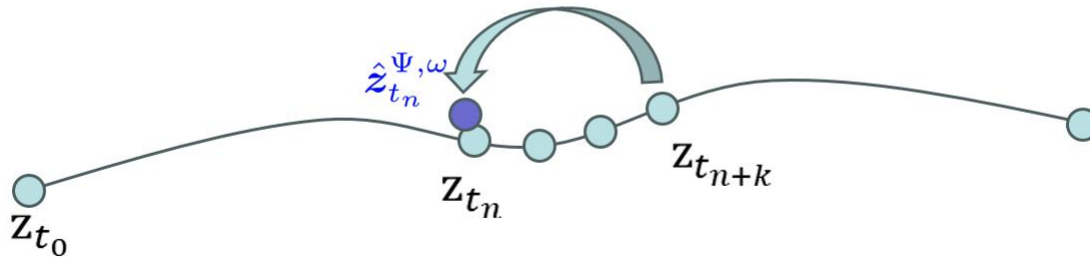
$\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}(\theta, \theta^-)$

$\theta^- \leftarrow \text{stopgrad}(\mu\theta^- + (1 - \mu)\theta)$

until convergence

CFG and more
advanced ODE-solver

Use existing ODE solver Ψ to approximate
 z_{t_n} (e.g. DDIM, DPM-solver, DPM-solver++)



$$\Psi_{\text{DDIM}}(z_{t_{n+k}}, t_{n+k}, t_n, c) = \underbrace{\frac{\alpha_{t_n}}{\alpha_{t_{n+k}}} z_{t_{n+k}} - \sigma_{t_n} \left(\frac{\sigma_{t_{n+k}} \cdot \alpha_{t_n}}{\alpha_{t_{n+k}} \cdot \sigma_{t_n}} - 1 \right) \hat{\epsilon}_\theta(z_{t_{n+k}}, c, t_{n+k})}_{\text{DDIM Estimated } z_{t_n}} - z_{t_{n+k}}$$

Inference/Sampling

Algorithm 3 Multistep Latent Consistency Sampling

Input: Latent Consistency Model $f_{\theta}(\cdot, \cdot, \cdot, \cdot)$, Sequence of timesteps $\tau_1 > \tau_2 > \dots > \tau_{N-1}$, Text condition c , Classifier-Free Guidance Scale ω , Noise schedule $\alpha(t), \sigma(t)$, Decoder $D(\cdot)$

Sample initial noise $\hat{z}_T \sim \mathcal{N}(\mathbf{0}; \mathbf{I})$

$z \leftarrow f_{\theta}(\hat{z}_T, \omega, c, T)$

for $n = 1$ to $N - 1$ **do**

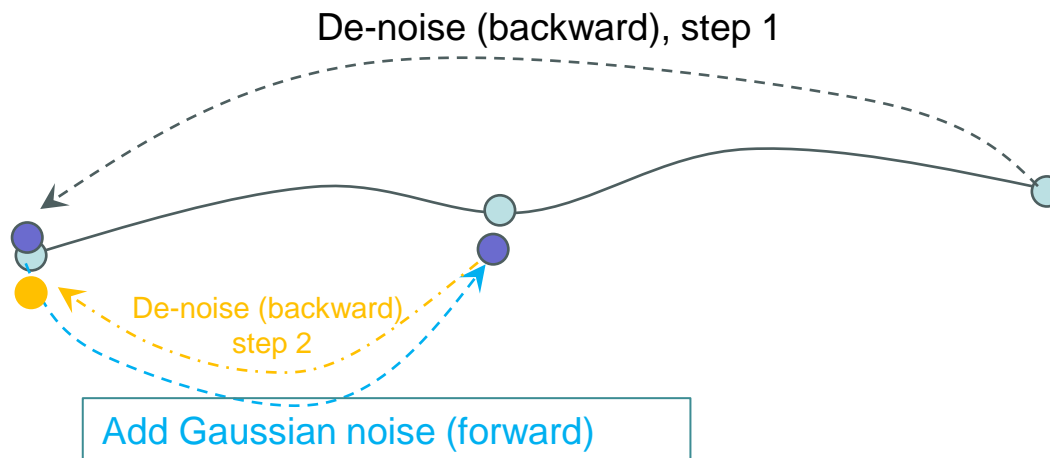
$\hat{z}_{\tau_n} \sim \mathcal{N}(\alpha(\tau_n)z; \sigma^2(\tau_n)\mathbf{I})$ Add noise (forward)

$z \leftarrow f_{\theta}(\hat{z}_{\tau_n}, \omega, c, \tau_n)$

end for

$x \leftarrow D(z)$

Output: x



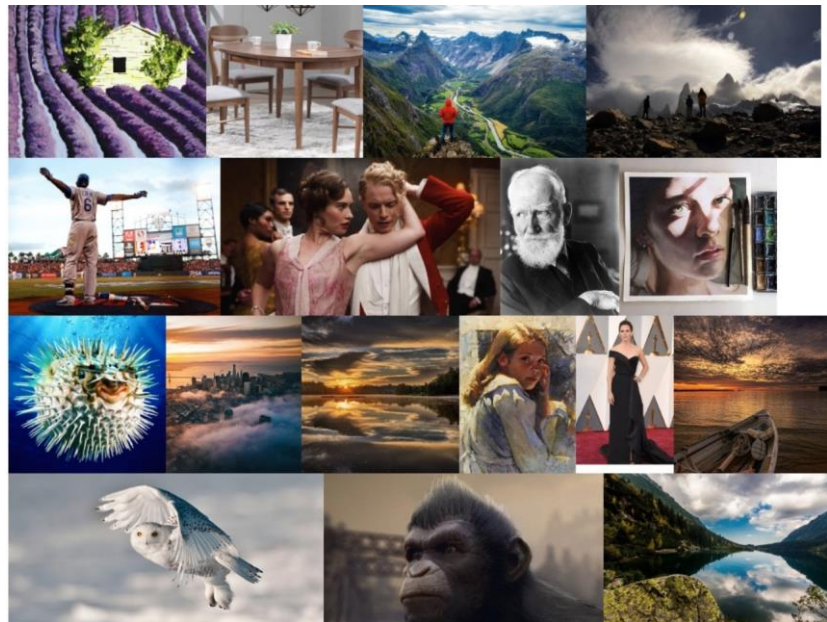
Two-step sampling

Experiments

- LAION-5B dataset: LAION-Aesthetics-6+ (12M) and LAION-Aesthetics-6.5+ (650K)
- Training: high-quality 768×768 2~4-step LCM takes only 32 A100 GPU hours for training



An orange cat is looking at its reflection in the mirror.



Experiments

| MODEL (512 × 512) RESO | FID ↓ | | | | CLIP SCORE ↑ | | | |
|------------------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | 1 STEP | 2 STEPS | 4 STEPS | 8 STEPS | 1 STEP | 2 STEPS | 4 STEPS | 8 STEPS |
| DDIM (Song et al., 2020a) | 183.29 | 81.05 | 22.38 | 13.83 | 6.03 | 14.13 | 25.89 | 29.29 |
| DPM (Lu et al., 2022a) | 185.78 | 72.81 | 18.53 | 12.24 | 6.35 | 15.10 | 26.64 | 29.54 |
| DPM++ (Lu et al., 2022b) | 185.78 | 72.81 | 18.43 | 12.20 | 6.35 | 15.10 | 26.64 | 29.55 |
| Guided-Distill (Meng et al., 2023) | 108.21 | 33.25 | 15.12 | 13.89 | 12.08 | 22.71 | 27.25 | 28.17 |
| LCM (Ours) | 35.36 | 13.31 | 11.10 | 11.84 | 24.14 | 27.83 | 28.69 | 28.84 |

Table 1: Quantitative results with $\omega = 8$ at 512×512 resolution. LCM significantly surpasses baselines in the 1-4 step region on LAION-Aesthetic-6+ dataset. For LCM, DDIM-Solver is used with a skipping step of $k = 20$.

| MODEL (768 × 768) RESO | FID ↓ | | | | CLIP SCORE ↑ | | | |
|------------------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | 1 STEP | 2 STEPS | 4 STEPS | 8 STEPS | 1 STEP | 2 STEPS | 4 STEPS | 8 STEPS |
| DDIM (Song et al., 2020a) | 186.83 | 77.26 | 24.28 | 15.66 | 6.93 | 16.32 | 26.48 | 29.49 |
| DPM (Lu et al., 2022a) | 188.92 | 67.14 | 20.11 | 14.08 | 7.40 | 17.11 | 27.25 | 29.80 |
| DPM++ (Lu et al., 2022b) | 188.91 | 67.14 | 20.08 | 14.11 | 7.41 | 17.11 | 27.26 | 29.84 |
| Guided-Distill (Meng et al., 2023) | 120.28 | 30.70 | 16.70 | 14.12 | 12.88 | 24.88 | 28.45 | 29.16 |
| LCM (Ours) | 34.22 | 16.32 | 13.53 | 14.97 | 25.32 | 27.92 | 28.60 | 28.49 |

Table 2: Quantitative results with $\omega = 8$ at 768×768 resolution. LCM significantly surpasses the baselines in the 1-4 step region on LAION-Aesthetic-6.5+ dataset. For LCM, DDIM-Solver is used with a skipping step of $k = 20$.

High resolution images generate by LCM



4 steps

High resolution images generate by LCM



4 steps

High resolution images generate by LCM



4 steps

High resolution images generate by LCM



2 steps

High resolution images generate by LCM



2 steps

High resolution images generate by LCM

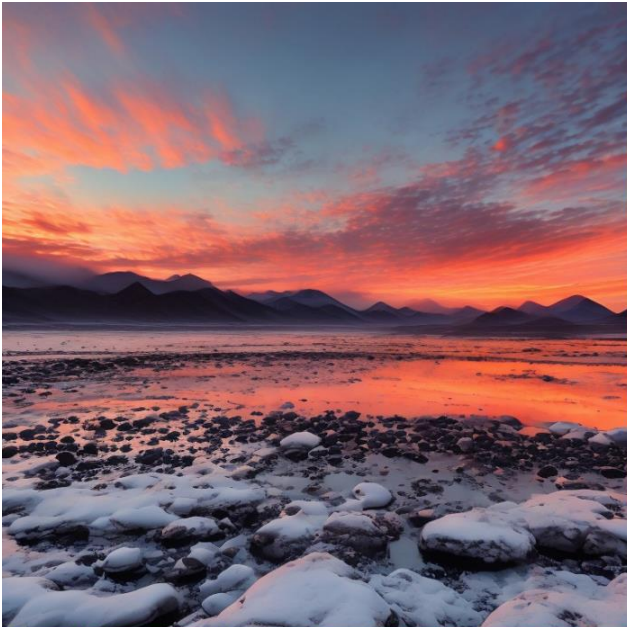


1 step

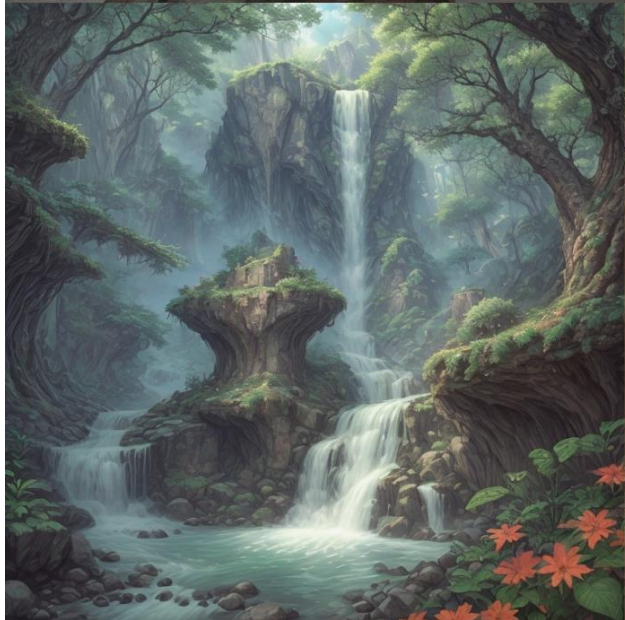
High resolution images generate by LCM



1 step



2 steps



2 steps

Finetuning on customized dataset

Algorithm 4 Latent Consistency Fine-tuning (LCF)

Input: customized dataset $\mathcal{D}^{(s)}$, pre-trained LCM parameter θ , learning rate η , distance metric $d(\cdot, \cdot)$, EMA rate μ , noise schedule $\alpha(t), \sigma(t)$, guidance scale $[w_{\min}, w_{\max}]$, skipping interval k , and encoder $E(\cdot)$

Encode training data into the latent space: $\mathcal{D}_z^{(s)} = \{(z, c) | z = E(x), (x, c) \in \mathcal{D}^{(s)}\}$

$\theta^- \leftarrow \theta$

repeat

 Sample $(z, c) \sim \mathcal{D}_z^{(s)}$, $n \sim \mathcal{U}[1, N - k]$ and $w \sim [w_{\min}, w_{\max}]$

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

$z_{t_{n+k}} \leftarrow \alpha(t_{n+k})z + \sigma(t_{n+k})\epsilon$, $z_{t_n} \leftarrow \alpha(t_n)z + \sigma(t_n)\epsilon$

$\mathcal{L}(\theta, \theta^-) \leftarrow d(\mathbf{f}_\theta(z_{t_{n+k}}, t_{n+k}, c, w), \mathbf{f}_{\theta^-}(z_{t_n}, t_n, c, w))$

$\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}(\theta, \theta^-)$

$\theta^- \leftarrow \text{stopgrad}(\mu\theta^- + (1 - \mu)\theta)$

until convergence

Finetuning on customized dataset

Origin LCM



1K Finetuning



10K Finetuning



30K Finetuning



Pokemon dataset

Thanks



Jian Li 李建

lapordge@gmail.com

Wechat id: lapordge

Latent Consistency Distillation

Algorithm 1 Latent Consistency Distillation (LCD)

Input: dataset \mathcal{D} , initial model parameter θ , learning rate η , ODE solver $\Psi(\cdot, \cdot, \cdot, \cdot)$, distance metric $d(\cdot, \cdot)$, EMA rate μ , noise schedule $\alpha(t), \sigma(t)$, guidance scale $[w_{\min}, w_{\max}]$, skipping interval k , and encoder $E(\cdot)$
 Encoding training data into latent space: $\mathcal{D}_z = \{(z, c) | z = E(x), (x, c) \in \mathcal{D}\}$

$\theta^- \leftarrow \theta$

repeat

 Sample $(z, c) \sim \mathcal{D}_z, n \sim \mathcal{U}[1, N - k]$ and $\omega \sim [\omega_{\min}, \omega_{\max}]$

 Sample $z_{t_{n+k}} \sim \mathcal{N}(\alpha(t_{n+k})z; \sigma^2(t_{n+k})\mathbf{I})$

$\hat{z}_{t_n}^{\Psi, \omega} \leftarrow z_{t_{n+k}} + (1 + \omega)\Psi(z_{t_{n+k}}, t_{n+k}, t_n, c) - \omega\Psi(z_{t_{n+k}}, t_{n+k}, t_n, \emptyset)$

$\mathcal{L}(\theta, \theta^-; \Psi) \leftarrow d(f_\theta(z_{t_{n+k}}, \omega, c, t_{n+k}), f_{\theta^-}(\hat{z}_{t_n}^{\Psi, \omega}, \omega, c, t_n))$

$\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}(\theta, \theta^-)$

$\theta^- \leftarrow \text{stopgrad}(\mu\theta^- + (1 - \mu)\theta)$

until convergence

Use existing ODE solver Ψ to approximate z_{t_n} (e.g. DDIM, DPM-solver, DPM-solver++)

