



Refining Generative Process with Discriminator Guidance in Score-based Diffusion Models

Dongjun Kim*, Yeongmin Kim*,
Se Jung Kwon, Wanmo Kang, Il-Chul Moon
KAIST

* Equal contribution

Before



- ✓ Discretization Error
- ✓ Score Estimation Error
- ✓ Prior Mismatch Error

Before



- ✓ Discretization Error
- ✓ **Score Estimation Error**
- ✓ Prior Mismatch Error

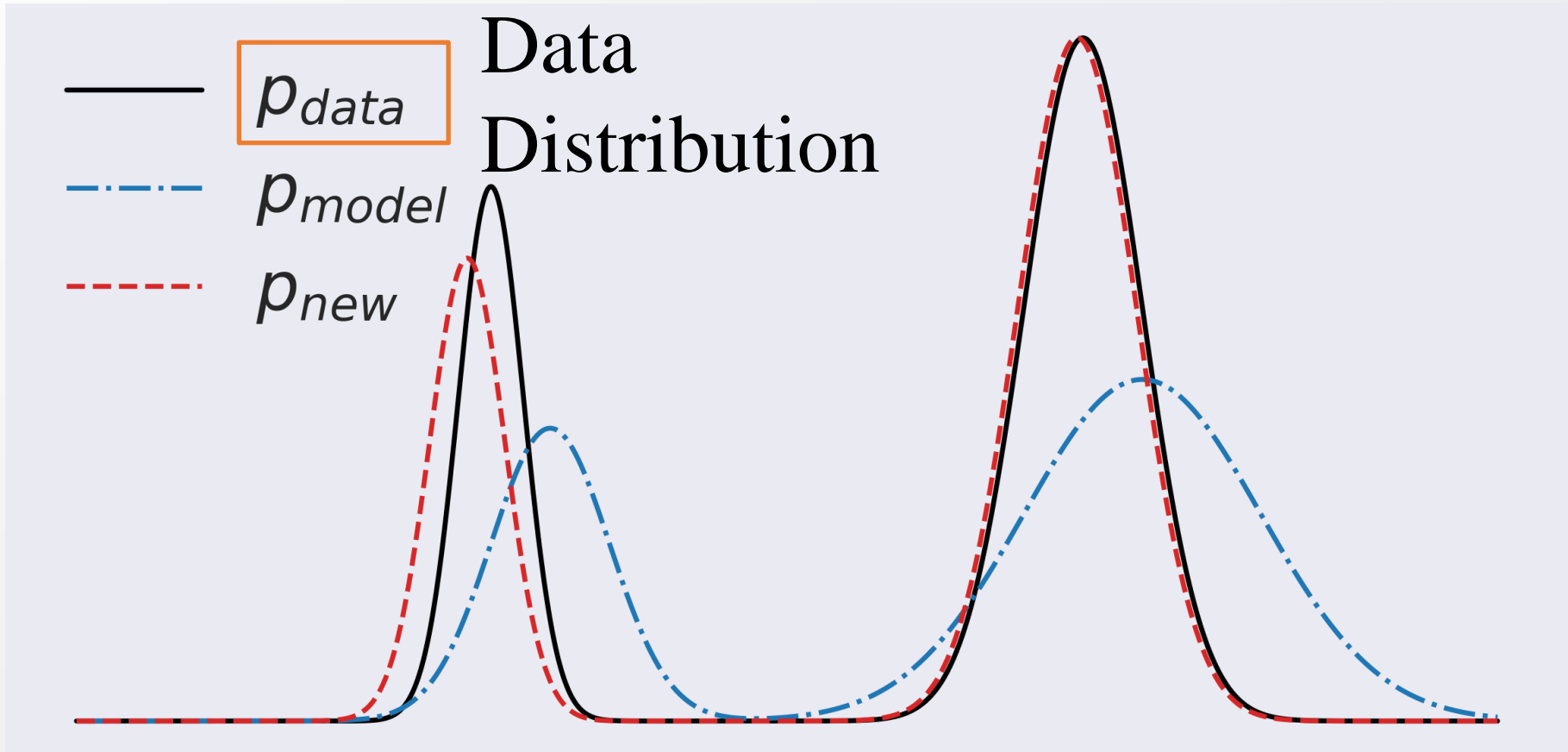
Our Method Improves Sample Quality

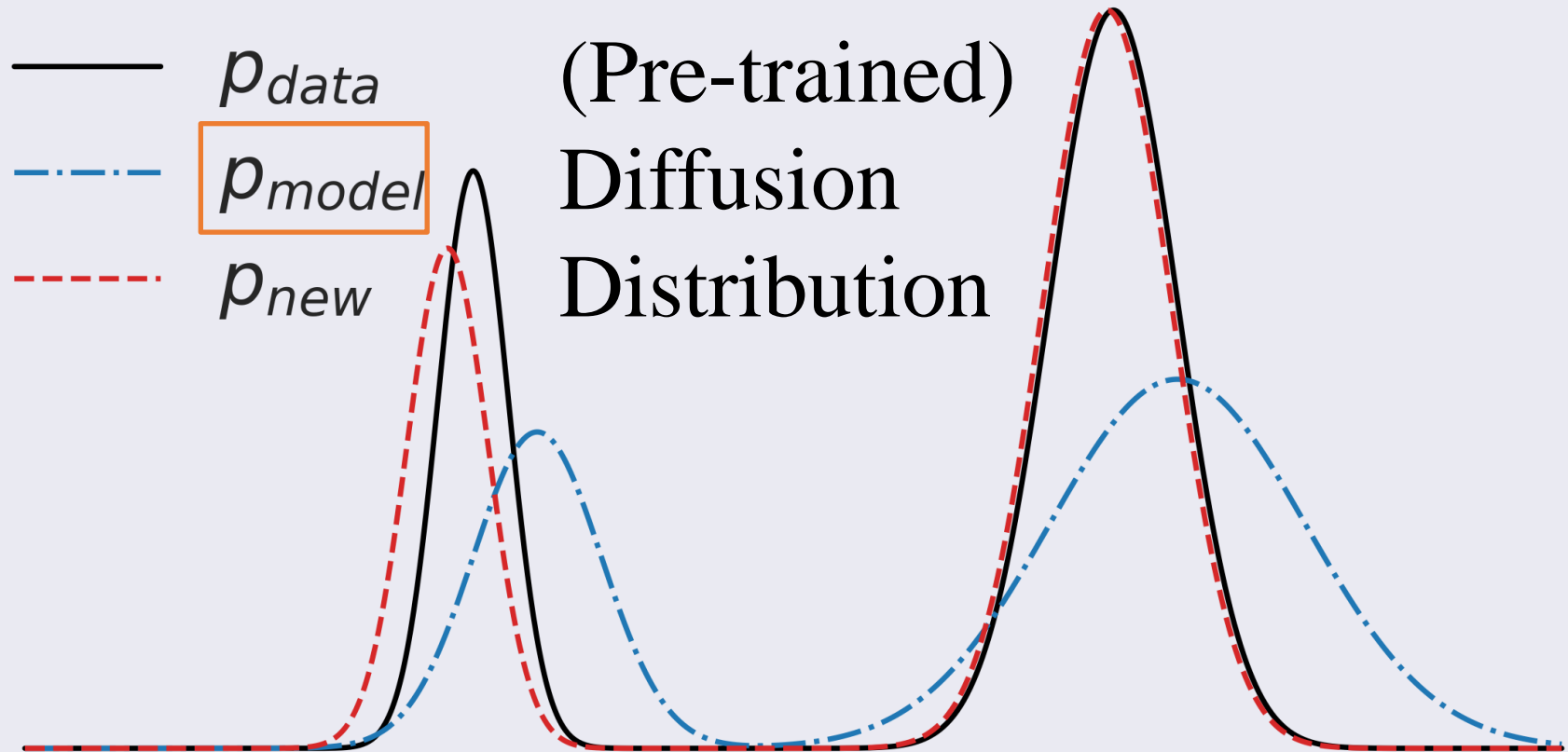
Before

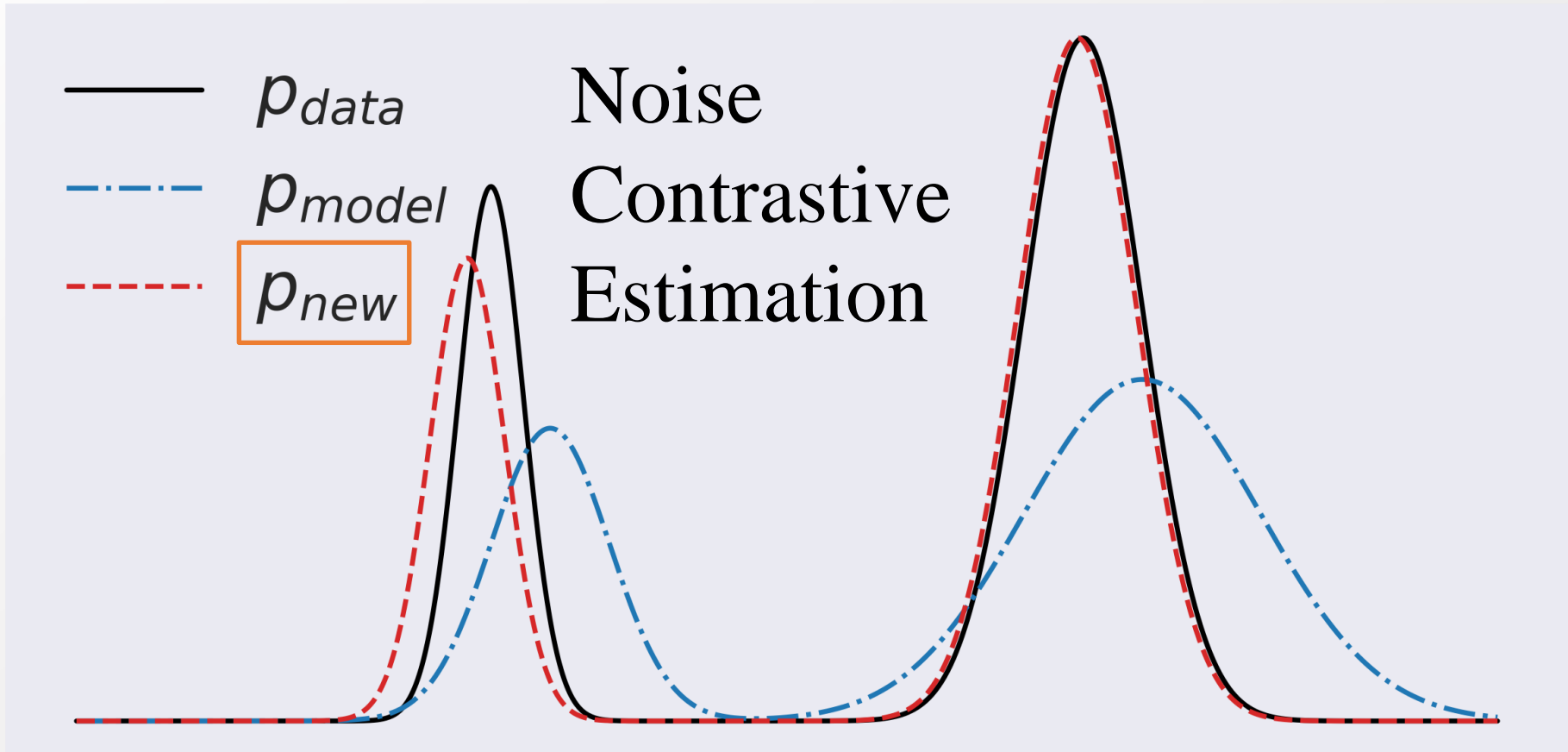


After

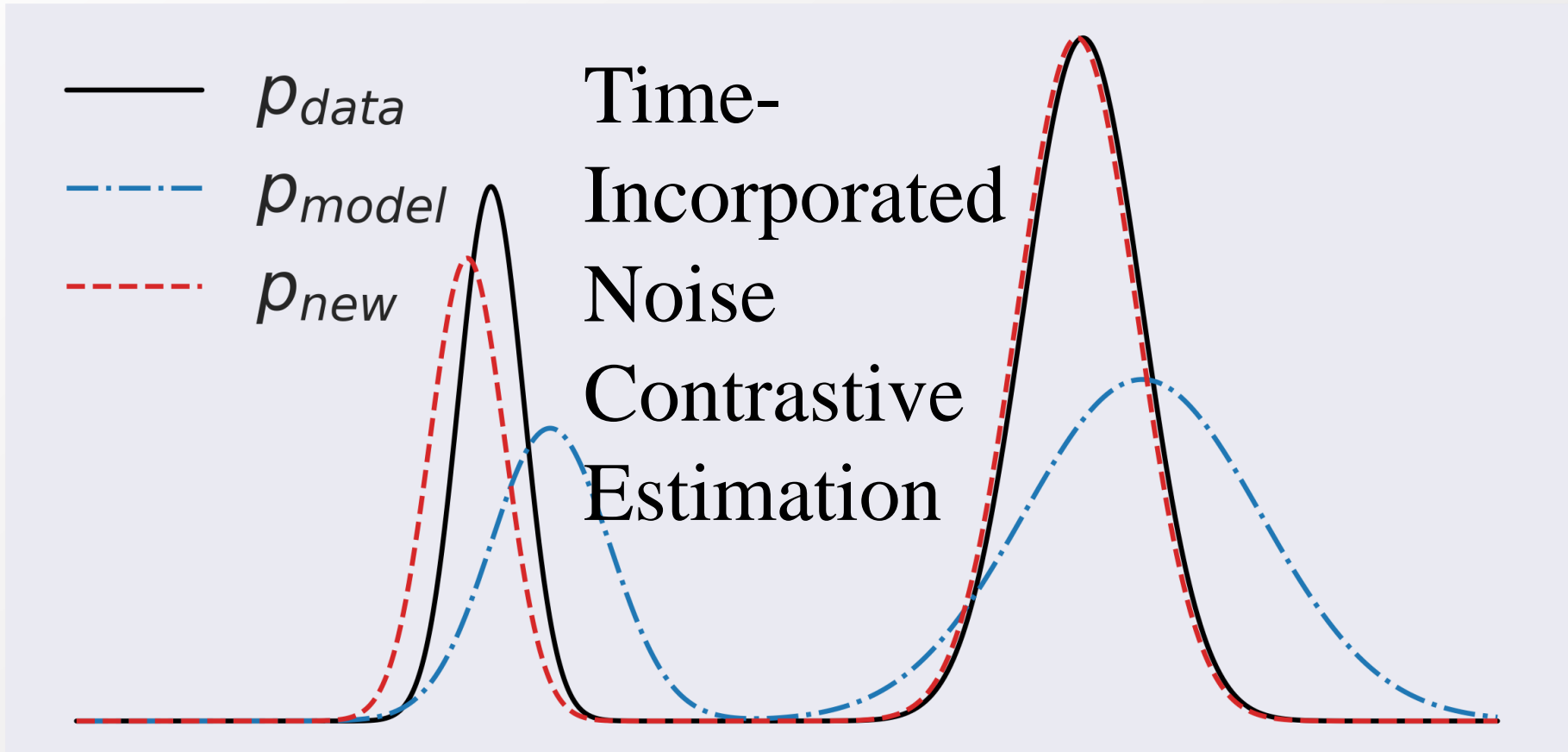




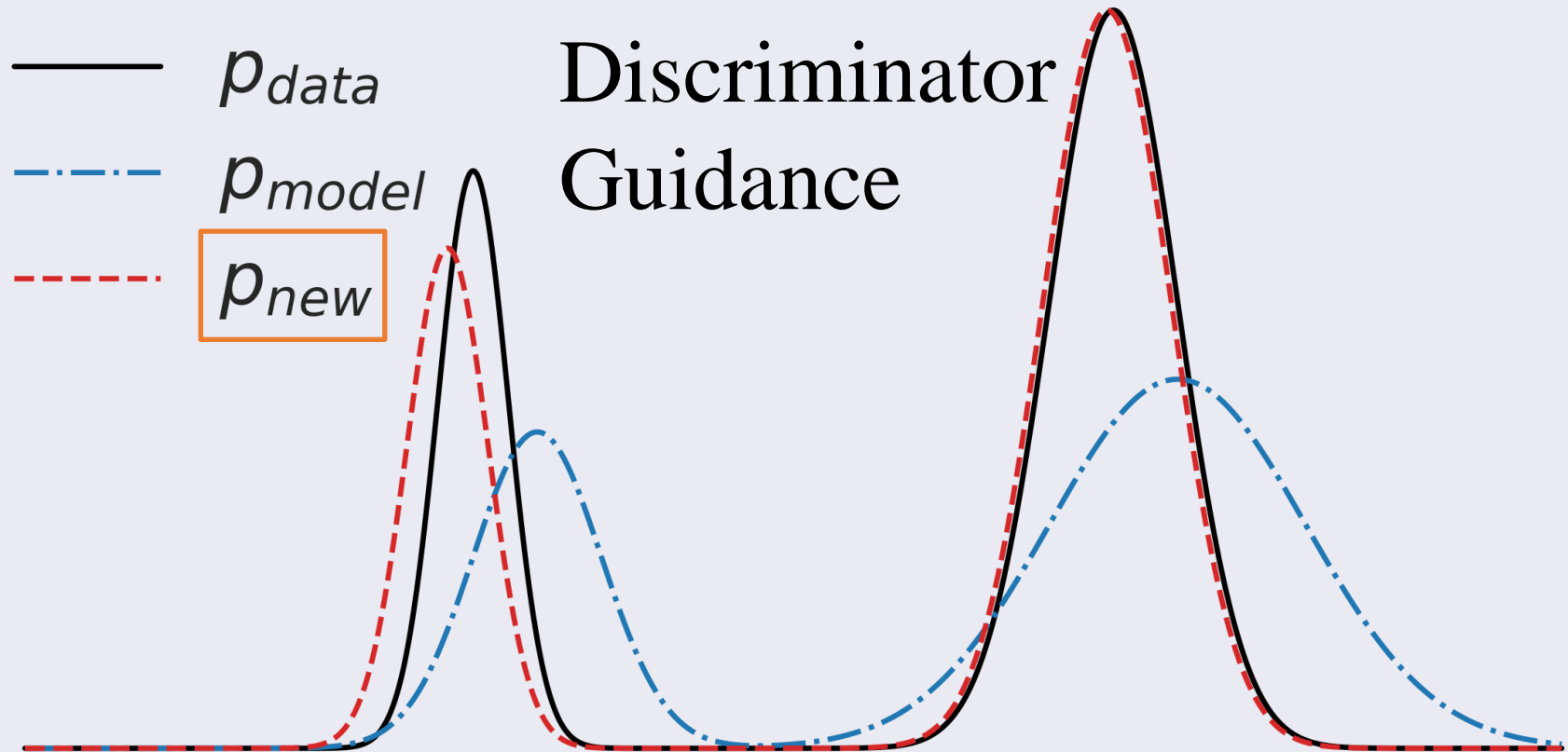




$$p_{new}(\mathbf{x}) \propto p_{model}(\mathbf{x}) \frac{d(\mathbf{x})}{1-d(\mathbf{x})}$$

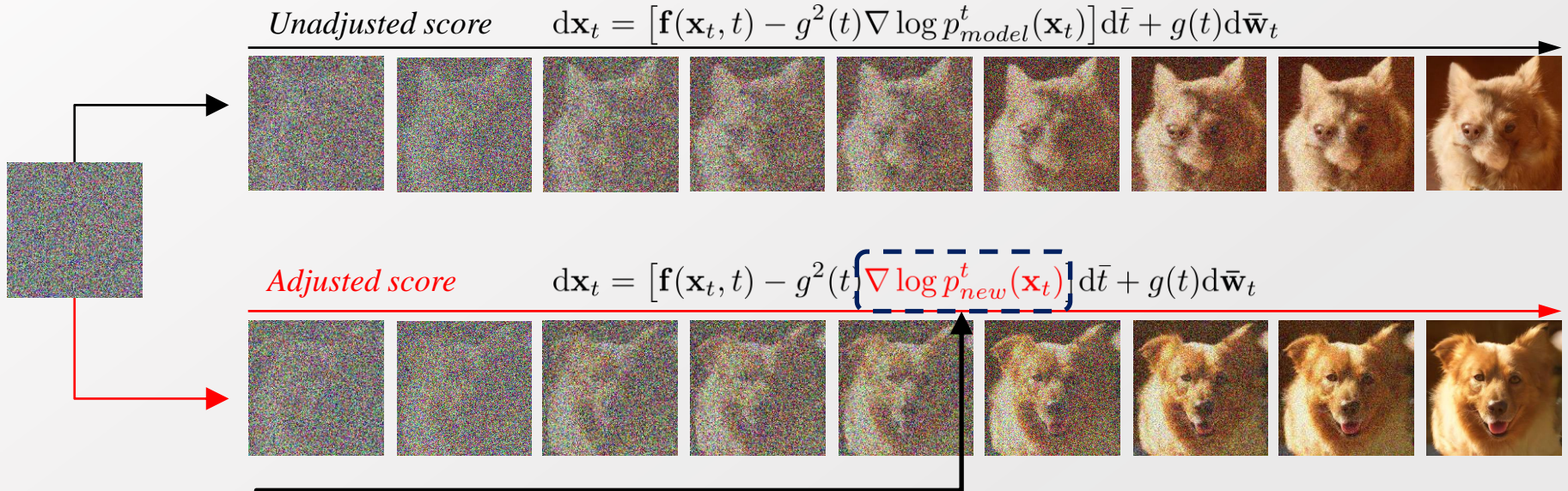


$$p_{new}^t(\mathbf{x}_t) \propto p_{model}^t(\mathbf{x}_t) \frac{d(\mathbf{x}_t, t)}{1 - d(\mathbf{x}_t, t)}$$



$$\nabla \log p_{new}^t(\mathbf{x}_t) = \underbrace{\nabla \log p_{model}^t(\mathbf{x}_t)}_{\text{Freeze}} + \underbrace{\nabla \log \frac{d(\mathbf{x}_t, t)}{1 - d(\mathbf{x}_t, t)}}_{\text{Disc. Guidance}}$$

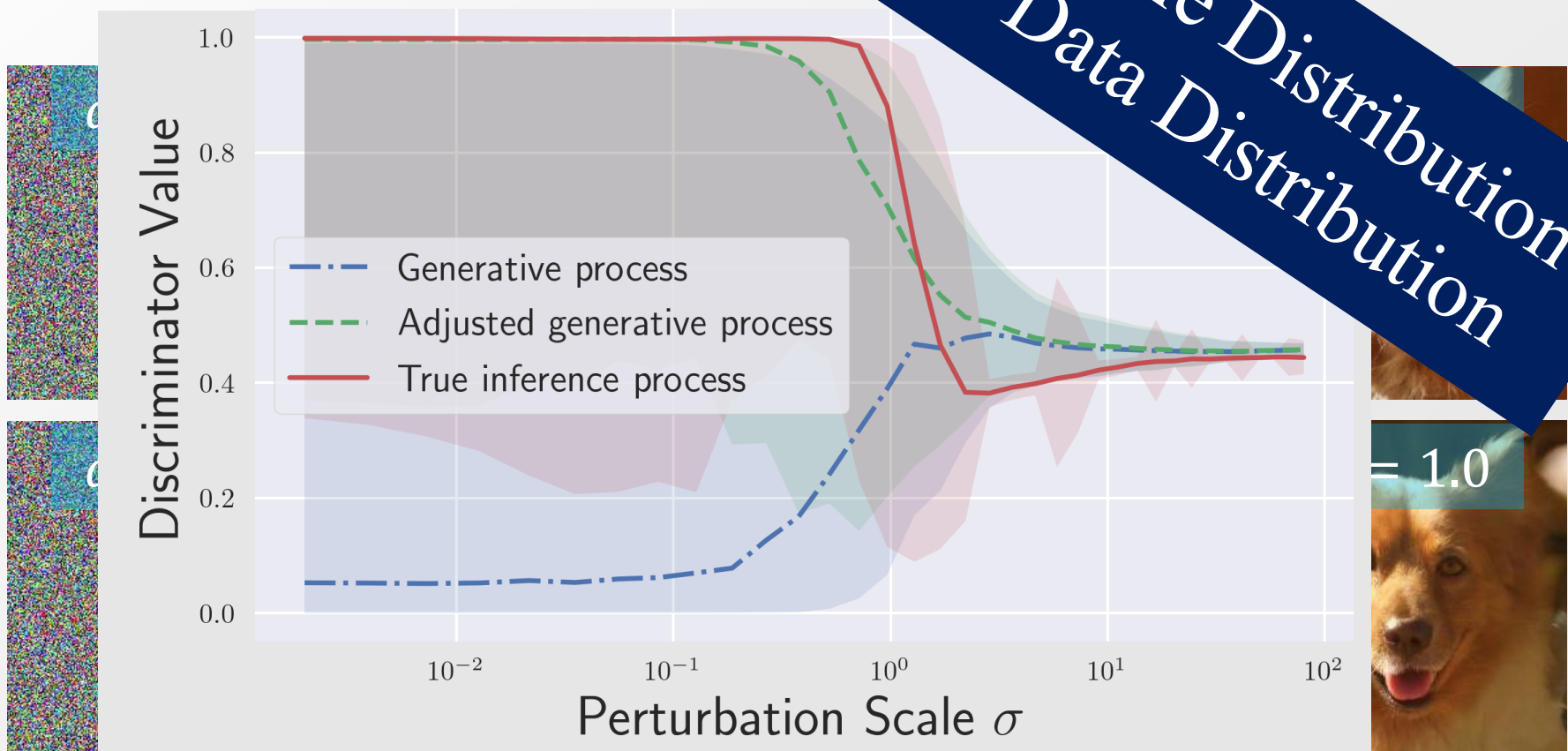
Sampling of Discriminator Guidance



$$\nabla \log p_{new}^t(\mathbf{x}_t) = \nabla \log p_{model}^t(\mathbf{x}_t) + \nabla \log \frac{d(\mathbf{x}_t, t)}{1 - d(\mathbf{x}_t, t)}$$

Meaning of Sampling

Sample Distribution = Data Distribution

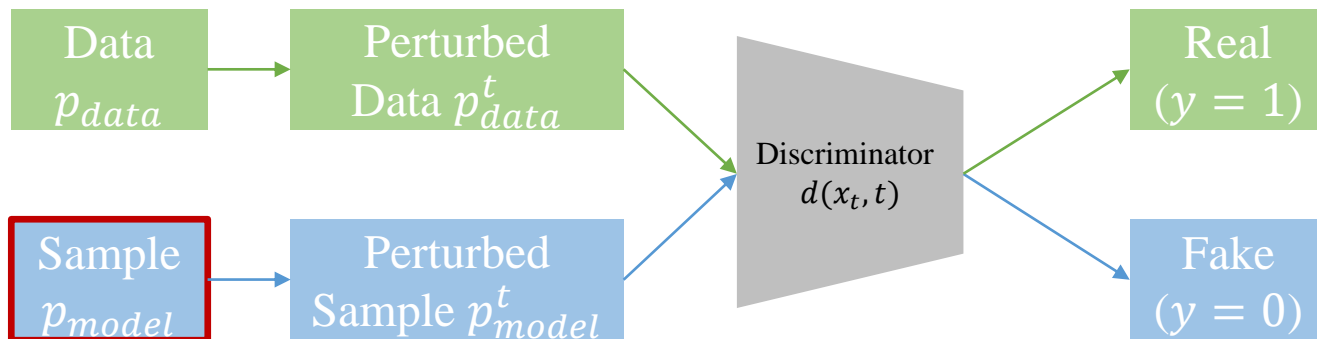


$$p_{new}^t(\mathbf{x}_t) \propto p_{model}^t(\mathbf{x}_t) \frac{d(\mathbf{x}_t, t)}{1 - d(\mathbf{x}_t, t)}$$

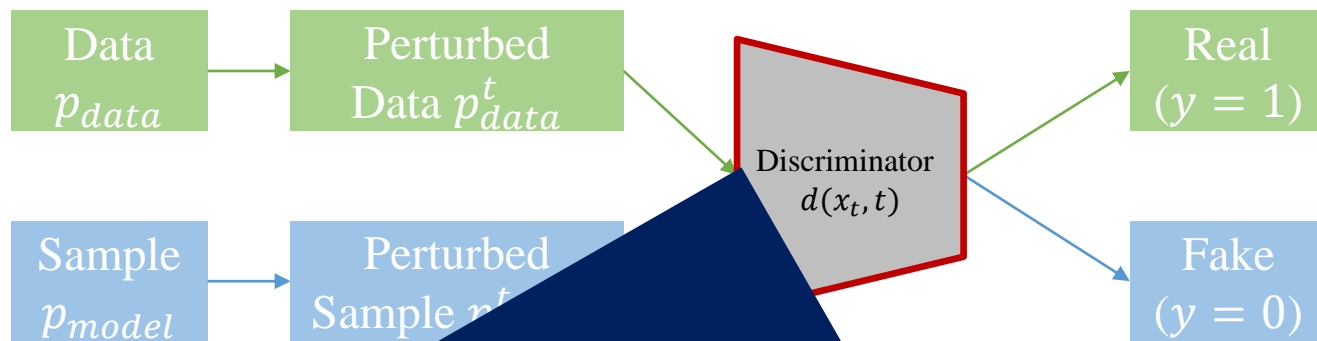
↑↑
↑↑
↑↑

Learning Details

Training Step 1



Training Step 2



Cheap Training

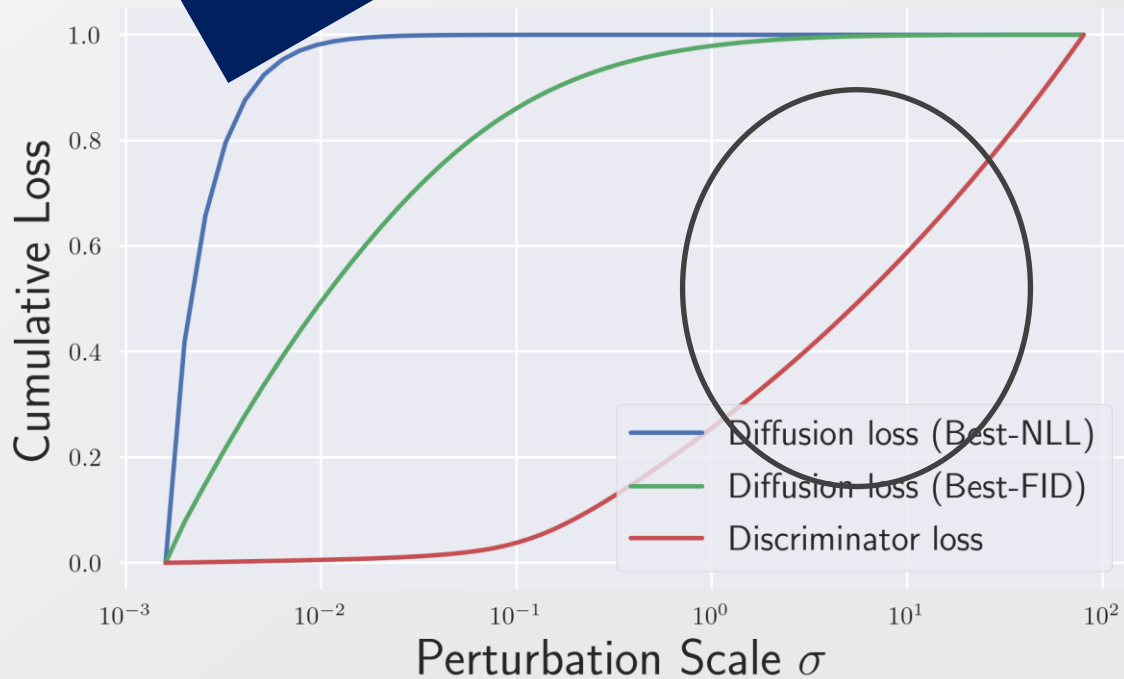
Dataset (Model)	Score Training	Discriminator Training (Step 2)
CIFAR-10 (EDM)	1M score eval (1hr)	1M discriminator eval (10min)
ImageNet 256x256 (ADM)	100M score eval	25.6M discriminator eval

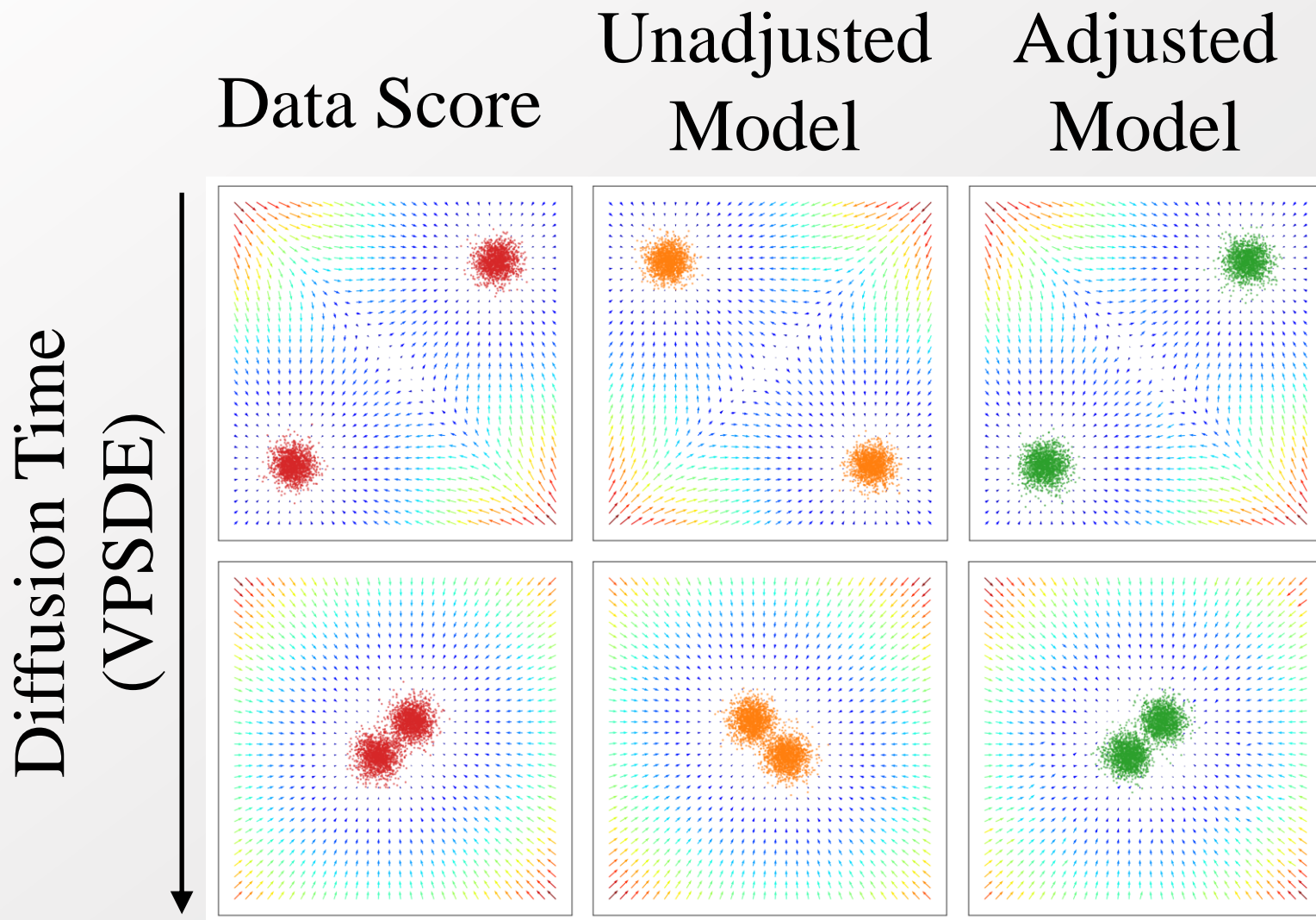
Binary

$$\mathcal{L}(\psi; \lambda) = \int_0^T \lambda dt$$

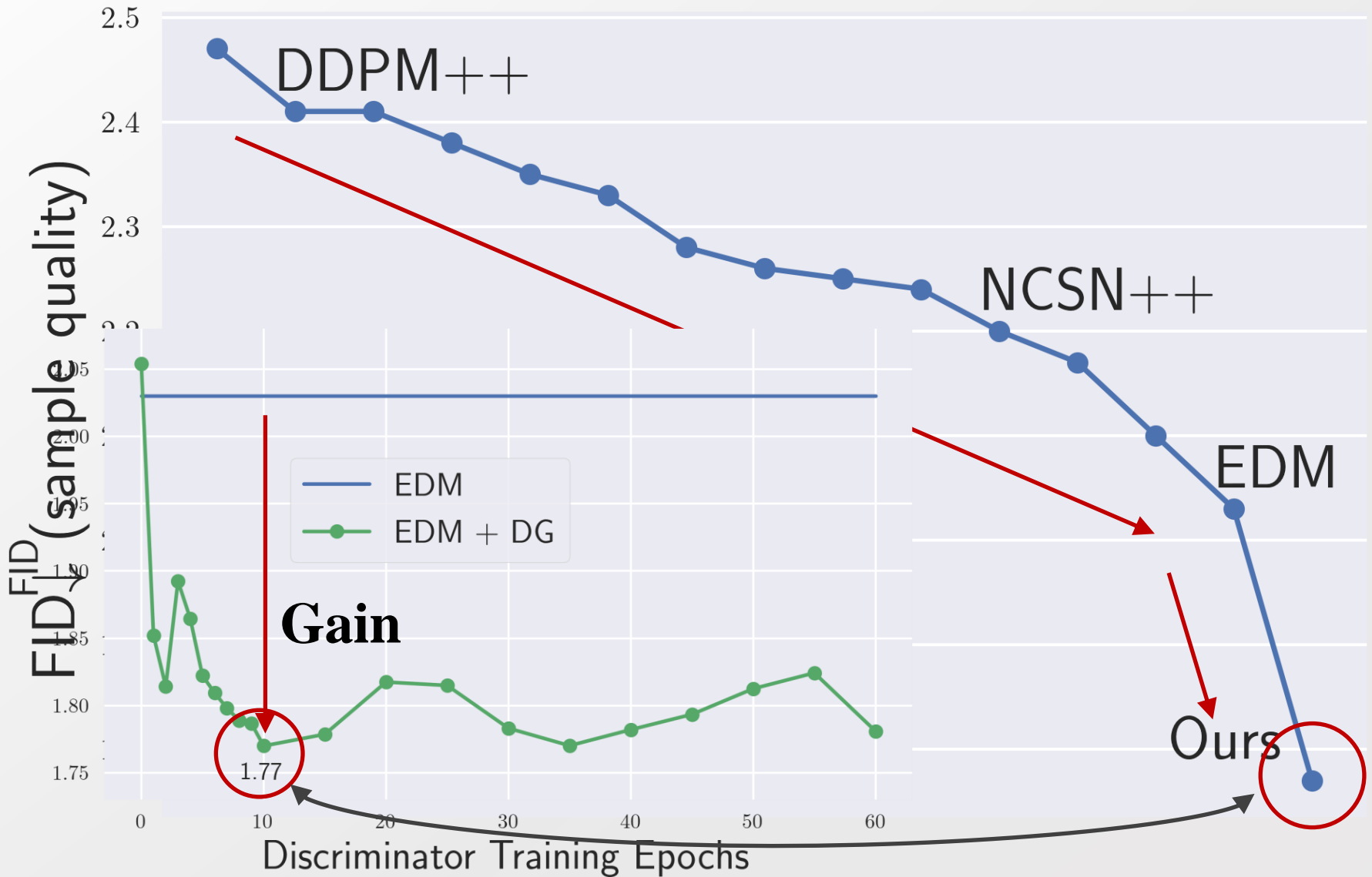
$$BCE_t(\psi) = \mathbb{E}_{p_{data}(x_t)} [-\log d_\psi(x_t, t)] + \mathbb{E}_{p_{\theta_\infty}(x_0)q(x_t|x_0)} [-\log (1 - d_\psi(x_t, t))]$$

Adv. Sampling but
Stable Training





Record-breaking Quality in CIFAR-10



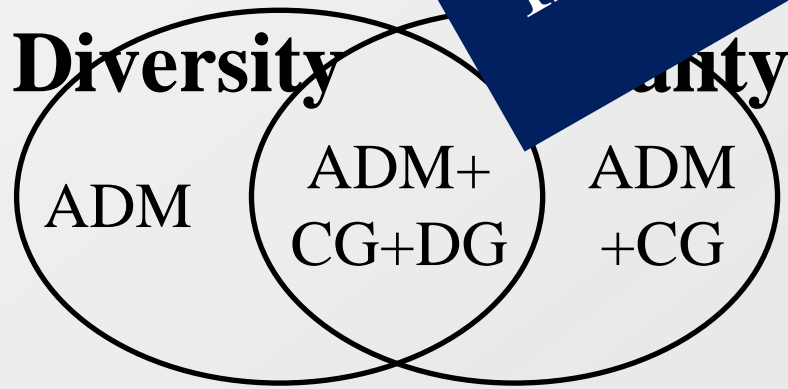
Record-breaking Quality & Diversity in ImageNet 256x256

Model	Diffusion Space	FID↓ (Quality)	Recall↑ (Diversity)
Validation Data		1.68	0.66
ADM	Data	10.94	0.63
ADM+CG	Data	4.59	0.52
ADM+CG+DG	Data	3.73	0.50
DiT-XL/2	Latent	9.62	
DiT-XL/2+CG	Latent	2.27	
DiT-XL/2+CG+DG	Latent		

ADM



Intra-Class Diversity

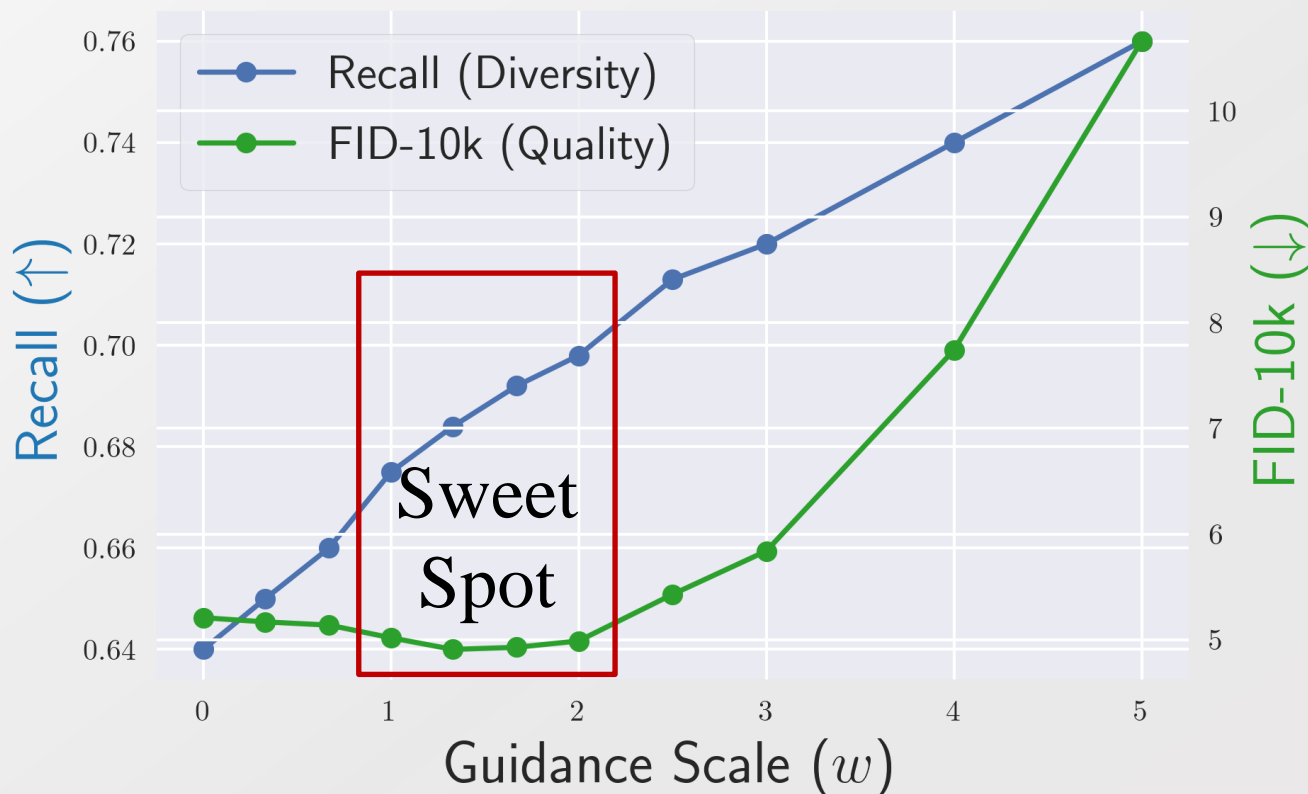


ADM
+CG+DG

Why Better Diversity?

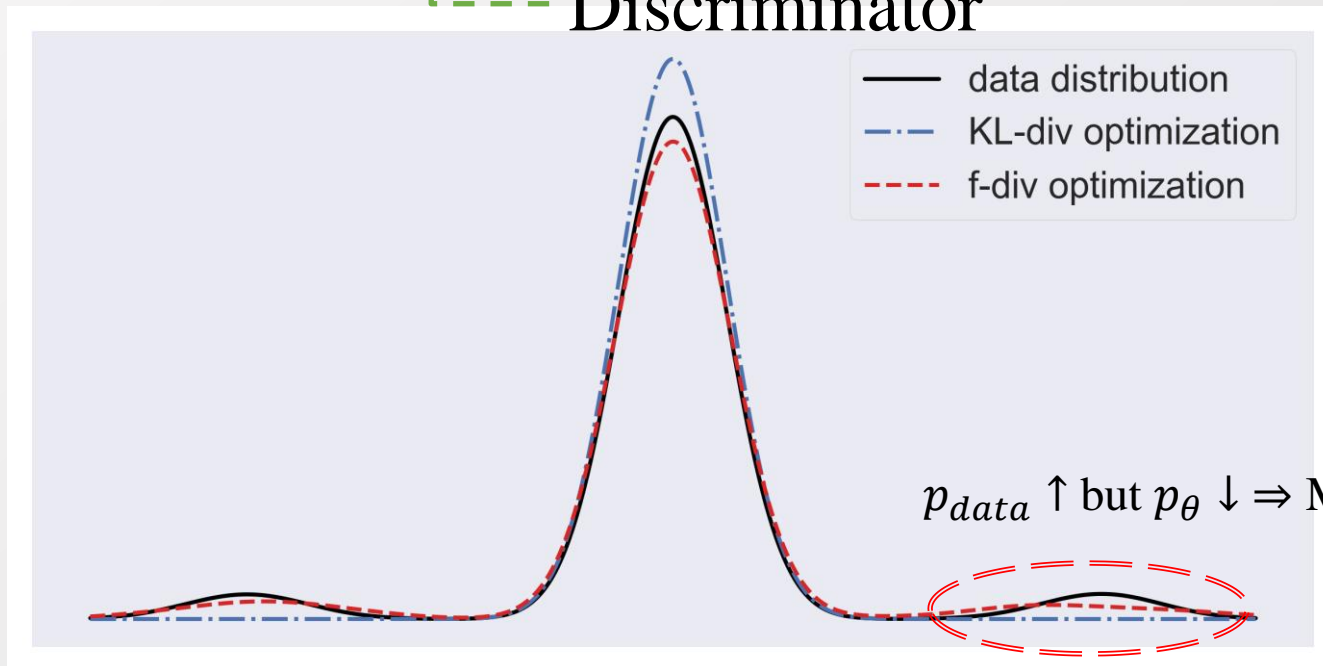
Disc. Guidance

$$\nabla \log p_{model}^t(\mathbf{x}_t|y) + w \nabla \log \frac{d^*(\mathbf{x}_t, t|y)}{1 - d^*(\mathbf{x}_t, t|y)}$$
$$= \nabla \log [(p_{data}^t(\mathbf{x}_t|y))^w (p_{model}^t(\mathbf{x}_t|y))^{1-w}]$$



$$\begin{aligned}
 & D_f(p_{data}(\mathbf{x}_0) \| p_{\theta}(\mathbf{x}_0)) \\
 &= \frac{1}{2} \int_0^T g^2(t) \mathbb{E} \left[f'' \left(\frac{p_{data}^t(\mathbf{x}_t)}{p_{\theta}^t(\mathbf{x}_t)} \right) \frac{p_{data}^t(\mathbf{x}_t)}{p_{\theta}^t(\mathbf{x}_t)} \left\| \nabla \log p_{data}^t(\mathbf{x}_t) - \mathbf{s}_{\theta}(\mathbf{x}_t, t) \right\|_2^2 \right] dt \\
 &\approx \frac{1}{2} \int_0^T g^2(t) \mathbb{E} \left[f'' \left(\frac{d_{\psi}(\mathbf{x}_t, t)}{1 - d_{\psi}(\mathbf{x}_t, t)} \right) \frac{d_{\psi}(\mathbf{x}_t, t)}{1 - d_{\psi}(\mathbf{x}_t, t)} \left\| \nabla \log p_{data}^t(\mathbf{x}_t) - \mathbf{s}_{\theta}(\mathbf{x}_t, t) \right\|_2^2 \right] dt
 \end{aligned}$$

- - - Score Mismatch Weight
- - - Score Mismatch
- - - Discriminator



Poster

Thu 1:30 – 3:00 #804

Thank You