

Your Diffusion Model is Secretly a Noise Classifier and Benefits from Contrastive Training

Presenter: Yunshu Wu

Yunshu Wu¹, Yingtao Luo², Xianghao Kong¹, Evangelos E. Papalexakis¹, Greg Ver Steeg¹

¹ University of California Riverside

² Carnegie Mellon University

Probability = (Optimal) Denoising

Guo, Shamai, Verdú, (IEEE 2005) – Information-MMSE relations

Kong, Brekelmans, Ver Steeg (ICLR 2023) "Information-theoretic Diffusion"

Gaussian noise channel $\mathbf{x}_\alpha \equiv \sqrt{\sigma(\alpha)}\mathbf{x} + \sqrt{\sigma(-\alpha)}\boldsymbol{\epsilon}$
 log-SNR $\mathbf{N}(\mathbf{0}, \mathbf{I})$



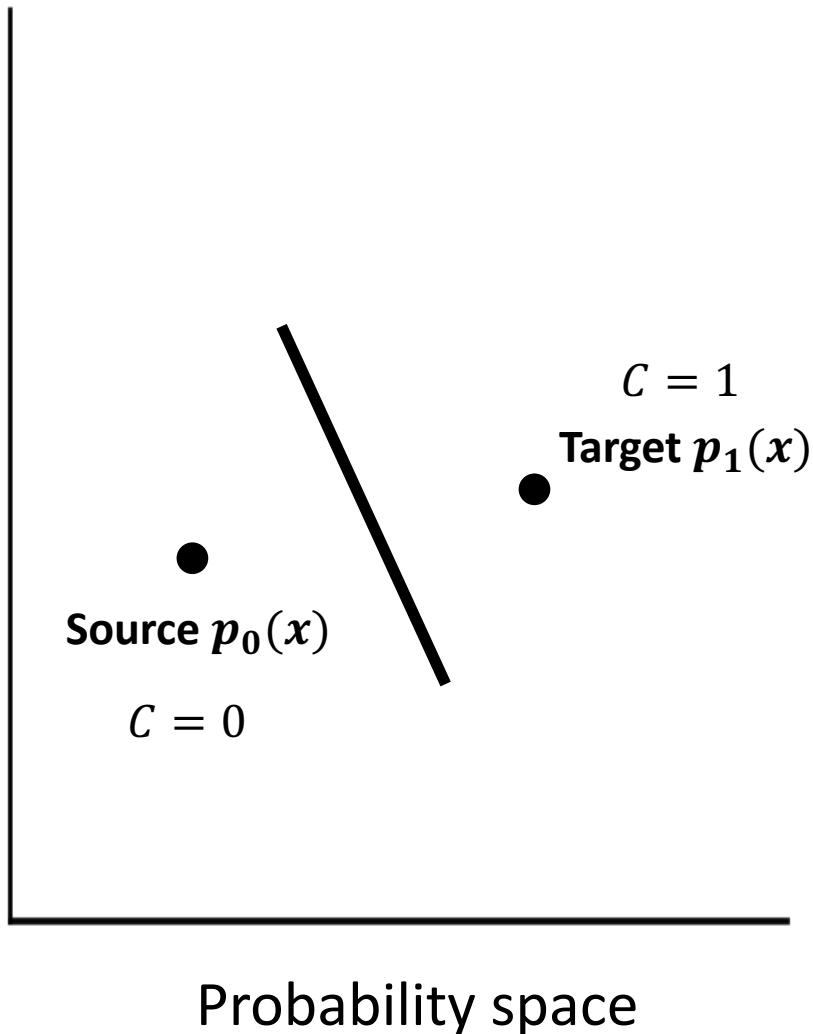
\mathbf{x}

\mathbf{x}_α

Negative Log-likelihood (NLL):

- NLL for data distribution:
$$-\log p(\mathbf{x}) = c + 1/2 \int_{-\infty}^{\infty} \mathbb{E}_{p(\boldsymbol{\epsilon})} [\|\boldsymbol{\epsilon} - \hat{\boldsymbol{\epsilon}}(\mathbf{x}_\alpha, \alpha)\|_2^2] d\alpha.$$
- NLL for noisy data distribution:
$$-\log p_\zeta(\mathbf{x}) = c + 1/2 \int_{-\infty}^{\infty} d\alpha \mathbb{E}_{p(\boldsymbol{\epsilon})} [\|\boldsymbol{\epsilon} - \hat{\boldsymbol{\epsilon}}_\zeta(\mathbf{x}_\alpha, \alpha)\|_2^2]$$
- $$-\log p_\zeta(\mathbf{x}) = c + 1/2 \int_{-\infty}^{\infty} d\alpha \mathbb{E}_{p(\boldsymbol{\epsilon})} [\|\boldsymbol{\epsilon} - b \cdot \hat{\boldsymbol{\epsilon}}(\mathbf{x}_\alpha, \beta)\|_2^2]$$

Log likelihood ratios and classification



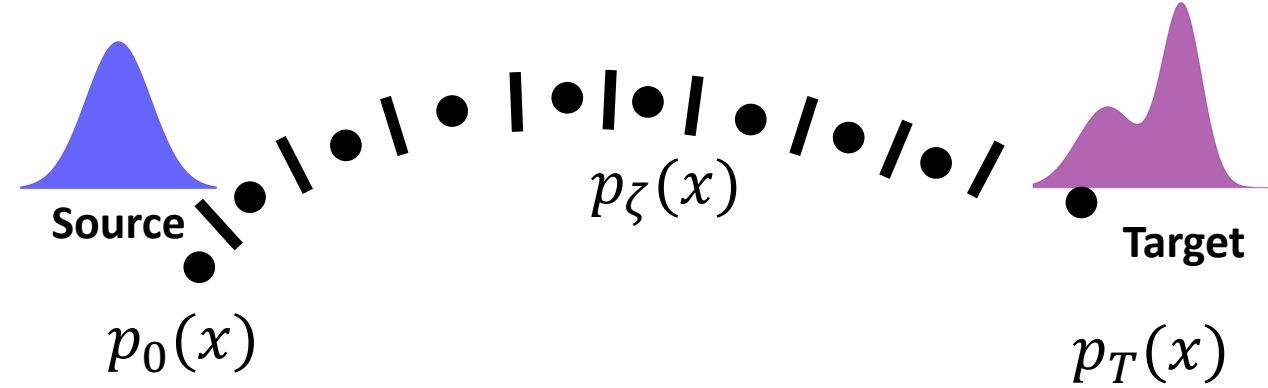
If we sample from either distribution with probability $P(C)=1/2$, then,

Bayes rule relates the log likelihood ratio to the optimal classifier:

$$LLR(x) \equiv \log \frac{p_1(x)}{p_0(x)} = \log \frac{p(C = 1 | x)}{p(C = 0 | x)}$$

Idea: train classifiers to learn densities

Denoising = Classifying



Sample from either distribution with probability $q(y = \pm 1) = \frac{1}{2}$, then, the optimal binary Bayes classifier is related to the log-likelihood ratio:

$$q(y|\mathbf{x}) = \frac{q(\mathbf{x}|y)q(y)}{q(\mathbf{x})} \quad \log q(y|\mathbf{x}) = -\text{softplus}\left(y \log \frac{q(\mathbf{x}|y = -1)}{q(\mathbf{x}|y = 1)}\right)$$

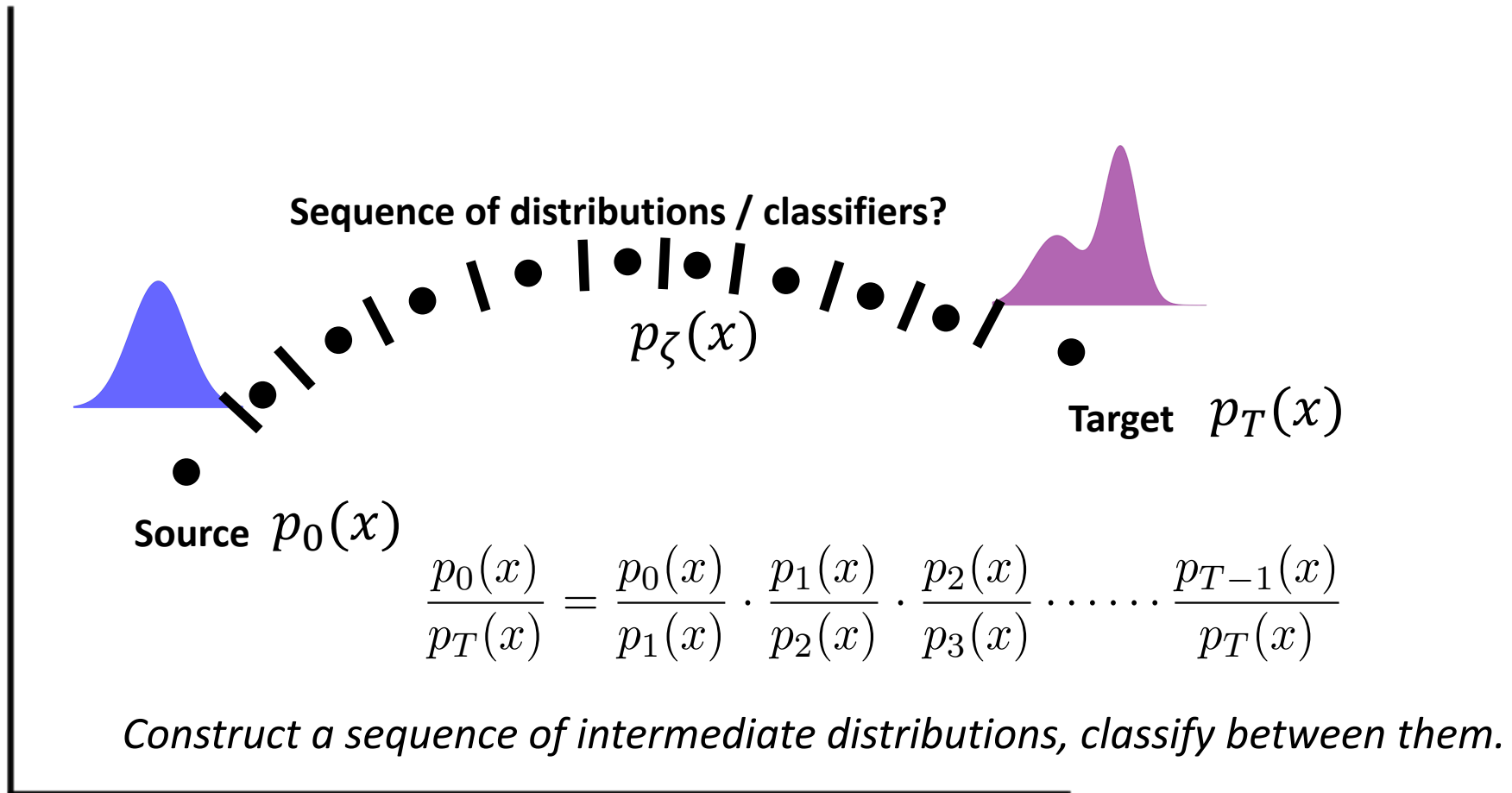
Let $q(\mathbf{x}|y=1) = p(\mathbf{x})$ be data distribution, $q(\mathbf{x}|y = -1) = p_\zeta(\mathbf{x})$ be noisy data distribution

LLR \rightarrow Optimal classifier: $LLR = \log p_\zeta(\mathbf{x}) - \log p(\mathbf{x})$

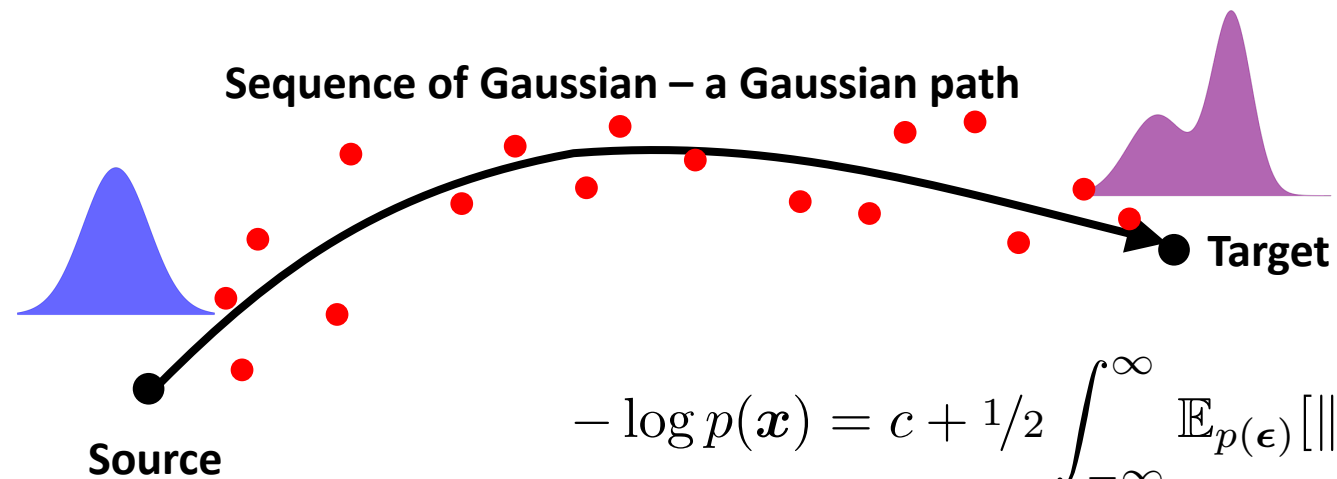
We have our Contrastive Diffusion Loss:

$$\mathcal{L}_{CDL} = \mathbb{E}_{q(\mathbf{x}, y)} [\text{softplus}(y(\log p_\zeta(\mathbf{x}) - \log p(\mathbf{x})))]$$

CDL – more general transport schemes



Profit : CDL training on OOD regions

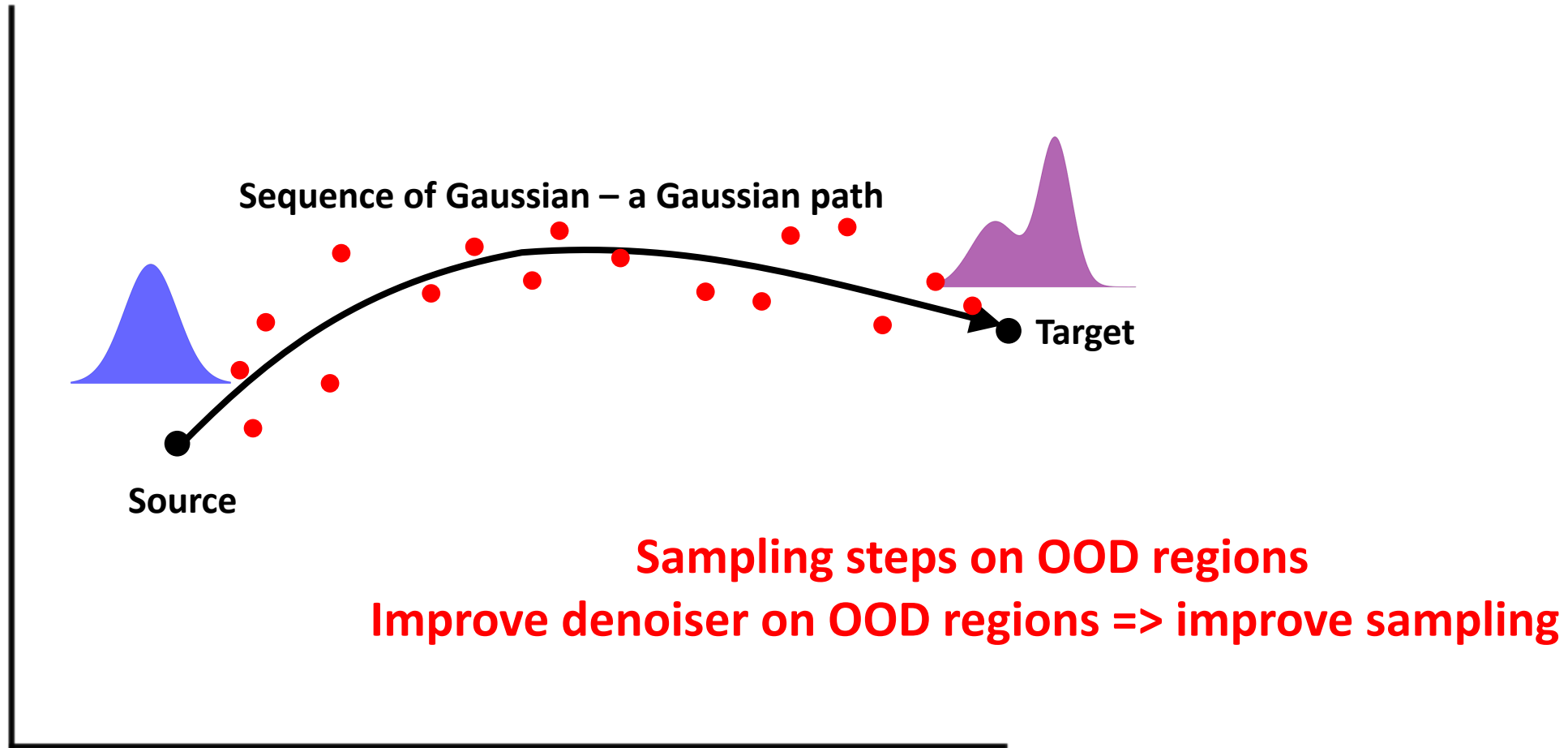


$$-\log p(\mathbf{x}) = c + 1/2 \int_{-\infty}^{\infty} \mathbb{E}_{p(\epsilon)} [\|\epsilon - \hat{\epsilon}(\mathbf{x}_\alpha, \alpha)\|_2^2] d\alpha.$$

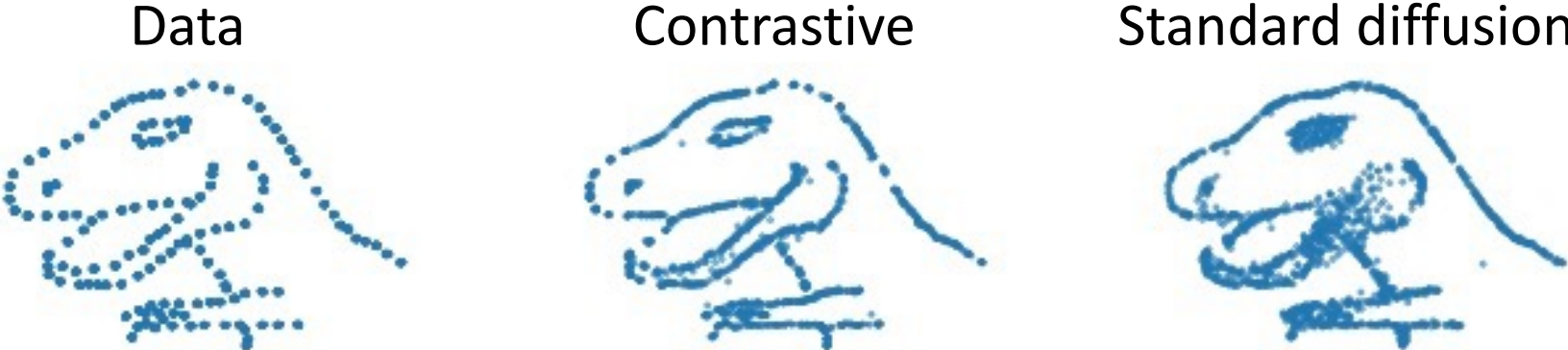
$$-\log p_\zeta(\mathbf{x}) = c + 1/2 \int_{-\infty}^{\infty} d\alpha \mathbb{E}_{p(\epsilon)} [\|\epsilon - b \cdot \hat{\epsilon}(\mathbf{x}_\alpha, \beta)\|_2^2]$$

Asynchronous input pairs = OOD regions for std diffusion loss

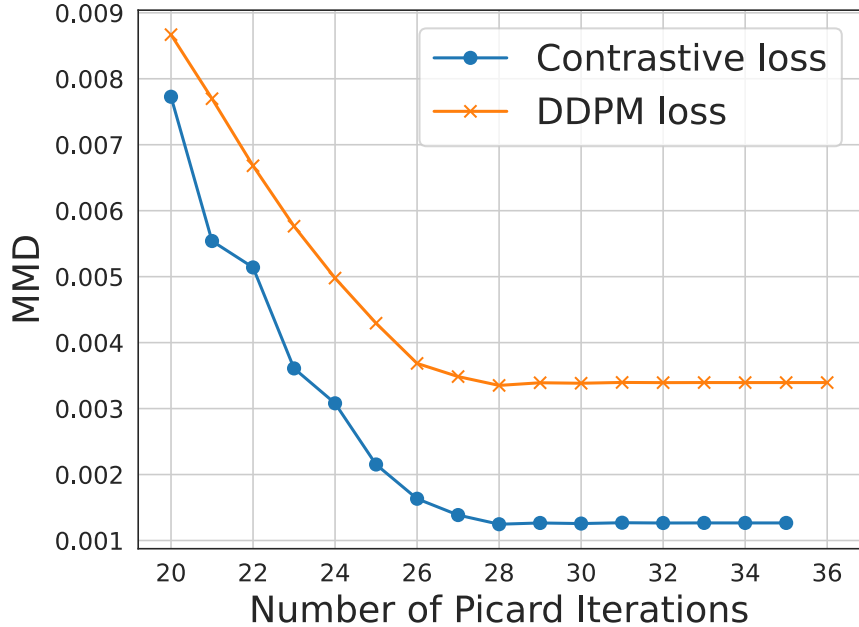
Profit : CDL training on OOD regions



Better distribution learning with hard constraints



Distribution similarity



Compute

Contrastive training improves “FID”

Models	CIFAR-10 at 32x32		AFHQv2 64x64	FFHQ 64x64
	unconditional	conditional	unconditional	unconditional
DDPM	9.43	NA	NA	NA
CDL-DDPM	9.06	NA	NA	NA
VP	3.24 ± 0.02	2.93 ± 0.02	2.95 ± 0.03	3.67 ± 0.04
CDL-VP	2.51 ± 0.01	2.41 ± 0.01	2.91 ± 0.02	3.33 ± 0.03
VE	3.00 ± 0.01	2.76 ± 0.01	2.98 ± 0.03	3.65 ± 0.02
CDL-VE	2.38 ± 0.01	2.25 ± 0.02	2.93 ± 0.01	3.29 ± 0.02

Table 2: Evaluating FID score (lower is better) of parallel DDPM sampler on real-world datasets using 5,000 samples. “NA” stands for “Not Applicable”. For reported FID scores, we run three sets of random seeds and reported the average with uncertainty.



Thank you!