
Consistency Models

Yang Song¹ Prafulla Dhariwal¹ Mark Chen¹ Ilya Sutskever¹

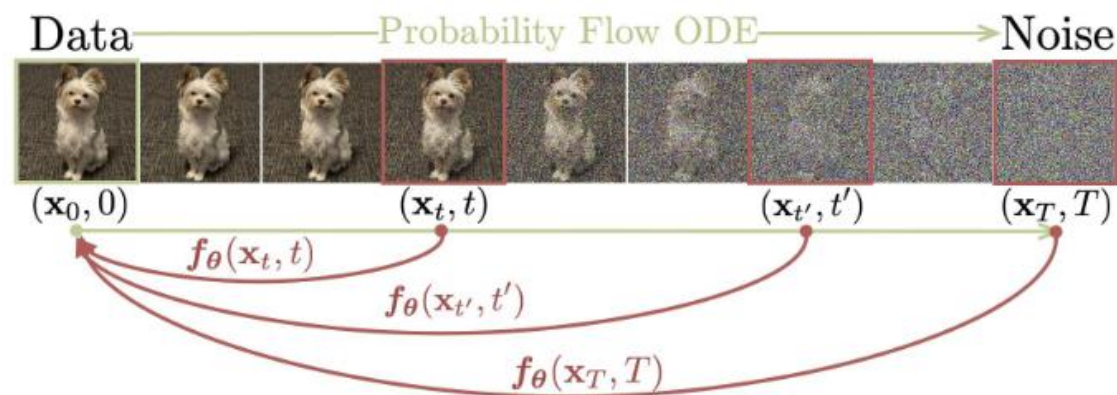
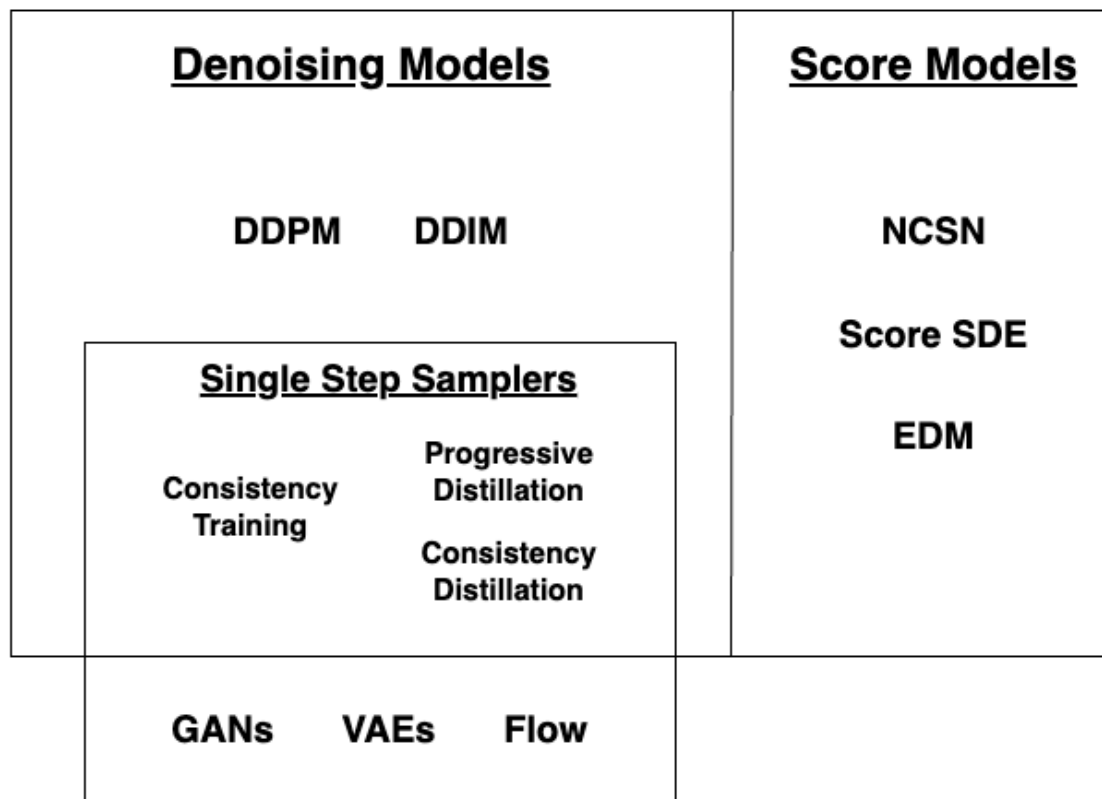


Table of Contents

- SDE/PF-ODE, Denoising models vs. Score models
 - Problem Definition
 - Consistency Distillation
 - Consistency Training
 - Metrics
 - Applications
 - Follow-up researches
-

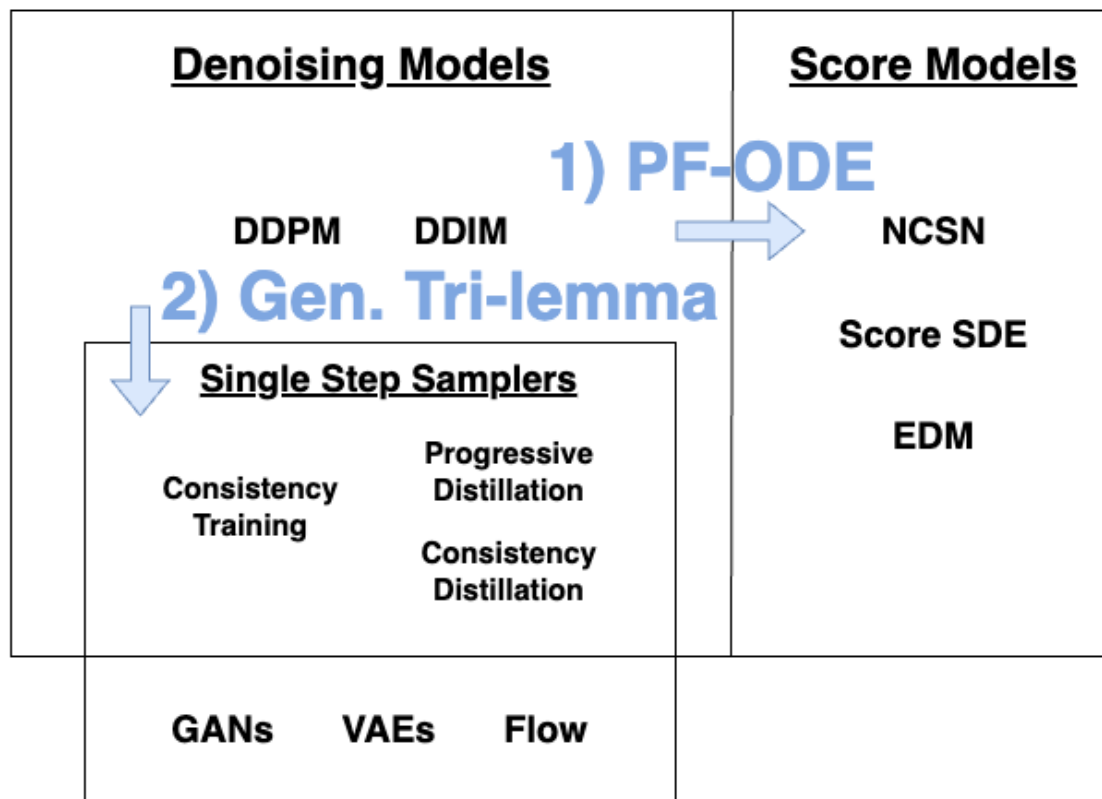
Introduction

Score Based Generative Models (Diffusion)



Introduction

Score Based Generative Models (Diffusion)



Diffusion Models

Task

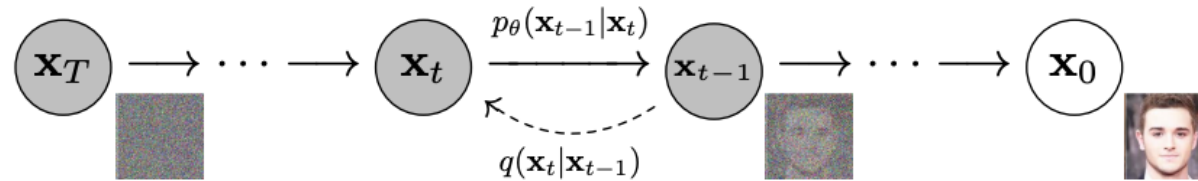


Figure 2: The directed graphical model considered in this work.

Model

$$f(x_{t_{n+1}}, t_{n+1}) \rightarrow x_{t_n}, x_T \rightarrow x_0$$

Parameterization

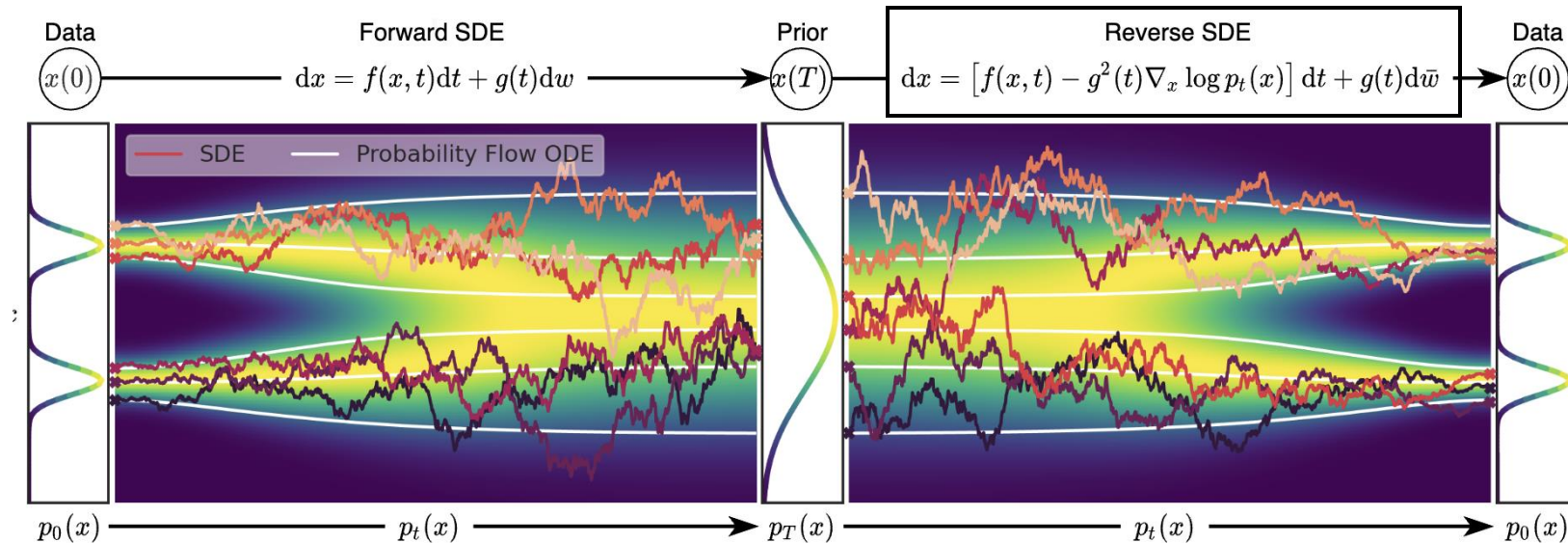
$$\text{noise} - \text{predict}, \quad \text{data} - \text{predict}$$

Score Function and ODE

$$d\mathbf{x} = \mathbf{f}(\mathbf{x}, t)dt + g(t)d\mathbf{w},$$

Denoising Diffusion

SDE

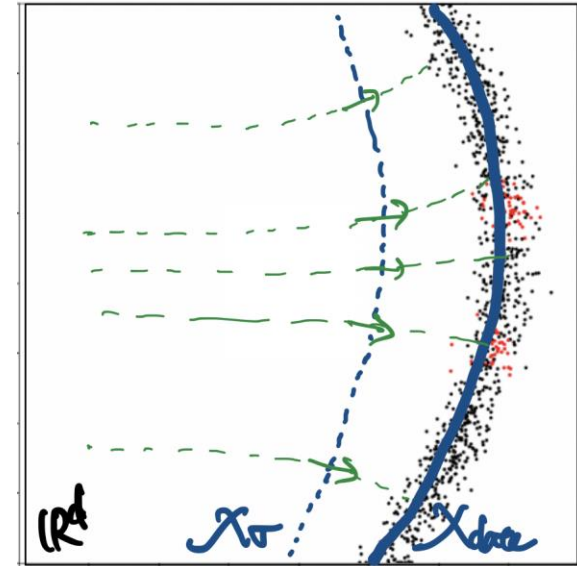
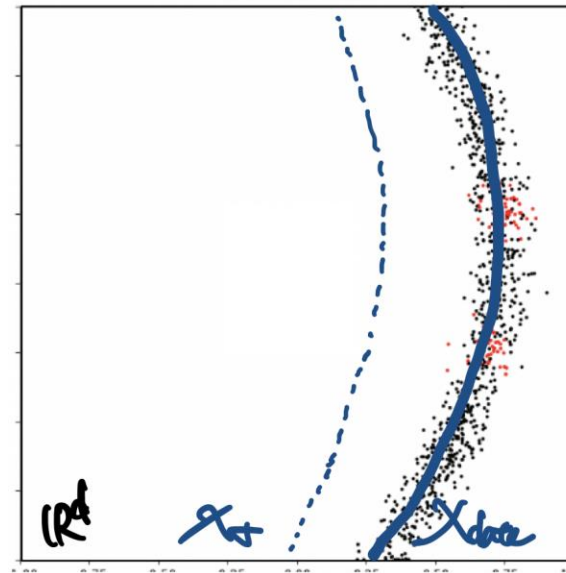
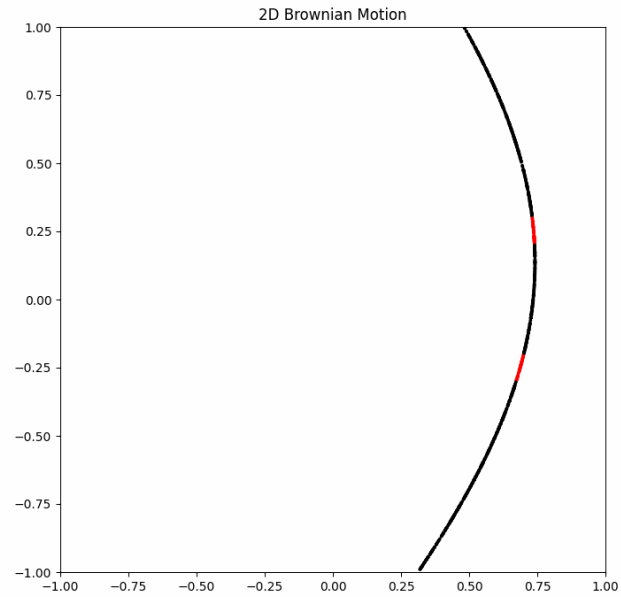


$$d\mathbf{x} = \left[\mathbf{f}(\mathbf{x}, t) - \frac{1}{2}g(t)^2\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) \right]dt,$$

PF ODE

→ Score function

SDE vs. PF-ODE



Score Models

$$d\mathbf{x} = \left[\mathbf{f}(\mathbf{x}, t) - \frac{1}{2}g(t)^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) \right] dt,$$

$\rightarrow s_{\theta}(\mathbf{x}, t)$

$$\frac{d\mathbf{x}}{d\sigma} = -\sigma \nabla_{\mathbf{x}} \log p_{\sigma}(\mathbf{x}) \quad \sigma \in [\sigma_{\min}, \sigma_{\max}],$$

$\rightarrow s_{\theta}(\mathbf{x}, \sigma)$

Model

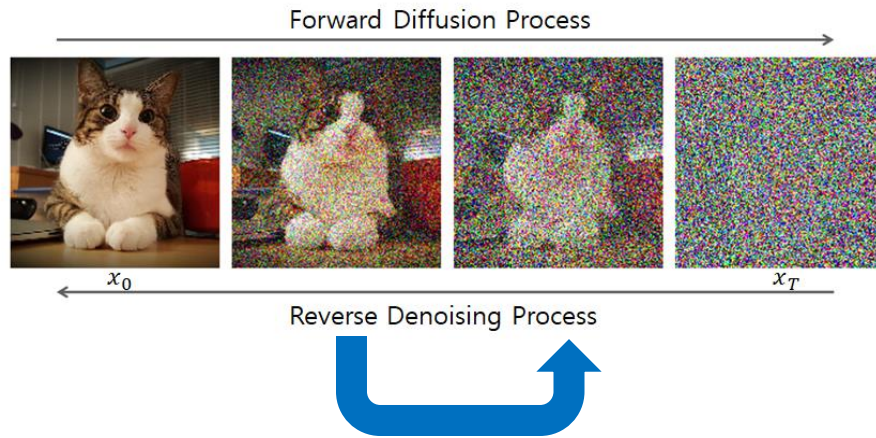
$$f(x_{t_{n+1}}, t_{n+1}) \rightarrow x_{t_n}, x_T \rightarrow x_0$$

Parameterization

score function – predict

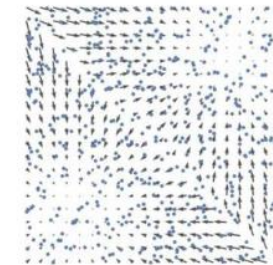
Denoising Process (Inference)

$$\frac{d\mathbf{x}}{d\sigma} = -\sigma \nabla_{\mathbf{x}} \log p_{\sigma}(\mathbf{x}) \quad \sigma \in [\sigma_{\min}, \sigma_{\max}],$$



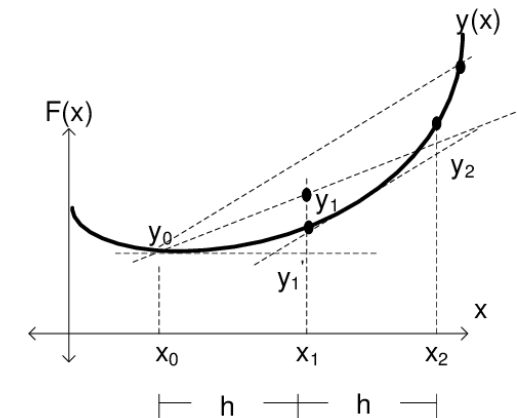
Direct Mapping

→ Empirical ODE solver



Follow noisy scores:
Langevin dynamics

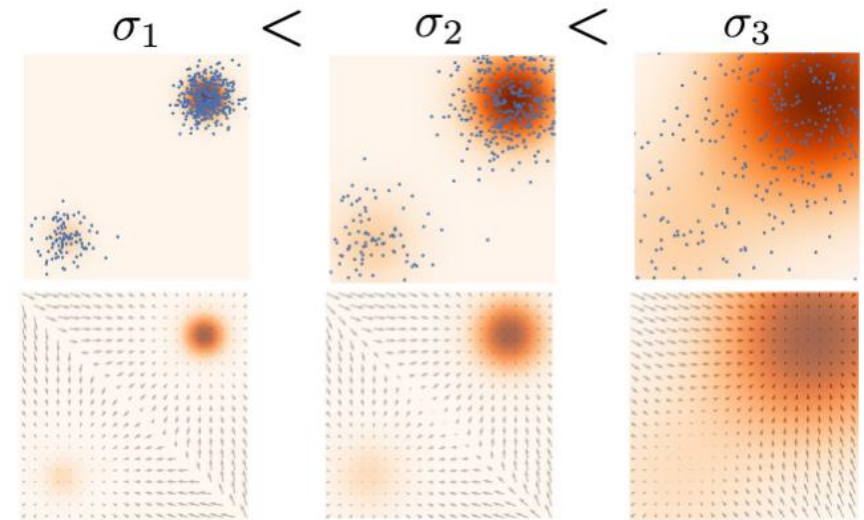
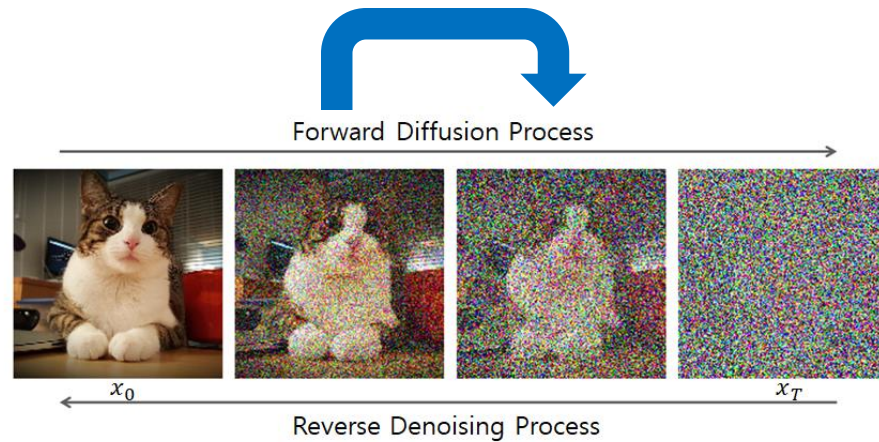
$$\mathbf{z}_t \sim \mathcal{N}(0, I)$$
$$\bar{\mathbf{x}}_{t+1} \leftarrow \bar{\mathbf{x}}_t + \frac{\epsilon}{2} s_{\theta}(\bar{\mathbf{x}}_t) + \sqrt{\epsilon} \mathbf{z}_t$$



Langevin Dynamics
Euler / Heun Estimation

→ Numerical ODE solver

Noising Process

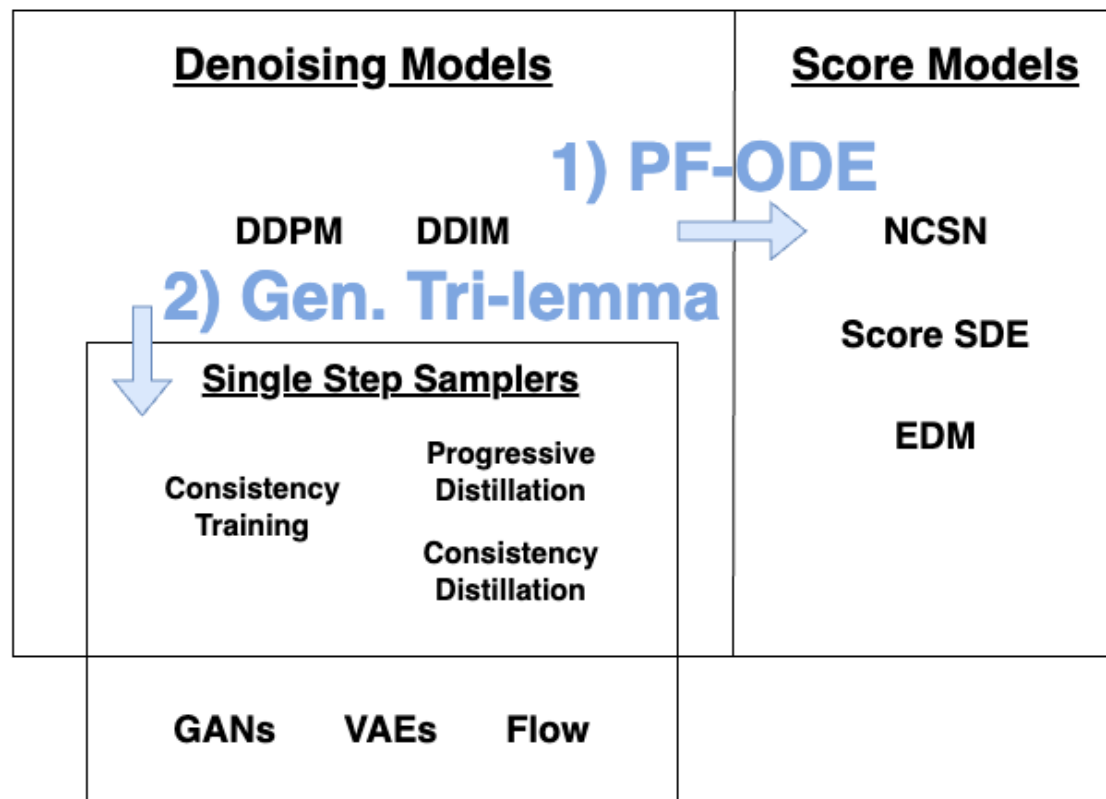


Data pairs for denoising training

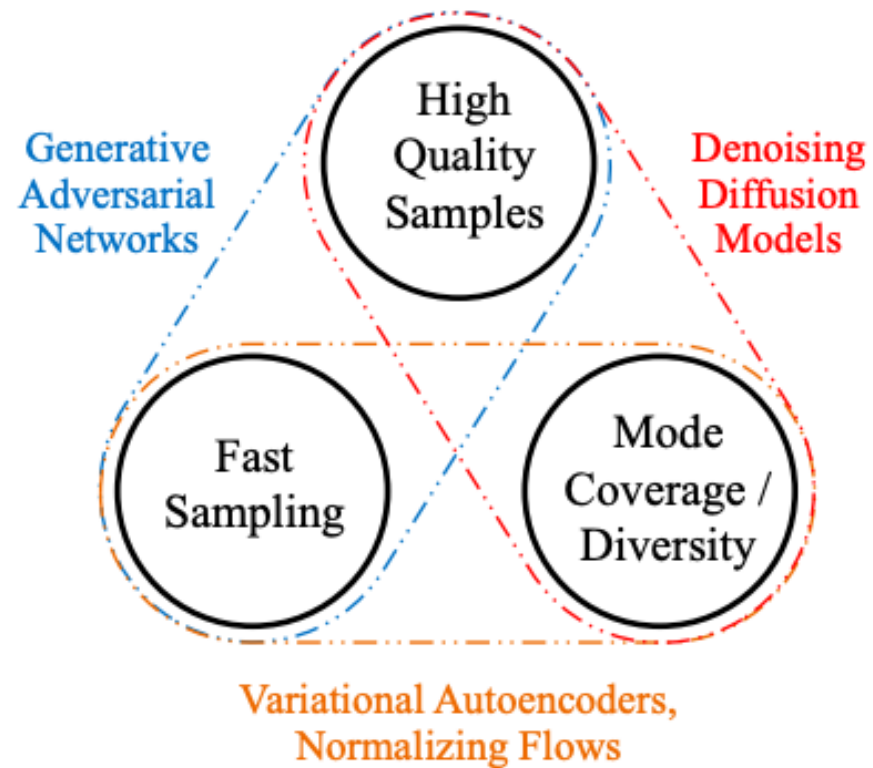
Data point sampling

$$X \sim f(x), x \in S, S \subseteq \mathbb{R}^d$$

Score Based Generative Models (Diffusion)

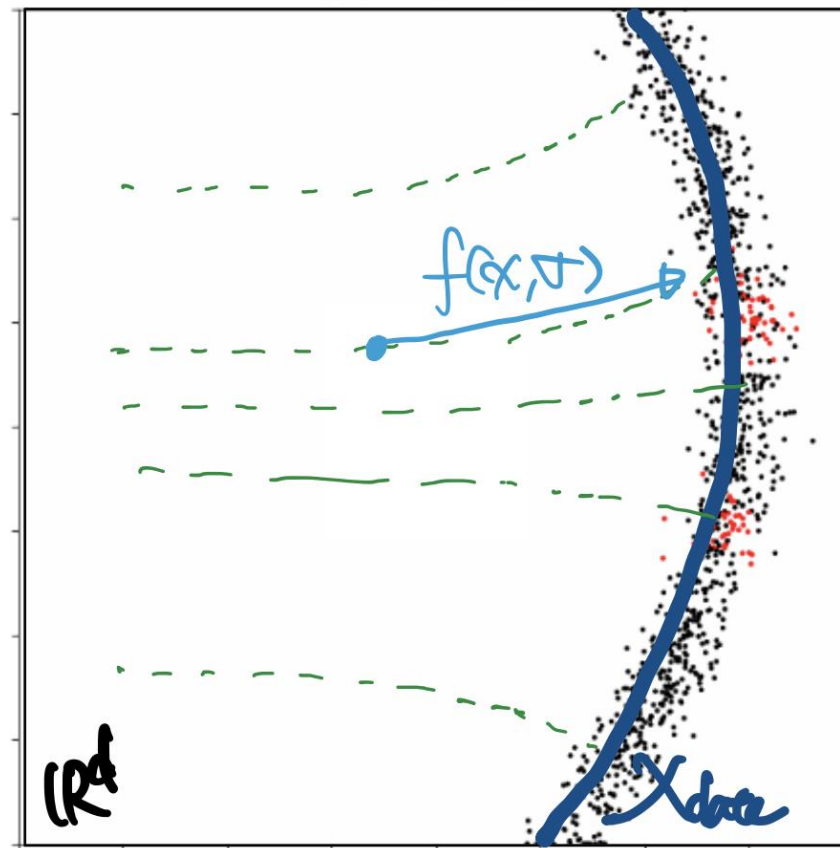


Generative Learning Trilemma



Problem Definition

- Sustain Quality of diffusion models
 - PF-ODE solver
- Single-step sampling
 - Direct mapping, $z \rightarrow x$ ($\mathbb{R}^d \rightarrow S, S \subseteq \mathbb{R}^d$)



Consistency Function

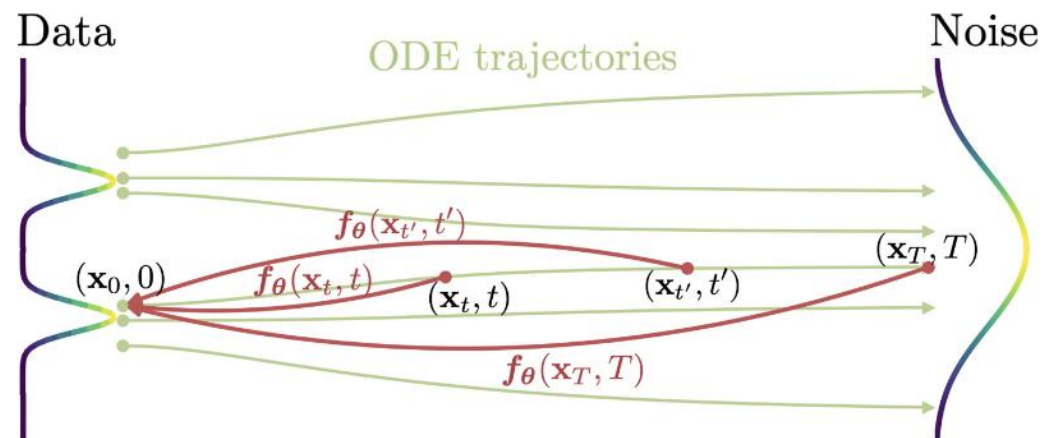


Figure 2: **Consistency models** are trained to map points on any trajectory of the **PF ODE** to the trajectory's origin.

Consistency function,

$$f : (x_t, t) \rightarrow x_\epsilon$$

must be self-consistent,

for any solution trajectory $\{x_t\}_{t \in [\epsilon, T]}$ on true PF ODE,

$$f(x_t, t) = f(x_{t'}, t'), \forall t, t' \in [\epsilon, T]$$

→ Approximation models can be trained via

- 1) CD (Consistency Distillation),
- 2) CT (Consistency Training)

Consistency Distillation

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

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

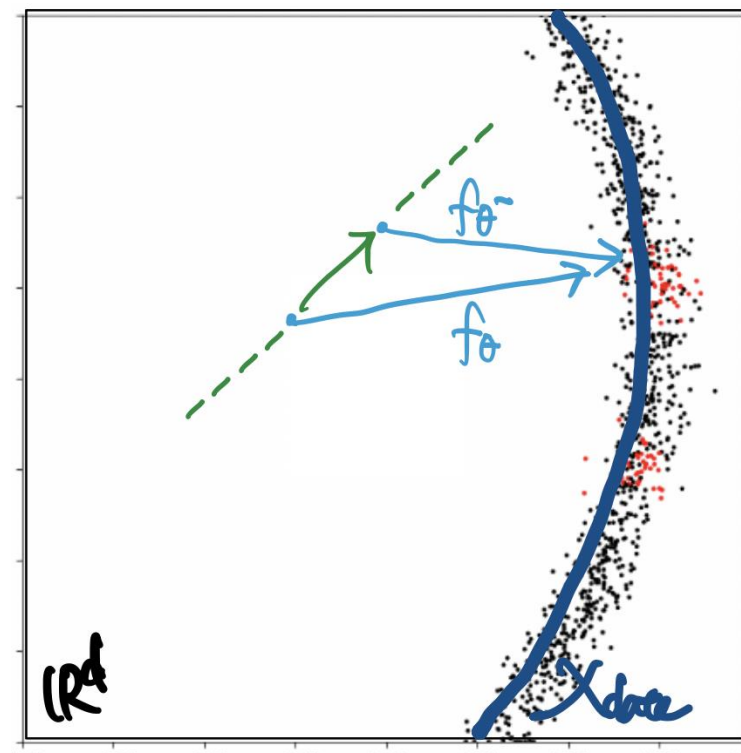
$\mathcal{L}(\theta, \theta^-; \phi) \leftarrow$

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

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

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

until convergence



Analogy to Contrastive Learning

- *Distillation* from Φ^φ
- *Data pair*: $(\widehat{x}_{t_n}^\varphi, t_n) - (x_{t_{n+1}}, t_{n+1})$
- Student: $\theta, x_{t_{n+1}}$, Teacher: θ^-, x_{t_n}
- Boundary condition $f(x_\epsilon, \epsilon) = x_\epsilon$

$$\mathcal{L}_{CD}^N(\theta, \theta^-; \phi) := \mathbb{E}[\lambda(t_n) d(\mathbf{f}_\theta(\mathbf{x}_{t_{n+1}}, t_{n+1}), \mathbf{f}_{\theta^-}(\hat{\mathbf{x}}_{t_n}^\phi, t_n))], \quad (7)$$

$$\mathbf{f}_\theta(\mathbf{x}, t) = c_{\text{skip}}(t)\mathbf{x} + c_{\text{out}}(t)F_\theta(\mathbf{x}, t), \quad (5)$$

$$c_{\text{skip}}(t) = \frac{\sigma_{\text{data}}^2}{(t - \epsilon)^2 + \sigma_{\text{data}}^2}, \quad c_{\text{out}}(t) = \frac{\sigma_{\text{data}}(t - \epsilon)}{\sqrt{\sigma_{\text{data}}^2 + t^2}},$$

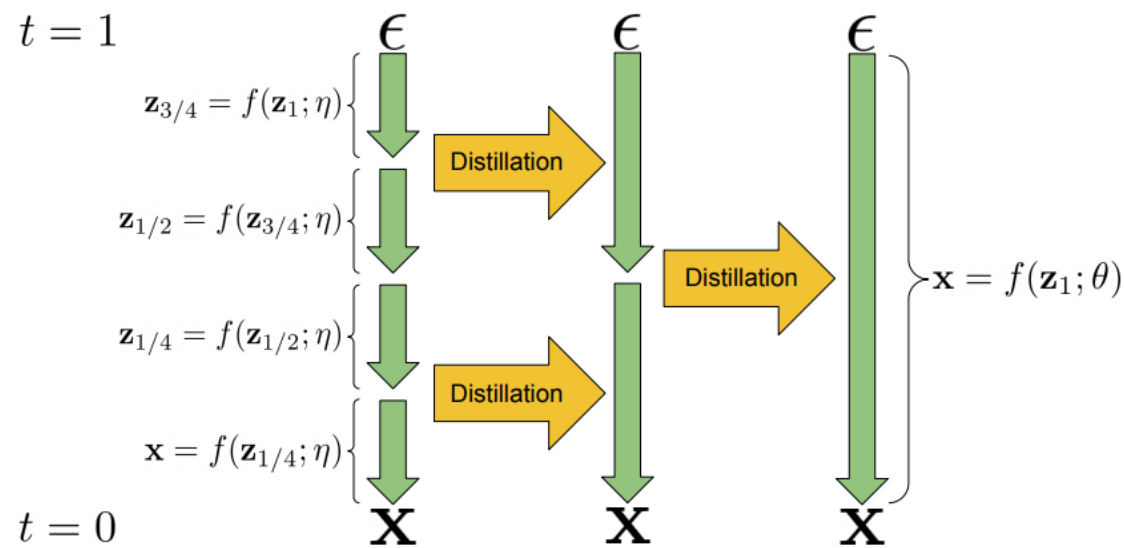
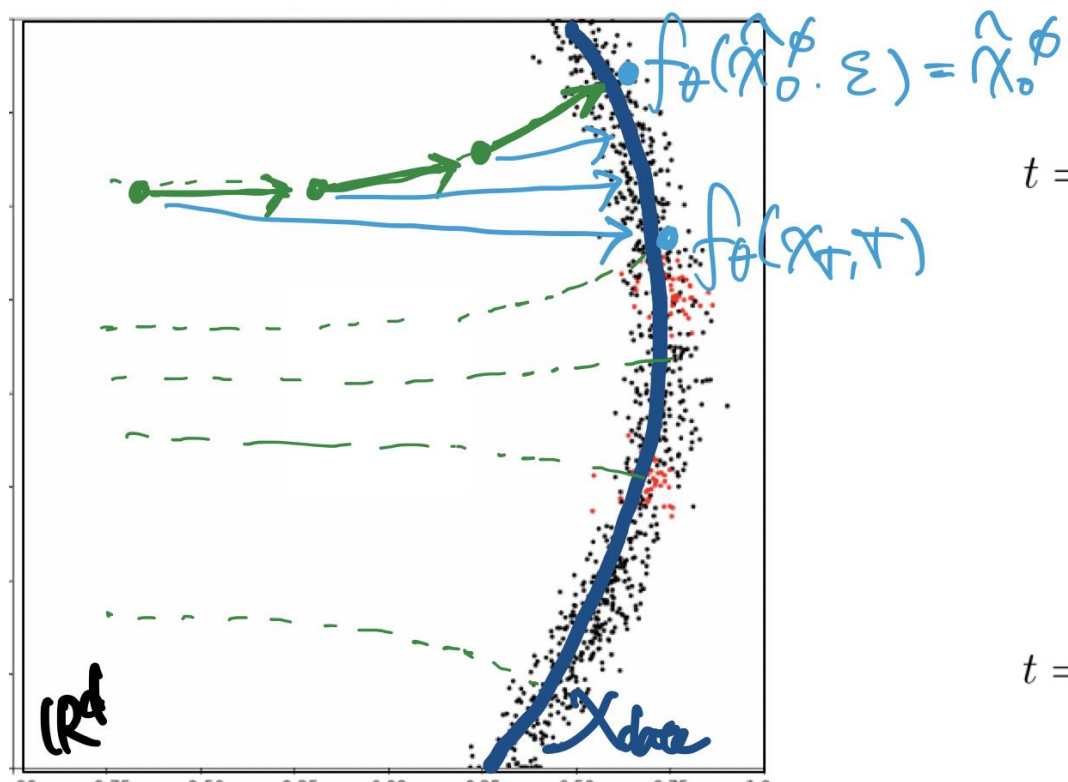
Why “Distillation”?

- Pretrained ODE solver for data pairs
 - Loss function is bound by the ODE solver
 - Model trains to be identical to ODE at $t = 1$
-

$$\lim_{iter \rightarrow \infty} f_{\theta} \neq f(x, t), \text{ rather } f(x, t; \varphi)$$

Theorem 1. Let $\Delta t := \max_{n \in \llbracket 1, N-1 \rrbracket} \{|t_{n+1} - t_n|\}$, and $\mathbf{f}(\cdot, \cdot; \phi)$ be the consistency function of the empirical PF ODE in Eq. (3). Assume \mathbf{f}_{θ} satisfies the Lipschitz condition: there exists $L > 0$ such that for all $t \in [\epsilon, T]$, \mathbf{x} , and \mathbf{y} , we have $\|\mathbf{f}_{\theta}(\mathbf{x}, t) - \mathbf{f}_{\theta}(\mathbf{y}, t)\|_2 \leq L \|\mathbf{x} - \mathbf{y}\|_2$. Assume further that for all $n \in \llbracket 1, N-1 \rrbracket$, the ODE solver called at t_{n+1} has local error uniformly bounded by $O((t_{n+1} - t_n)^{p+1})$ with $p \geq 1$. Then, if $\mathcal{L}_{CD}^N(\theta, \theta; \phi) = 0$, we have

$$\sup_{n, \mathbf{x}} \|\mathbf{f}_{\theta}(\mathbf{x}, t_n) - \mathbf{f}(\mathbf{x}, t_n; \phi)\|_2 = O((\Delta t)^p).$$



Consistency Training

Algorithm 3 Consistency Training (CT)

Input: dataset \mathcal{D} , initial model parameter θ , learning rate η , step schedule $N(\cdot)$, EMA decay rate schedule $\mu(\cdot)$, $d(\cdot, \cdot)$, and $\lambda(\cdot)$

$\theta^- \leftarrow \theta$ and $k \leftarrow 0$

repeat

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

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

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

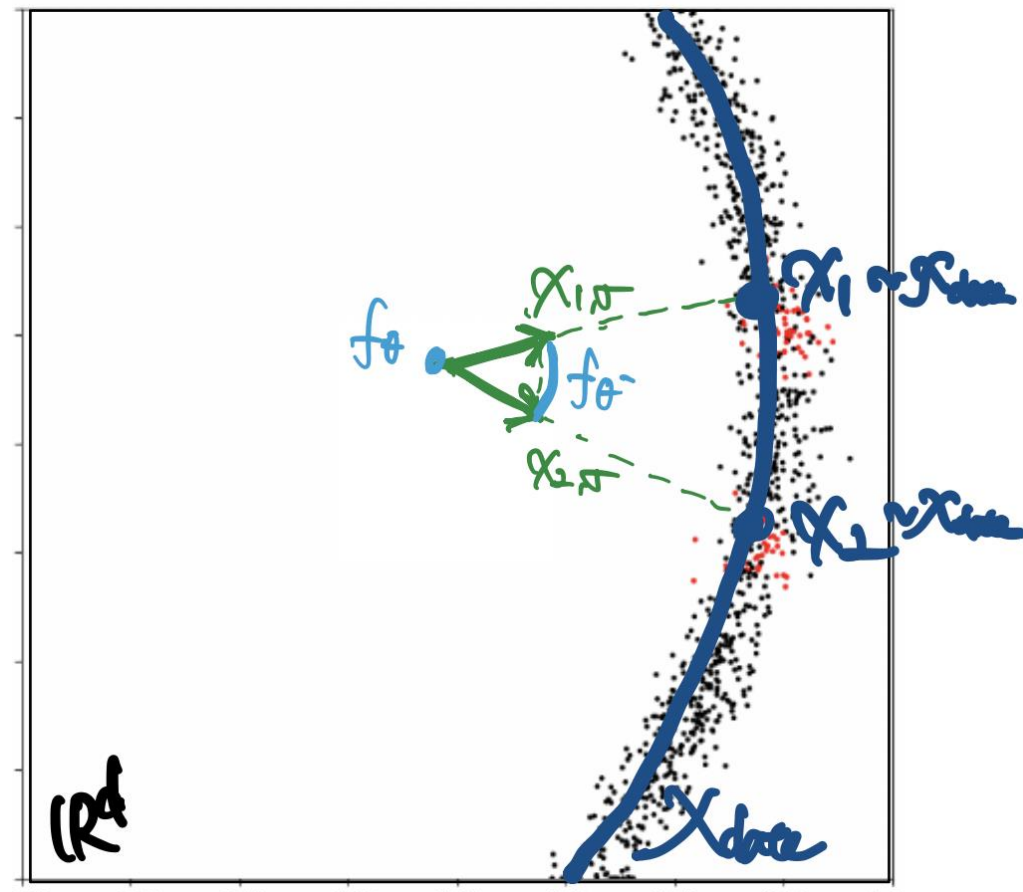
$\lambda(t_n)d(\mathbf{f}_\theta(\mathbf{x} + t_{n+1}\mathbf{z}, t_{n+1}), \mathbf{f}_{\theta^-}(\mathbf{x} + t_n\mathbf{z}, t_n))$

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

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

$k \leftarrow k + 1$

until convergence



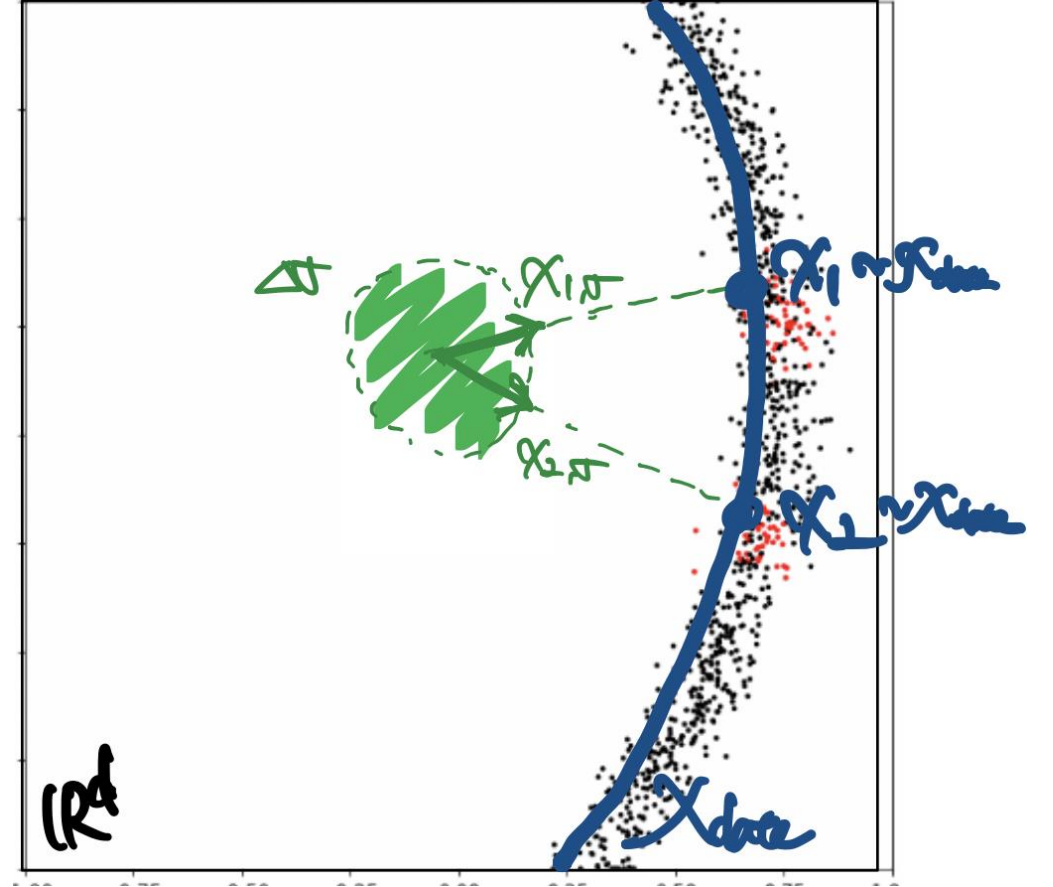
Theorem 2. Let $\Delta t := \max_{n \in \llbracket 1, N-1 \rrbracket} \{|t_{n+1} - t_n|\}$. Assume d and $\mathbf{f}_{\boldsymbol{\theta}^-}$ are both twice continuously differentiable with bounded second derivatives, the weighting function $\lambda(\cdot)$ is bounded, and $\mathbb{E}[\|\nabla \log p_{t_n}(\mathbf{x}_{t_n})\|_2^2] < \infty$. Assume further that we use the Euler ODE solver, and the pre-trained score model matches the ground truth, i.e., $\forall t \in [\epsilon, T] : \mathbf{s}_{\phi}(\mathbf{x}, t) \equiv \nabla \log p_t(\mathbf{x})$. Then,

$$\mathcal{L}_{CD}^N(\boldsymbol{\theta}, \boldsymbol{\theta}^-; \phi) = \mathcal{L}_{CT}^N(\boldsymbol{\theta}, \boldsymbol{\theta}^-) + o(\Delta t), \quad (9)$$

where the expectation is taken with respect to $\mathbf{x} \sim p_{\text{data}}$, $n \sim \mathcal{U}[\llbracket 1, N-1 \rrbracket]$, and $\mathbf{x}_{t_{n+1}} \sim \mathcal{N}(\mathbf{x}; t_{n+1}^2 \mathbf{I})$. The consistency training objective, denoted by $\mathcal{L}_{CT}^N(\boldsymbol{\theta}, \boldsymbol{\theta}^-)$, is defined as

$$\mathbb{E}[\lambda(t_n) d(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x} + t_{n+1} \mathbf{z}, t_{n+1}), \mathbf{f}_{\boldsymbol{\theta}^-}(\mathbf{x} + t_n \mathbf{z}, t_n))], \quad (10)$$

where $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Moreover, $\mathcal{L}_{CT}^N(\boldsymbol{\theta}, \boldsymbol{\theta}^-) \geq O(\Delta t)$ if $\inf_N \mathcal{L}_{CD}^N(\boldsymbol{\theta}, \boldsymbol{\theta}^-; \phi) > 0$.



$$N(k) = \left\lceil \sqrt{\frac{k}{K} ((s_1 + 1)^2 - s_0^2) + s_0^2 - 1} \right\rceil + 1$$

$$\mu(k) = \exp \left(\frac{s_0 \log \mu_0}{N(k)} \right),$$

Metrics

Table 1: Sample quality on CIFAR-10. *Methods that require synthetic data construction for distillation.

METHOD	NFE (↓)	FID (↓)	IS (↑)
Diffusion + Samplers			
DDIM (Song et al., 2020)	50	4.67	
DDIM (Song et al., 2020)	20	6.84	
DDIM (Song et al., 2020)	10	8.23	
DPM-solver-2 (Lu et al., 2022)	10	5.94	
DPM-solver-fast (Lu et al., 2022)	10	4.70	
3-DEIS (Zhang & Chen, 2022)	10	4.17	
Diffusion + Distillation			
Knowledge Distillation* (Luhman & Luhman, 2021)	1	9.36	
DFNO* (Zheng et al., 2022)	1	4.12	
1-Rectified Flow (+distill)* (Liu et al., 2022)	1	6.18	9.08
2-Rectified Flow (+distill)* (Liu et al., 2022)	1	4.85	9.01
3-Rectified Flow (+distill)* (Liu et al., 2022)	1	5.21	8.79
PD (Salimans & Ho, 2022)	1	8.34	8.69
CD	1	3.55	9.48
PD (Salimans & Ho, 2022)	2	5.58	9.05
CD	2	2.93	9.75

METHOD	NFE (↓)	FID (↓)	IS (↑)
Direct Generation			
BigGAN (Brock et al., 2019)	1	14.7	9.22
Diffusion GAN (Xiao et al., 2022)	1	14.6	8.93
AutoGAN (Gong et al., 2019)	1	12.4	8.55
E2GAN (Tian et al., 2020)	1	11.3	8.51
ViTGAN (Lee et al., 2021)	1	6.66	9.30
TransGAN (Jiang et al., 2021)	1	9.26	9.05
StyleGAN2-ADA (Karras et al., 2020)	1	2.92	9.83
<u>StyleGAN-XL</u> (Sauer et al., 2022)	1	1.85	
Score SDE (Song et al., 2021)	2000	2.20	9.89
DDPM (Ho et al., 2020)	1000	3.17	9.46
LSGM (Vahdat et al., 2021)	147	2.10	
PFGM (Xu et al., 2022)	110	2.35	9.68
<u>EDM</u> (Karras et al., 2022)	35	2.04	9.84
1-Rectified Flow (Liu et al., 2022)	1	378	1.13
Glow (Kingma & Dhariwal, 2018)	1	48.9	3.92
Residual Flow (Chen et al., 2019)	1	46.4	
GLFlow (Xiao et al., 2019)	1	44.6	
DenseFlow (Grcić et al., 2021)	1	34.9	
DC-VAE (Parmar et al., 2021)	1	17.9	8.20
CT	1	8.70	8.49
CT	2	5.83	8.85

Metrics

Table 1: Sample quality on CIFAR-10. *Methods that require synthetic data construction for distillation.

METHOD	NFE (↓)	FID (↓)	IS (↑)
Diffusion + Samplers			
DDIM (Song et al., 2020)	50	4.67	
DDIM (Song et al., 2020)	20	6.84	
DDIM (Song et al., 2020)	10	8.23	
DPM-solver-2 (Lu et al., 2022)	10	5.94	
DPM-solver-fast (Lu et al., 2022)	10	4.70	
3-DEIS (Zhang & Chen, 2022)	10	4.17	
Diffusion + Distillation			
Knowledge Distillation* (Luhman & Luhman, 2021)	1	9.36	
DFNO* (Zheng et al., 2022)	1	4.12	
1-Rectified Flow (+distill)* (Liu et al., 2022)	1	6.18	9.08
2-Rectified Flow (+distill)* (Liu et al., 2022)	1	4.85	9.01
3-Rectified Flow (+distill)* (Liu et al., 2022)	1	5.21	8.79
PD (Salimans & Ho, 2022)	1	8.34	8.69
CD	1	3.55	9.48
PD (Salimans & Ho, 2022)	2	5.58	9.05
CD	2	2.93	9.75

METHOD	NFE (↓)	FID (↓)	IS (↑)
Direct Generation			
BigGAN (Brock et al., 2019)	1	14.7	9.22
Diffusion GAN (Xiao et al., 2022)	1	14.6	8.93
AutoGAN (Gong et al., 2019)	1	12.4	8.55
E2GAN (Tian et al., 2020)	1	11.3	8.51
ViTGAN (Lee et al., 2021)	1	6.66	9.30
TransGAN (Jiang et al., 2021)	1	9.26	9.05
StyleGAN2-ADA (Karras et al., 2020)	1	2.92	9.83
<u>StyleGAN-XL</u> (Sauer et al., 2022)	1	1.85	
Score SDE (Song et al., 2021)	2000	2.20	9.89
DDPM (Ho et al., 2020)	1000	3.17	9.46
LSGM (Vahdat et al., 2021)	147	2.10	
PFGM (Xu et al., 2022)	110	2.35	9.68
<u>EDM</u> (Karras et al., 2022)	35	2.04	9.84
1-Rectified Flow (Liu et al., 2022)	1	378	1.13
Glow (Kingma & Dhariwal, 2018)	1	48.9	3.92
Residual Flow (Chen et al., 2019)	1	46.4	
GLFlow (Xiao et al., 2019)	1	44.6	
DenseFlow (Grcić et al., 2021)	1	34.9	
DC-VAE (Parmar et al., 2021)	1	17.9	8.20
CT	1	8.70	8.49
CT	2	5.83	8.85

Metrics

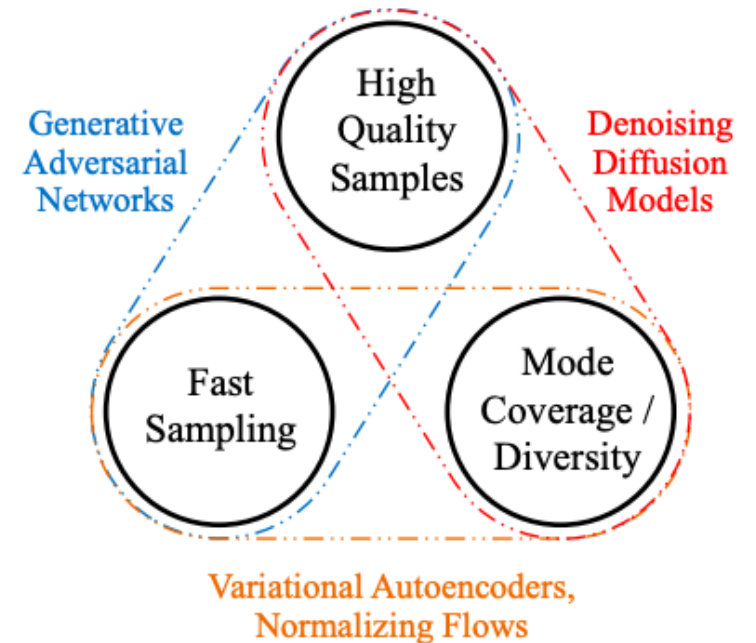
Table 1: Sample quality on CIFAR-10. *Methods that require synthetic data construction for distillation.

METHOD	NFE (↓)	FID (↓)	IS (↑)
Diffusion + Samplers			
DDIM (Song et al., 2020)	50	4.67	
DDIM (Song et al., 2020)	20	6.84	
DDIM (Song et al., 2020)	10	8.23	
DPM-solver-2 (Lu et al., 2022)	10	5.94	
DPM-solver-fast (Lu et al., 2022)	10	4.70	
3-DEIS (Zhang & Chen, 2022)	10	4.17	
Diffusion + Distillation			
Knowledge Distillation* (Luhman & Luhman, 2021)	1	9.36	
DFNO* (Zheng et al., 2022)	1	4.12	
1-Rectified Flow (+distill)* (Liu et al., 2022)	1	6.18	9.08
2-Rectified Flow (+distill)* (Liu et al., 2022)	1	4.85	9.01
3-Rectified Flow (+distill)* (Liu et al., 2022)	1	5.21	8.79
PD (Salimans & Ho, 2022)	1	8.34	8.69
CD	1	3.55	9.48
PD (Salimans & Ho, 2022)	2	5.58	9.05
CD	2	2.93	9.75

METHOD	NFE (↓)	FID (↓)	IS (↑)
Direct Generation			
BigGAN (Brock et al., 2019)	1	14.7	9.22
Diffusion GAN (Xiao et al., 2022)	1	14.6	8.93
AutoGAN (Gong et al., 2019)	1	12.4	8.55
E2GAN (Tian et al., 2020)	1	11.3	8.51
ViTGAN (Lee et al., 2021)	1	6.66	9.30
TransGAN (Jiang et al., 2021)	1	9.26	9.05
StyleGAN2-ADA (Karras et al., 2020)	1	2.92	9.83
<u>StyleGAN-XL</u> (Sauer et al., 2022)	1	1.85	
Score SDE (Song et al., 2021)	2000	2.20	9.89
DDPM (Ho et al., 2020)	1000	3.17	9.46
LSGM (Vahdat et al., 2021)	147	2.10	
PFGM (Xu et al., 2022)	110	2.35	9.68
<u>EDM</u> (Karras et al., 2022)	35	2.04	9.84
1-Rectified Flow (Liu et al., 2022)	1	378	1.13
Glow (Kingma & Dhariwal, 2018)	1	48.9	3.92
Residual Flow (Chen et al., 2019)	1	46.4	
GLFlow (Xiao et al., 2019)	1	44.6	
DenseFlow (Grcić et al., 2021)	1	34.9	
DC-VAE (Parmar et al., 2021)	1	17.9	8.20
CT	1	8.70	8.49
CT	2	5.83	8.85

Implications and Applications

- Consistency Distillation
 - Valid distillation method for score-based models
 - Enable single-step sampling
 - Sustain sample quality
 - Dependency on pre-trained ODE solvers
- Consistency Training
 - Direct generation, without dependency
 - Relatively poor generation performance



Characteristics from Diffusion

- Semantic Continuity
 - Multistep sampling, Inpainting, Interpolation, Denoising, Upscaling...
- Controlled Data Generation
 - Classifier Free Guidance (*LCM*)

Algorithm 1 Multistep Consistency Sampling

Input: Consistency model $f_{\theta}(\cdot, \cdot)$, sequence of time points $\tau_1 > \tau_2 > \dots > \tau_{N-1}$, initial noise $\hat{\mathbf{x}}_T$

$\mathbf{x} \leftarrow f_{\theta}(\hat{\mathbf{x}}_T, T)$

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

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

$\hat{\mathbf{x}}_{\tau_n} \leftarrow \mathbf{x} + \sqrt{\tau_n^2 - \epsilon^2} \mathbf{z}$

$\mathbf{x} \leftarrow f_{\theta}(\hat{\mathbf{x}}_{\tau_n}, \tau_n)$

end for

Output: \mathbf{x}

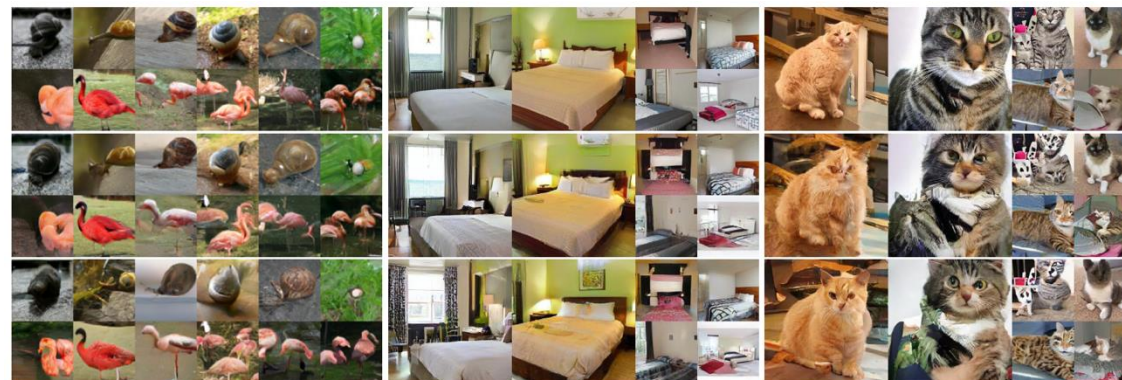


Figure 5: Samples generated by EDM (*top*), CT + single-step generation (*middle*), and CT + 2-step generation (*Bottom*). All corresponding images are generated from the same initial noise.

Algorithm 4 Zero-Shot Image Editing

- 1: **Input:** Consistency model $f_\theta(\cdot, \cdot)$, sequence of time points $t_1 > t_2 > \dots > t_N$, reference image \mathbf{y} , invertible linear transformation \mathbf{A} , and binary image mask Ω
 - 2: $\mathbf{y} \leftarrow \mathbf{A}^{-1}[(\mathbf{A}\mathbf{y}) \odot (1 - \Omega) + \mathbf{0} \odot \Omega]$
 - 3: Sample $\mathbf{x} \sim \mathcal{N}(\mathbf{y}, t_1^2 \mathbf{I})$
 - 4: $\mathbf{x} \leftarrow f_\theta(\mathbf{x}, t_1)$
 - 5: $\mathbf{x} \leftarrow \mathbf{A}^{-1}[(\mathbf{A}\mathbf{y}) \odot (1 - \Omega) + (\mathbf{A}\mathbf{x}) \odot \Omega]$
 - 6: **for** $n = 2$ **to** N **do**
 - 7: Sample $\mathbf{x} \sim \mathcal{N}(\mathbf{x}, (t_n^2 - \epsilon^2) \mathbf{I})$
 - 8: $\mathbf{x} \leftarrow f_\theta(\mathbf{x}, t_n)$
 - 9: $\mathbf{x} \leftarrow \mathbf{A}^{-1}[(\mathbf{A}\mathbf{y}) \odot (1 - \Omega) + (\mathbf{A}\mathbf{x}) \odot \Omega]$
 - 10: **end for**
 - 11: **Output:** \mathbf{x}
-



(a) *Left:* The gray-scale image. *Middle:* Colorized images. *Right:* The ground-truth image.



(b) *Left:* The downsampled image (32×32). *Middle:* Full resolution images (256×256). *Right:* The ground-truth image (256×256).



(c) *Left:* A stroke input provided by users. *Right:* Stroke-guided image generation.

IMPROVED TECHNIQUES FOR TRAINING CONSISTENCY MODELS

Yang Song & Prafulla Dhariwal
OpenAI
{songyang, prafulla}@openai.com

Table 2: Comparing the quality of unconditional samples on CIFAR-10.

METHOD	NFE (↓)	FID (↓)	IS (↑)
Direct Generation			
Score SDE (Song et al., 2021)	2000	2.38	9.83
Score SDE (deep) (Song et al., 2021)	2000	2.20	9.89
DDPM (Ho et al., 2020)	1000	3.17	9.46
LSGM (Vahdat et al., 2021)	147	2.10	
PFGM (Xu et al., 2022)	110	2.35	9.68
EDM* (Karras et al., 2022)	35	2.04	9.84
EDM-G++ (Kim et al., 2023)	35	1.77	
IGEBM (Du & Mordatch, 2019)	60	40.6	6.02
NVAE (Vahdat & Kautz, 2020)	1	23.5	7.18
Glow (Kingma & Dhariwal, 2018)	1	48.9	3.92
Residual Flow (Chen et al., 2019)	1	46.4	
BigGAN (Brock et al., 2019)	1	14.7	9.22
StyleGAN2 (Karras et al., 2020b)	1	8.32	9.21
StyleGAN2-ADA (Karras et al., 2020a)	1	2.92	9.83
CT (LPIPS) (Song et al., 2023)	1	8.70	8.49
	2	5.83	8.85
iCT (ours)	1	2.83	9.54
	2	2.46	9.80
iCT-deep (ours)	1	2.51	9.76
	2	2.24	9.89

Table 1: Comparing the design choices for CT in Song et al. (2023) versus our modifications.

	Design choice in Song et al. (2023)	Our modifications
EMA decay rate for the teacher network	$\mu(k) = \exp(\frac{s_0 \log \mu_0}{N(k)})$	$\mu(k) = 0$
Metric in consistency loss	$d(\mathbf{x}, \mathbf{y}) = \text{LPIPS}(\mathbf{x}, \mathbf{y})$	$d(\mathbf{x}, \mathbf{y}) = \sqrt{\ \mathbf{x} - \mathbf{y}\ _2^2 + c^2} - c$
Discretization curriculum	$N(k) = \left\lceil \sqrt{\frac{k}{K}} ((s_1 + 1)^2 - s_0^2) + s_0^2 - 1 \right\rceil + 1$	$N(k) = \min(s_0 2^{\lfloor \frac{k}{K'} \rfloor}, s_1) + 1,$ where $K' = \left\lfloor \frac{K}{\log_2 \lceil s_1/s_0 \rceil + 1} \right\rfloor$
Noise schedule	σ_i , where $i \sim \mathcal{U}[1, N(k) - 1]$	σ_i , where $i \sim p(i)$, and $p(i) \propto \text{erf}(\frac{\log(\sigma_{i+1}) - P_{\text{mean}}}{\sqrt{2}P_{\text{std}}}) - \text{erf}(\frac{\log(\sigma_i) - P_{\text{mean}}}{\sqrt{2}P_{\text{std}}})$
Weighting function	$\lambda(\sigma_i) = 1$	$\lambda(\sigma_i) = \frac{1}{\sigma_{i+1} - \sigma_i}$
Parameters	$s_0 = 2, s_1 = 150, \mu_0 = 0.9$ on CIFAR-10 $s_0 = 2, s_1 = 200, \mu_0 = 0.95$ on ImageNet 64×64	$s_0 = 10, s_1 = 1280$ $c = 0.00054\sqrt{d}$, d is data dimensionality $P_{\text{mean}} = -1.1, P_{\text{std}} = 2.0$
	$k \in [0, K]$, where K is the total training iterations $\sigma_i = (\sigma_{\min}^{1/\rho} + \frac{i-1}{N(k)-1}(\sigma_{\max}^{1/\rho} - \sigma_{\min}^{1/\rho}))^\rho$, where $i \in [1, N(k)]$, $\rho = 7, \sigma_{\min} = 0.002, \sigma_{\max} = 80$	

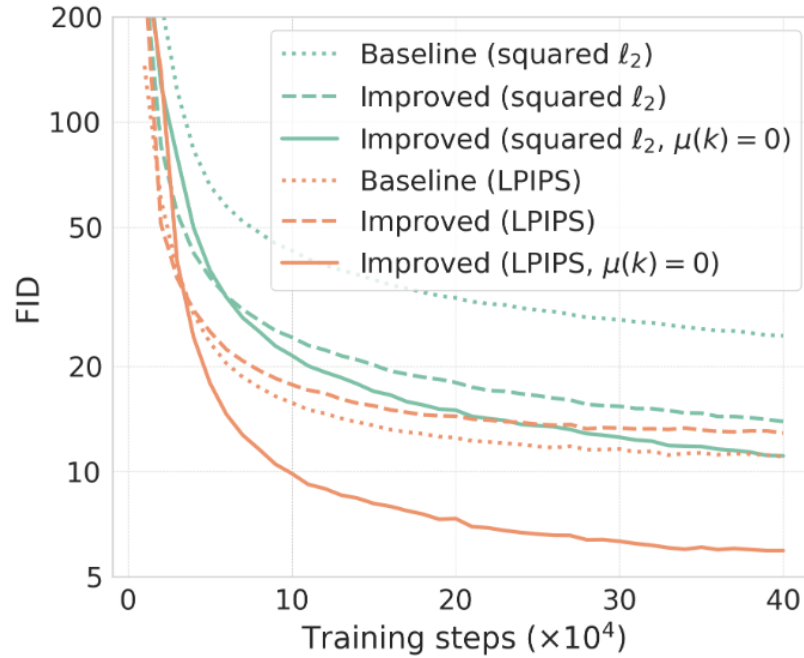
[5] Song, Yang, and Prafulla Dhariwal. "Improved techniques for training consistency models." *arXiv preprint arXiv:2310.14189* (2023).

- Smaller loss weighting on higher noise level
 - Smaller noise level embedding
 - Dropout
 - Discard EMA
 - Pseudo-Huber Loss
 - Discretization step scheduling
 - Noise sampling scheduling
-

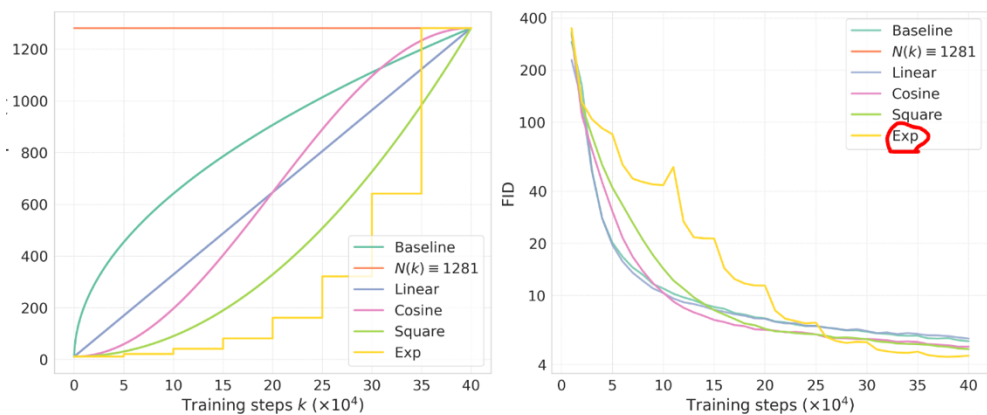
Proposition 1. *Given the notations introduced earlier, and using the uniform weighting function $\lambda(\sigma) = 1$ along with the squared ℓ_2 metric, we have*

$$\lim_{N \rightarrow \infty} \mathcal{L}^N(\theta, \theta^-) = \lim_{N \rightarrow \infty} \mathcal{L}_{CT}^N(\theta, \theta^-) = \mathbb{E} \left[\left(1 - \frac{\sigma_{min}}{\sigma_i} \right)^2 (\theta - \theta^-)^2 \right] \quad \text{if } \theta^- \neq \theta \quad (6)$$

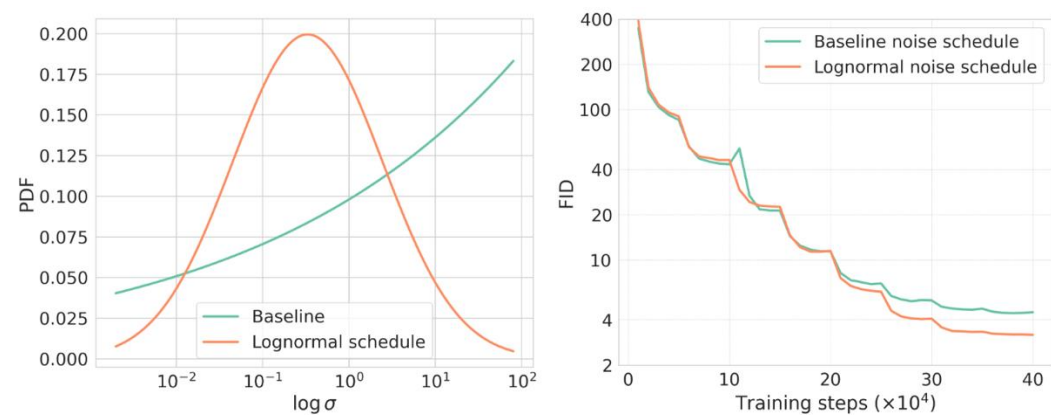
$$\lim_{N \rightarrow \infty} \frac{1}{\Delta\sigma} \frac{d\mathcal{L}^N(\theta, \theta^-)}{d\theta} = \begin{cases} \frac{d}{d\theta} \mathbb{E} \left[\frac{\sigma_{min}}{\sigma_i^2} \left(1 - \frac{\sigma_{min}}{\sigma_i} \right) (\theta - \xi)^2 \right], & \theta^- = \theta \\ +\infty, & \theta^- < \theta \\ -\infty, & \theta^- > \theta \end{cases} \quad (7)$$



(a) LPIPS & squared ℓ_2 metrics.



(b) Various curriculums for $N(k)$. (c) FIDs vs. $N(k)$ curriculums.



(a) PDF of $\log \sigma$ (b) Lognormal vs. default schedules.

$$N(k) = \min(s_0 2^{\lfloor \frac{k}{K'} \rfloor}, s_1) + 1,$$

$$\text{where } K' = \left\lfloor \frac{K}{\log_2[s_1/s_0] + 1} \right\rfloor$$

$$p(\sigma_i) \propto \text{erf} \left(\frac{\log(\sigma_{i+1}) - P_{\text{mean}}}{\sqrt{2}P_{\text{std}}} \right) - \text{erf} \left(\frac{\log(\sigma_i) - P_{\text{mean}}}{\sqrt{2}P_{\text{std}}} \right),$$

LATENT CONSISTENCY MODELS: SYNTHESIZING HIGH-RESOLUTION IMAGES WITH FEW-STEP INFERENCE

Simian Luo* **Yiqin Tan*** **Longbo Huang[†]** **Jian Li[†]** **Hang Zhao[†]**
 Institute for Interdisciplinary Information Sciences, Tsinghua University
 {luosm22, tyq22}@mails.tsinghua.edu.cn
 {longbohuang, lijian83, hangzhao}@tsinghua.edu.cn

$$\tilde{\epsilon}_{\theta}(z_t, w, c, t) = (1 + w)\epsilon_{\theta}(z_t, c, t) - w\epsilon_{\theta}(z_t, \emptyset, t)$$

CFG Conditional diffusion model Unconditional diffusion model

PF ODE $\frac{dx}{dt} = f(t)z_t + \frac{g^2(t)}{2\sigma_t} \epsilon_{\theta}(z_t, c, t)$

Augmented PF ODE $\frac{dx}{dt} = f(t)z_t + \frac{g^2(t)}{2\sigma_t} \tilde{\epsilon}_{\theta}(z_t, w, c, t)$

CFG

$$\mathcal{L}_{CD}(\theta, \theta^-; \Psi) = \mathbb{E}_{\mathbf{z}, \mathbf{c}, n} \left[d \left(\mathbf{f}_{\theta}(\mathbf{z}_{t_{n+1}}, \mathbf{c}, t_{n+1}), \mathbf{f}_{\theta^-}(\hat{\mathbf{z}}_{t_n}^{\Psi}, \mathbf{c}, t_n) \right) \right]$$



→ Efficient sampling / High resolution / Text Conditioning

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

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

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

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

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

until convergence



Finetuning, CD \rightarrow CT conversion

