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

Yunshu Wu

University of California Riverside
ywu380@ucr.edu

Yingtao Luo

Carnegie Mellon University
yingtao1@andrew.cmu.edu

Xianghao Kong

University of California Riverside
xkong016@ucr.edu

Evangelos E. Papalexakis

University of California Riverside
epapalex@cs.ucr.edu

Greg Ver Steeg

University of California Riverside
greg.versteeg@ucr.edu

Abstract

Diffusion models learn to denoise data and the trained denoiser is then used to generate new samples from the data distribution. In this paper, we revisit the diffusion sampling process and identify a fundamental cause of sample quality degradation: the denoiser is poorly estimated in regions that are far Outside Of the training Distribution (OOD), and the sampling process inevitably evaluates in these OOD regions. This can become problematic for all sampling methods, especially when we move to *parallel sampling* which requires us to initialize and update the entire sample trajectory of dynamics in parallel, leading to many OOD evaluations. To address this problem, we introduce a new self-supervised training objective that differentiates the levels of noise added to a sample, leading to improved OOD denoising performance. The approach is based on our observation that diffusion models implicitly define a log-likelihood ratio that distinguishes distributions with different amounts of noise, and this expression depends on denoiser performance outside the standard training distribution. We show by diverse experiments that the proposed contrastive diffusion training is effective for both sequential and parallel settings, and it improves the performance and speed of parallel samplers significantly.¹

1 Introduction

Denosing diffusion models [28] achieve state-of-the-art performance on various unsupervised learning tasks and have intriguing theoretical connections to methods like denoising autoencoders [35], VAEs [6], stochastic differential equations [19, 33], information theory [13], and score matching [30, 31]. Diffusion models are presented with data samples corrupted by a *forward* dynamical process that progressively adds more Gaussian noise and trained to *reverse* this dynamics or denoise the corrupted samples. Samples are then generated by applying the reverse dynamics on images of pure Gaussian noise to produce high-quality samples from the target distribution.

The key to the success of diffusion models is the dynamics that gradually bridges the source and a target distribution, but it suffers from slow sampling, as sequentially simulating these dynamics

¹Code can be found at <https://github.com/yunshuwu/ContrastiveDiffusionLoss.git>

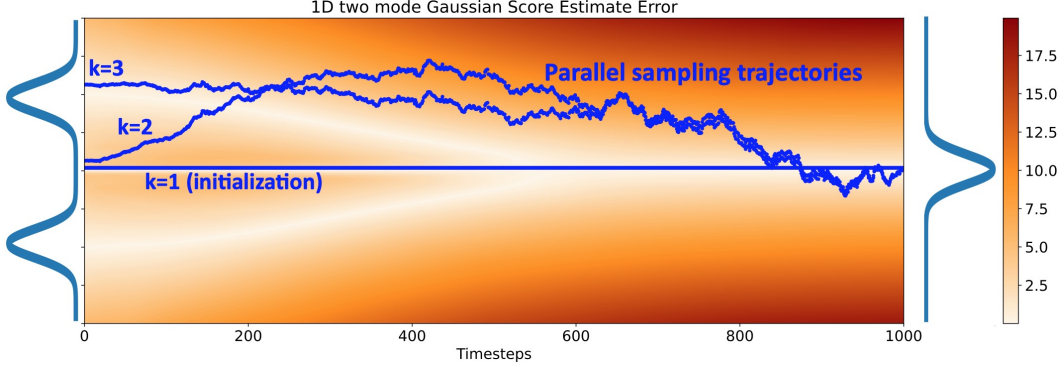


Figure 1: We plot the error in the score estimate for an 1D two mode Gaussian example where diffusion dynamics bridge between a Gaussian and a mixture (see Appendix A.3). Regions near the standard forward training data paths have lower error magnitude (light), whereas other areas have higher error magnitude (dark). While sequential samplers adhere as closely as possible to low-error regions, parallel samplers initialize and update the entire sample trajectory (blue trajectories), leading to evaluations in high-error regions. When the sampling trajectory is initialized, most are inevitably in the OOD regions and will update to the low-error regions gradually.

can take thousands of denoising steps for one sample. Most recent works attempt to expedite the sequential dynamics by taking fewer, larger steps [29, 11, 17, 21]. However, the complexity of these samplers and the need for expensive sampling hyper-parameter grid searches tailored to specific datasets makes them difficult to generalize.

Shih et al. 2024 suggests a different approach by randomly initializing the entire path of the reverse dynamics and then updating all the steps in the path in parallel. The parallel sampling approach promises to drastically reduce wall-clock time at the cost of increased parallel computation. However, it encounters a problem that has largely gone unnoticed in the sequential sampling literature: while sequential paths sampled during generation are designed carefully to stay as close as possible to the forward paths that add noise to data, the parallel sampler often evaluates in regions far from where the score estimate (the denoiser) is trained, as illustrated in Fig. 1. The shading shows that the error of the score estimate is large in these regions, leading to poor performance for parallel samplers. Discretization errors can lead to similar issues, even for standard sequential samplers.

We propose to improve the training of diffusion models so that the error of the denoiser is reduced in OOD regions, and we hypothesize that this should significantly improve the performance for parallel samplers as they require more OOD evaluations. Our approach starts with an unexpected connection: the optimal MSE denoiser that defines the diffusion dynamics *also defines an optimal noise classifier* that distinguishes between samples with different amounts of noise. This provides a useful additional signal for training, because optimizing for the noise classification task involves evaluating the denoiser for one noise level on samples from distributions at different noise levels, while standard MSE optimization only evaluates the denoiser on samples from the matching noisy distribution. Accurate denoiser evaluation in regions that are OOD for standard diffusion training is important for robust sampling dynamics.

Contributions:

- We use the information-theoretic formulation of diffusion to draw connections between diffusion, log-likelihood ratio estimation, and classification. This reveals that optimal diffusion denoisers are also implicitly optimal classifiers for predicting the amount of noise added to an image.
- We leverage the noise classifier interpretation to introduce a novel self-supervised loss function for regularizing diffusion model training, which we call the Contrastive Diffusion Loss (CDL). CDL provides training signal in regions that are OOD for the standard MSE diffusion loss.
- We show that CDL improves the trade-off between generation speed and sample quality, and that this advantage is consistent across different models, hyper-parameters, and sampling schemes. The improvement is especially substantial for parallel diffusion samplers [27] which rely heavily on OOD denoiser evaluations.

2 Diffusion Model Background: Optimal Denoisers are Density Estimators

The defining feature of diffusion models is a sequence of distributions that progressively add noise to the data, from which we then learn to recover the original data. The (“variance preserving” [33]) channel that mixes the signal \mathbf{x} with Gaussian noise is defined as $\mathbf{x}_\alpha \equiv \sqrt{\sigma(\alpha)}\mathbf{x} + \sqrt{\sigma(-\alpha)}\epsilon$ with $\epsilon \sim \mathcal{N}(0, \mathbb{I})$, $\mathbf{x} \sim p(\mathbf{x})$, where α represents the log of the Signal-to-Noise Ratio (SNR), $p(\mathbf{x})$ is the unknown data distribution for $\mathbf{x} \in \mathbb{R}^d$, and $\sigma(\cdot)$ is the sigmoid function. We define the sequence of intermediate distributions drawn according to this channel with a subscript as $p_\alpha(\mathbf{x})$. By definition, we express $\lim_{\alpha \rightarrow \infty} p_\alpha(\mathbf{x}) = p(\mathbf{x})$ in this paper. Note that we use a different scaling convention for noise from [11] and [6], where the former one takes $\mathbf{x} + \sigma\epsilon$ as the forward noising channel and the latter one takes $\sqrt{\alpha_t}\mathbf{x} + \sqrt{1 - \alpha_t}\epsilon$ as the forward noising channel. For further detailed relationships among these scaling conventions, please check App. B.3.

The minimum mean square error (MMSE) estimator $\hat{\epsilon}$ for recovering ϵ from the noisy channel that mixes \mathbf{x} and ϵ can be derived via variational calculus and written as follows.

$$\hat{\epsilon}(\mathbf{x}_\alpha, \alpha) \equiv \mathbb{E}_{\epsilon \sim p(\epsilon|\mathbf{x}_\alpha)}[\epsilon] = \arg \min_{\tilde{\epsilon}(\cdot, \cdot)} \mathbb{E}_{p(\epsilon)p(\mathbf{x})}[\|\epsilon - \tilde{\epsilon}(\mathbf{x}_\alpha, \alpha)\|_2^2]. \quad (1)$$

Sampling from the true posterior is typically intractable, but by using a neural network to approximate the solution to the regression optimization problem, we can get an approximation for $\hat{\epsilon}$. From [13], we see that log-likelihood can be written *exactly* in terms of an expression that depends only on the MMSE solution to the Gaussian denoising problem, i.e.

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

The constant, $c = d/2 \log(2\pi e) - \frac{d}{2} \int_0^\infty d\bar{\alpha} \sigma(\bar{\alpha})$ does not depend on data and will play no role in our approach, as it cancels out in our derivations in Sec. 3.

3 What Your Diffusion Model is Hiding: Noise Classifiers

We now introduce our first main result, which shows that diffusion models implicitly define optimal noise classifiers. Eq. (2) expresses the probability density of the data directly in terms of the denoising function. If we apply Eq. (2) to the noisy distributions that bridge the data and a Gaussian, $p_\zeta(\mathbf{x})$, we can see that all mixture densities can be written in terms of the same optimal denoising function, $\hat{\epsilon}(\cdot, \cdot)$. The complete derivation is presented in App. A.2.

$$-\log p_\zeta(\mathbf{x}) = c + 1/2 \int_{-\infty}^{\infty} d\alpha \mathbb{E}_{p(\epsilon)}[\|\epsilon - b \cdot \hat{\epsilon}(\mathbf{x}_\alpha, \beta)\|_2^2] \quad (3)$$

$$\mathbf{x}_\alpha \equiv \sqrt{\sigma(\alpha)}\mathbf{x} + \sqrt{\sigma(-\alpha)}\epsilon \quad (4)$$

$$\beta \equiv \sigma^{-1}(\sigma(\zeta)\sigma(\alpha)), \quad b \equiv \sqrt{\sigma(-\alpha)/\sigma(-\beta)} \quad (5)$$

Intuitively, if we find the optimal denoising function for the data distribution, it may be hypothesized that it can denoise an already *noisy* version of the data distribution. Using Eq. 2, this directly translates into an expression for density of mixture distributions. Differences in log likelihoods lead to cancellation of constants, and these Log Likelihood Ratios (LLR) are related to the optimal classifiers [4] as we show below.

To connect LLRs with classification, consider the following generative model. We generate a random binary label $q(y = \pm 1) = 1/2$. Then, conditioned on y , we sample from some distribution $q(\mathbf{x}|y)$. Given samples $(\mathbf{x}, y) \sim q(\mathbf{x}, y) = q(\mathbf{x}|y)q(y)$, the Bayes optimal classifier is:

$$\begin{aligned} q(y|\mathbf{x}) &= \frac{q(\mathbf{x}|y)q(y)}{q(\mathbf{x})} = \frac{q(\mathbf{x}|y)q(y)}{q(\mathbf{x}|y=1)q(y=1) + q(\mathbf{x}|y=-1)q(y=-1)} \\ &= 1/(1 + \frac{q(\mathbf{x}|y=-1)}{q(\mathbf{x}|y=1)}) = 1/(1 + \exp(y(\log q(\mathbf{x}|y=-1) - \log q(\mathbf{x}|y=1)))) \\ \log q(y|\mathbf{x}) &= -\log(1 + \exp(y \log \frac{q(\mathbf{x}|y=-1)}{q(\mathbf{x}|y=1)})) = -\text{softplus}(y \log \frac{q(\mathbf{x}|y=-1)}{q(\mathbf{x}|y=1)}) \end{aligned} \quad (6)$$

In the second line, because $\forall y, q(y) = 1/2$, these constants cancel out. Then we can just expand definitions and re-arrange to write in terms of log probabilities.

Contrastive Diffusion Loss (CDL) Our next contribution is to use the new connection between diffusion denoisers and noise classifiers to define a new training objective. We set the distributions $q(\mathbf{x}|y = 1)$ and $q(\mathbf{x}|y = -1)$ to be two distributions at different noise levels that we can write in terms of the optimal diffusion denoiser from Eq. 3. So we have $q(\mathbf{x}|y = 1) \equiv p(\mathbf{x})$, the data distribution, and $q(\mathbf{x}|y = -1) \equiv p_\zeta(\mathbf{x})$, for some noise level, ζ . Then given a sample $(\mathbf{x}, y) \sim q(\mathbf{x}, y)$ the per-sample cross-entropy loss for the noise classifier, Eq. (6), is as follows.

$$\mathcal{L}_{CDL} = \mathbb{E}_{q(\mathbf{x}, y)} [\text{softplus}(y(\log p_\zeta(\mathbf{x}) - \log p(\mathbf{x})))] \quad (7)$$

We can estimate both densities directly from our denoising model using Eq. (3), with the constants canceling out in the process. This loss differs significantly from the standard diffusion loss. Intuitively, to distinguish between a sample from the data distribution, $p(\mathbf{x})$, versus a noisy version of the data distribution, $p_\zeta(\mathbf{x})$, we need to evaluate denoisers on points from both distributions. In standard diffusion training, denoisers at noise level ζ are only trained on samples from $p_\zeta(\mathbf{x})$.

Limitations: We highlight that CDL is more expensive to compute than the standard diffusion loss, significantly increasing the total cost of diffusion model training. Implementation details appear in App. B.4 and training cost details appear in App. B.2.

Choice of noise to contrast Next, let’s break the Log-Likelihood Ratio (LLR) term in Eq. (7) down to see how to choose ζ to maximize the benefit of CDL. Combining Eq. (2) and Eq. (3) we have Eq. (8), where the constant cancels out.

$$LLR = \log p_\zeta(\mathbf{x}) - \log p(\mathbf{x}) = \int_{-\infty}^{\infty} d\alpha \mathbb{E}_{p(\epsilon)} [\|\epsilon - \hat{\epsilon}(\mathbf{z}_\alpha, \alpha)\|_2^2] - \mathbb{E}_{p(\epsilon)} [\|\epsilon - b\hat{\epsilon}(\mathbf{z}_\alpha, \beta)\|_2^2] \quad (8)$$

with: $\mathbf{z}_\alpha \equiv \sqrt{\sigma(\alpha)}\mathbf{x} + \sqrt{\sigma(-\alpha)}\epsilon$

From Eq. (8) we see that $\hat{\epsilon}(\cdot, \cdot)$ is trained on two pairs: $(\mathbf{z}_\alpha, \alpha)$ and $(\mathbf{z}_\alpha, \beta)$, where $\beta \equiv \sigma^{-1}(\sigma(\alpha)\sigma(\zeta)) < \min(\alpha, \zeta)$ (Eq. 5). During standard training, only the first pair is trained (Eq. 1). This means that our CDL objective trains the denoiser to perform correctly even for samples from distributions that are noisier than the specified noise level (a pair like $(\mathbf{z}_\alpha, \beta)$). This can be useful during sampling as extra noise added due to discretization errors can be corrected by the denoiser trained with CDL.

In practice, diffusion training pipelines are highly tuned on popular datasets like CIFAR10 and ImageNet, so the amplitude of discretization errors during sampling is small, meaning that errors won’t nudge points too far away from the true trajectory. Therefore, we choose a large value for ζ , corresponding to classifying only small differences in noise levels. Empirically $\zeta \sim \text{Uniform}[6, 15]$ or $\zeta \sim \text{logistic}[6, 15]$ performed best in our experiments.

Denoising, sampling dynamics, and the score connection We have focused so far on denoising and density estimation, but we now want to connect this discussion to the primary use case for diffusion models and the focus of Sec. 4, *sampling*. There are many choices in how to implement sampling dynamics [11], but all of them rely on the *score function*, $\nabla_x \log p_\alpha(\mathbf{x})$. The score function points toward regions of space with high likelihood, and by slowly transitioning (or annealing), from the score function of a noisy distribution to one closer to the data distribution, we can build reliable sampling dynamics. To connect denoisers with sampling we must show that a denoising function, $\hat{\epsilon}$, that is optimal according to Eq. 1 also specifies the score function.

$$\nabla_x \log p_\alpha(\mathbf{x}) = -\frac{\hat{\epsilon}(\mathbf{x}, \alpha)}{\sqrt{\sigma(-\alpha)}} \quad (9)$$

The derivation is straightforward and is given in Appendix A.1.

4 Sequential and Parallel Sampling with Diffusion Models

Sampling dynamics are typically presented in terms of a stochastic process $\{\mathbf{x}_t\}_{t=1}^T$ with timestep, t , rather than in terms of log SNR, α . We will denote $\mathbf{x}_t \equiv \mathbf{x}_{\alpha(t)}$, $p_t(\mathbf{x}) \equiv p_{\alpha(t)}(\mathbf{x})$, to connect to our previous notation, with $\alpha(t)$ representing a monotonic relationship described in App. B.3. Note

that decreasing log-SNR α corresponds to increase timestep t , since smaller log-SNR means there is more noise added to the data.

The general form of sampling dynamics is a process of slowly transitioning from samples of a simple and tractable distribution to the target distribution. Specifically, start with an isotropic Gaussian $\mathbf{x}_T \sim \mathcal{N}(0, \mathbb{I})$, the sampler steps through a series of intermediate distributions with noise levels $\{T, T-1, \dots, 1\}$ following the score estimates. Many works [33, 11] interpret diffusion models as stochastic differential equations (SDEs). The forward process is in the form of

$$d\mathbf{x}_t = \underbrace{f(\mathbf{x}_t, t) dt}_{\text{drift } s} + \underbrace{g(t) d\mathbf{w}_t}_{\text{diffusion}}, \quad \mathbf{x}_0 \sim p(\mathbf{x}) \quad (10)$$

where \mathbf{w}_t is the standard Wiener process/Brownian motion, f and g are drift coefficient and diffusion coefficient of \mathbf{x}_t separately. The reverse process of Eq. 10 is then used to generate samples

$$d\mathbf{x}_t = \underbrace{(f(\mathbf{x}_t, t)\mathbf{x} - g^2(t)\nabla_{\mathbf{x}} \log p_t(\mathbf{x})) dt}_{\text{drift } s} + \underbrace{g(t) d\mathbf{w}_t}_{\text{diffusion}}, \quad \mathbf{x}_T \sim p(\mathbf{x}) \quad (11)$$

Depending on choices for f, g , we can get either a stochastic or ordinary (deterministic) differential equation. Either way, numerical differential equation solvers are used to approximate the true dynamics. The solver introduces discretization errors at each step, causing the trajectory to deviate into the OOD region where the score (or denoiser) is poorly estimated, further compounding the errors. More discretization steps reduce accumulated error and leads to better sample quality, at the expense of more sequential computation. As a result, a significant limitation of diffusion models is that they require many iterations to produce high quality samples.

Sequential Sampling The influential diffusion sampler DDPM [6] iterates over thousands of discretization steps in simulating the dynamics. Recently, many sequential sampling methods have been developed to take fewer and larger steps while introducing less error [29, 21, 11]. Specifically, Karras et al. 2022 studies the curvature shape of SDE/ODE trajectory and suggests a discretization technique where the resulting tangent of the solution trajectory always points towards the denoiser output. However, speeding up the sequential sampling sacrifices generation quality. Furthermore, the SOTA sequential samplers [17, 11, 29] require hyperparameter tuning and grid search on specific datasets, which poses challenges to the generalization of these samplers to other datasets.

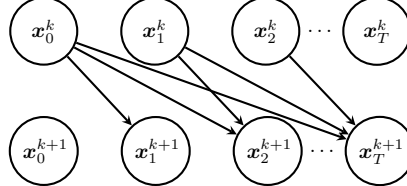


Figure 2: The computation graph of Picard iteration for parallel sampling [27]

Parallel Sampling Shih et al. 2024 explores a parallel sampling scheme, where the entire reverse process path is randomly initialized and then all steps in the path are updated in parallel. Parallel sampling is based on the method of Picard iteration, an old technique for solving ODEs through fixed-point iteration. An ODE is defined by a drift function $s(\mathbf{x}, t)$ and initial value \mathbf{x}_0 . In the integral form, the value at time t can be written as

$$\mathbf{x}_t = \mathbf{x}_0 + \int_0^t s(\mathbf{x}_u, u) du$$

In other words, the value at time t must be initial value plus the integral of the derivative along the path of the solution. This formula suggests a natural way of solving the ODE by starting with a guess of the solution $\{\mathbf{x}_t^{k+1} : 0 \leq t \leq 1\}$ at initial iteration $k = 0$, and iteratively refining by updating the value at every time t until convergence ²

(Picard Iteration)
$$\mathbf{x}_t^{k+1} = \mathbf{x}_0^k + \int_0^t s(\mathbf{x}_u^k, u) du \quad (12)$$

²For detailed convergence proof, we refer to Shih et al. 2024 section 3.

To perform Picard iterations numerically, which is shown in Fig. 2, we can write the discretized form of Eq. 12 with step size $1/T$, for $t \in [0, T]$:

$$\mathbf{x}_t^{k+1} = \mathbf{x}_0^k + \frac{1}{T} \sum_{i=0}^{t-1} s(\mathbf{x}_i^k, i/T) \quad (13)$$

We see that the expensive computations $\{s(\mathbf{x}_i^k, i/T) : i \in [0, T)\}$ can be performed in parallel. After some number of Picard iterations, the error difference between two iterates $\|\mathbf{x}_t^{k+1} - \mathbf{x}_t^k\|_2^2$ drops below some convergence threshold. This converged trajectory, \mathbf{x}_t^* , should be close to the sequential sampler trajectory. Looking at the example in Fig. 1, we show the trajectories of three iterations $k = 1, 2, 3$. The trajectories before convergence are consistently appearing in the regions with high score error.

5 Experiments

We sample from models fine-tuned on Contrastive Diffusion Loss (CDL) via both parallel and sequential diffusion samplers across a variety of generation tasks, including complex low-dimensional manifold 2D synthetic data and real-world image generation. Our results demonstrate that employing CDL as a regularizer in models trained with standard diffusion loss enhances density estimation and sample quality while also accelerating convergence in parallel sampling. All sampling tests are done on A6000 GPUs. We visualize CDL image generation examples in App. C.

Training configuration Our method is architecture-agnostic. In synthetic experiments, we adopt a simple MLP architecture with positional encoding for timesteps³, as it is one of the most versatile models in the literature. In real-world experiments, we consider the standard diffusion loss with two training configurations: DDPM by Ho et al. 2020 and EDM by Karras et al. 2022. For more details on model training, data split, and hyper-parameters, please refer to App. B.2.

Generation quality metrics For real-world data, our intrinsic metric is Fréchet Inception Distance (FID) score [5]. The number of images we generated for FID computation follows their baseline models’ FID settings. For parallel experiments, the FID scores are computed between 5,000 generated images and another 5,000 randomly selected real images; For sequential experiments, FID scores are computed between 50,000 generated images and all available real images.

For synthetic data, to measure how well the generated samples resemble samples from the ground truth distribution, we use the (unbiased) kernel estimator of the squared Maximum Mean Discrepancy (MMD), with Gaussian kernel with bandwidth set empirically as described in App. B.1.

Sampling speed metrics We adopt the following three metrics: (1) Neural function evaluations (NFE) for all settings, i.e. how many times the denoiser is evaluated to produce a sample; (2) For the parallel setting, we report the number of parallel Picard Iterations; (3) Furthermore for the parallel setting, the wall-clock time is reported. While parallel sampling can use fewer total iterations and less wall-clock time than a sequential sampler, this may come at the cost of an increase in the total number of function evaluations. This gap is called the algorithm inefficiency. In the subsequent section, we use contrastive diffusion loss as a training regularizer for standard diffusion losses and refer to the corresponding models as CDL-regularized models.

5.1 Parallel Sampling

In the parallel setting, we use Parallel DDPM sampler [27] with 1000-step diffusion sampling. And for synthetic experiment, to reflect sampling speed by only number of Picard iteration and wall-clock time, we don’t use the sliding window technique, and the 2D data is small enough to fit the whole sampler trajectory in GPU memory. While for real-world experiment, sliding window is still applied.

Synthetic Dataset We consider the 2D Dino dataset [18], characterized by its highly nonlinear density concentrated on a low-dimensional manifold. For baselines, we employ the standard DDPM loss [6], as all standard diffusion losses similarly minimize a sum of MSE losses between the actual and estimated denoisers. Both CDL-regularized and DDPM-objective-trained models are trained with a MLP where timestep is encoded by positional encoding. We train it for 2000 epochs to ensure convergence and check the training and validation loss curve to avoid overfitting.

³Architecture adopted from: <https://github.com/Jmkernes/Diffusion>

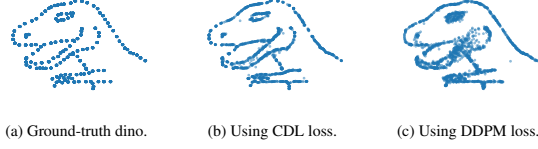


Figure 4: Parallel DDPM sampler generated Dino data. Comparing to Dino sampled from DDPM loss, CDL-loss sampled Dino has better sample quality and density estimate around hard areas.

Model Loss	MMD	#Iter	NFE	Time (ns)
DDPM	0.0031	36	14,397	1,870
CDL	0.0012	27	13,983	1,368

Table 1: Parallel DDPM sampling speed results. We generate 2,000 samples. Here we set MMD threshold= 0.002, and #Iter refers to number of picard iterations till MMD threshold. Both NFE and Time are counted till parallel convergence.

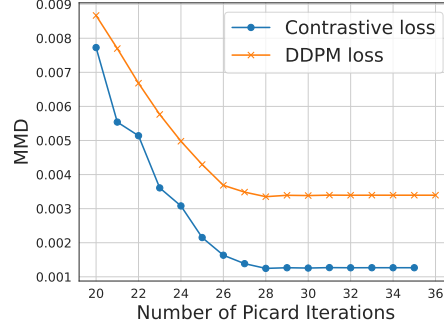


Figure 5: Parallel sampling MMD plot. We can see that the CDL-regularized model and DDPM model converge themselves at 35 and 36 Picard iterations, separately.

In Fig. 4, it is clearly demonstrated that the parallel generated samples from the CDL-regularized model is much better than the model trained only with the standard DDPM loss, especially around the chin, eyes and paws, where the manifolds are close to each other and difficult to distinguish and learn. From the MMD plot, we see that comparing to the baseline curve trained only with standard DDPM loss, the CDL-regularized curve converges faster with smaller number of Picard iterations and better sample quality (lower MMD scores). Table. 1 shows the sampling speed results, from where we see that CDL-regularized model converges faster with lower final MMD and better sample quality.

Real-world Datasets We select "DDPM" model by [6] trained on unconditional CIFAR-10 at 32×32 , also "DDPM++ cont. (VP)" and "NCSN++ cont. (VE)" models by [11] trained on CIFAR-10, unconditional FFHQ, and unconditional AFHQv2 [14, 10, 3] as baselines, comparing to the corresponding CDL-regularized models. We adopt the pre-trained models from Ho et al. 2020⁴ and Karras et al. 2022⁵. As shown in Tab. 2, CDL-regularized models always outperformed baselines with respect to FID scores.

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	7.93 ± 0.07	7.67 ± 0.07	4.58 ± 0.07	6.26 ± 0.07
CDL-VP	7.47 ± 0.07	7.27 ± 0.07	4.51 ± 0.04	5.89 ± 0.07
VE	7.81 ± 0.07	7.59 ± 0.07	4.65 ± 0.10	6.33 ± 0.07
CDL-VE	7.35 ± 0.07	7.19 ± 0.07	4.54 ± 0.07	5.94 ± 0.07

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.

5.2 Sequential Sampling

While CDL clearly improves parallel sampling quality and convergence speed, we also show that it improves the trade-off between generation speed and sample quality in the sequential setting. As for sequential diffusion sampling choices, we consider the DDPM sampler from Ho et al. 2020, and both the deterministic and stochastic samplers from Karras et al. 2022. To ensure fair comparisons, we adopt the original sampling hyper-parameter settings for all baselines.

Deterministic samplers For FID test, we follow the exact sampling settings outlined in Karras et al. 2022 for each dataset. FID scores are reported in Tab. 3, for sequential deterministic EDM samplers, CDL objective ensures that the generation quality is consistently similar or better.

⁴https://github.com/pesser/pytorch_diffusion

⁵<https://github.com/NVlabs/edm>

Models	CIFAR-10 at 32x32		AFHQv2 64x64	FFHQ 64x64
	unconditional	conditional	unconditional	unconditional
VP	2.00 ± 0.02	1.84 ± 0.02	2.04 ± 0.00	2.38 ± 0.01
CDL-VP	1.99 ± 0.04	1.82 ± 0.03	2.00 ± 0.00	2.29 ± 0.02
VE	2.01 ± 0.01	1.81 ± 0.01	2.17 ± 0.00	2.56 ± 0.03
CDL-VE	2.01 ± 0.01	1.81 ± 0.01	2.11 ± 0.01	2.47 ± 0.02
NFE (EDM/CDL)	35	35	79	79

Table 3: Evaluating sequential deterministic EDM samplers generation quality. For reported FID scores, we run three sets of random seeds and reported the average with uncertainty.

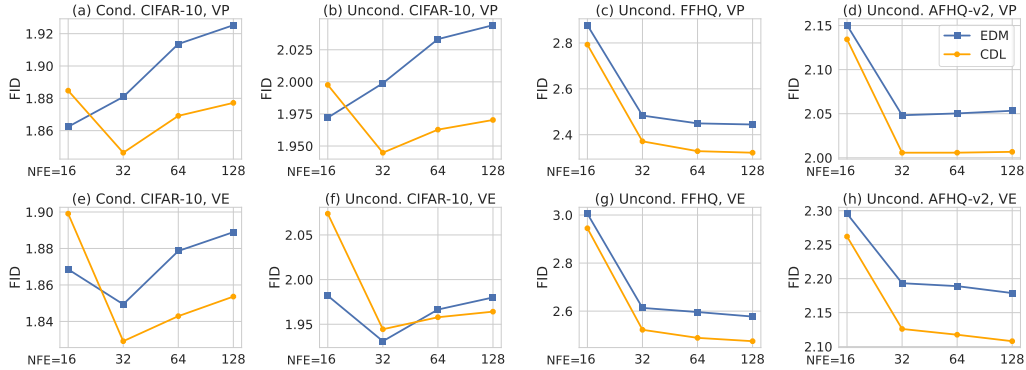


Figure 6: The FID comparison between our CDL and the baseline EDM in the deterministic sampler experiment.

In principle, increasing NFE has the potential to decrease the overall discretization errors, consequently leading to improved sample quality. However, in practice we observed an unusual behavior⁶ with the Karras deterministic sampler – as NFE increases, the FID score deteriorates (Fig. 6). In contrast to EDM models, CDL-regularized models exhibit a more stable FID score. This partially resolves the deterministic sampler sensitivity while improving the quality.

Stochastic samplers In practice, stochastic samplers often yield superior performance compared to deterministic ones. However this is not true in Karras et al. 2022: stochastic samplers outperform deterministic ones only on challenging datasets, for simpler datasets, the introduction of stochasticity not only fails to enhance performance but exhibits image degradation issues, characterized by a loss of detail. They attribute this phenomenon to L2-trained denoisers excessively removing noise at each step, and propose to slightly increase the standard deviation (S_{noise})⁷ of newly added noise to 1.007. We argue that this approach may not resolve the issue and instead complicates the hyperparameter grid search process by introducing an additional parameter, S_{noise} .

In this experiment, we conducted two stochastic sampling configurations for our baseline EDM-trained models. The first configuration, referred to as EDM-opt, operated at the EDM optimal setting with $S_{\text{noise}} = 1.007$. The second, named as EDM-sub-opt, used a setting with $S_{\text{noise}} = 1.00$, effectively disabling S_{noise} . As for CDL configuration, we exclusively examined the scenario with $S_{\text{noise}} = 1.00$ to determine whether CDL could address the problem of excessive noise removal.

The results, as visualized in Figure 7, indicate that CDL outperforms EDM in both S_{noise} configurations. Notably, CDL not only improves upon the EDM-sub-opt configuration (dark blue line) but also surpasses the performance of the EDM-opt configuration (light blue line), even at its optimal setting. This not only demonstrates that CDL robustly provides a better sample quality, but also suggests that CDL can eliminate the need for the hyperparameter S_{noise} . This reduction enables a more efficient grid search for the optimal EDM sampling settings, potentially enhancing the practicality of using such a sampler for other applications.

⁶The same issue is also reported in <https://github.com/NVlabs/edm/issues/4>

⁷We refer to Karras et al. [11] for details

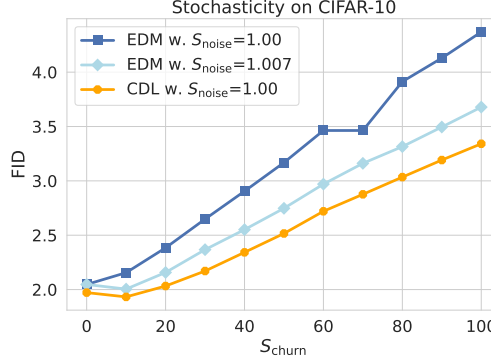


Figure 7: The FID comparison between our CDL and the baselines EDM in the stochastic sampler experiment on CIFAR-10. CDL’s performance is strictly better for all S_{churn} , outperforming the optimal setting of EDM which inflates the standard deviation S_{churn} of the newly added noise.

6 Related Work

The generative modeling trilemma [37] seeks generative models that produce samples (i) quickly, (ii) with high quality, and (iii) good mode coverage. Diffusion models excel at the latter two but a large amount of research has attempted to address the problem of slow sampling speed. From the inception of diffusion models [28], dynamics has been at the forefront, so most work has focused on the interpretation of the dynamics as a differential equation that can be sped up through accelerated numerical solvers [29, 33, 39, 9, 16]. Our approach is compatible with any of these approaches as we are sampler agnostic, seeking only to improve the input to the sampler, which is the denoiser or score function estimator, through regularization during the diffusion model training. A separate line of work instead attempts to distill a diffusion model into a faster model that achieves similar sampling quality [26, 34, 20, 36, 37]. In principle, these methods could also benefit from distilling based on a more robust base diffusion model trained with CDL.

Diffusion models admit a surprisingly diverse array of mathematical perspectives, like variational perspectives [6, 7, 12], differential equations [33, 19], and nonequilibrium thermodynamics [28]. Our approach is mostly inspired by connections between the information-theoretic perspective [13] and the score matching perspective [30, 32, 31, 8, 35]. In particular, we point out that score function estimates in traditional diffusion training are sub-optimal, and the information-theoretic perspective leads to a new objective (CDL) that can improve the score estimate.

Learning based on contrastive losses builds on the seminal works [4, 22], with influential applications including SimCLR [1] and CLIP [24]. A few works have considered contrastive learning inspired modifications to diffusion either to enforce multimodal data relationships [15, 40], for style transfer [38], or for guidance during generation [23], but none use the diffusion model as a noise classifier to improve diffusion training as we do. Most similar to our approach are [25, 2] that both train a classifier to distinguish between samples from distributions on a path between a source and target distribution. Our contribution is to relate this approach to diffusion models by noting that standard diffusion models implicitly already implement the required classifiers for distinguishing distributions on the path from the data distribution to a Gaussian distribution.

7 Conclusion

In this paper, we introduced a novel connection between diffusion models and optimal noise classifiers. While this relationship has a variety of potential applications that could be explored in future work, we used the connection to propose a new self-supervised loss regularizer for diffusion models, the Contrastive Diffusion Loss (CDL). CDL reduces the error of the learned denoiser in regions that are OOD for the standard loss. We showed that CDL improves the robustness of diffusion models across all types of sampling dynamics, and leads to significant speed-ups for a promising new generation of parallel samplers.

Broader Impacts: This paper aims to innovate on the methodology of diffusion models. We anticipate no direct potential societal consequences of our work, as the main focus of this work is in theory and algorithm design. However, it is important to acknowledge the potential risks associated with diffusion models, as misuse can contribute to the spread of disinformation and deepfakes.

References

- [1] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.
- [2] Kristy Choi, Chenlin Meng, Yang Song, and Stefano Ermon. Density ratio estimation via infinitesimal classification. *arXiv preprint arXiv:2111.11010*, 2021.
- [3] Yunjei Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. Stargan v2: Diverse image synthesis for multiple domains. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8188–8197, 2020.
- [4] Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 297–304. JMLR Workshop and Conference Proceedings, 2010.
- [5] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.
- [6] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- [7] Chin-Wei Huang, Jae Hyun Lim, and Aaron C Courville. A variational perspective on diffusion-based generative models and score matching. *Advances in Neural Information Processing Systems*, 34:22863–22876, 2021.
- [8] Aapo Hyvärinen and Peter Dayan. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(4), 2005.
- [9] Alexia Jolicoeur-Martineau, Ke Li, Rémi Piché-Taillefer, Tal Kachman, and Ioannis Mitliagkas. Gotta go fast when generating data with score-based models. *arXiv preprint arXiv:2105.14080*, 2021.
- [10] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019.
- [11] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *arXiv preprint arXiv:2206.00364*, 2022.
- [12] Diederik P Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *arXiv preprint arXiv:2107.00630*, 2021.
- [13] Xianghao Kong, Rob Brekelmans, and Greg Ver Steeg. Information-theoretic diffusion. In *International Conference on Learning Representations*, 2023. URL <https://arxiv.org/abs/2302.03792>.
- [14] Alex Krizhevsky. Learning multiple layers of features from tiny images. pages 32–33, 2009. URL <https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf>.
- [15] Chaejeong Lee, Jayoung Kim, and Noseong Park. Codi: Co-evolving contrastive diffusion models for mixed-type tabular synthesis. *arXiv preprint arXiv:2304.12654*, 2023.
- [16] Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. Pseudo numerical methods for diffusion models on manifolds. *arXiv preprint arXiv:2202.09778*, 2022.
- [17] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models. *arXiv preprint arXiv:2211.01095*, 2022.
- [18] Justin Matejka and George Fitzmaurice. Same stats, different graphs: generating datasets with varied appearance and identical statistics through simulated annealing. In *Proceedings of the 2017 CHI conference on human factors in computing systems*, pages 1290–1294, 2017.

- [19] David McAllester. On the mathematics of diffusion models. *arXiv preprint arXiv:2301.11108*, 2023.
- [20] Chenlin Meng, Robin Rombach, Ruiqi Gao, Diederik Kingma, Stefano Ermon, Jonathan Ho, and Tim Salimans. On distillation of guided diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14297–14306, 2023.
- [21] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021.
- [22] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- [23] Yidong Ouyang, Liyan Xie, and Guang Cheng. Improving adversarial robustness by contrastive guided diffusion process. *arXiv preprint arXiv:2210.09643*, 2022.
- [24] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [25] Benjamin Rhodes, Kai Xu, and Michael U Gutmann. Telescoping density-ratio estimation. *Advances in neural information processing systems*, 33:4905–4916, 2020.
- [26] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. In *International Conference on Learning Representations*, 2021.
- [27] Andy Shih, Suneel Belkhale, Stefano Ermon, Dorsa Sadigh, and Nima Anari. Parallel sampling of diffusion models. *Advances in Neural Information Processing Systems*, 36, 2024.
- [28] Jascha Sohl-Dickstein, Eric A Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. *arXiv preprint arXiv:1503.03585*, 2015.
- [29] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.
- [30] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in Neural Information Processing Systems*, 32, 2019.
- [31] Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. *Advances in neural information processing systems*, 33:12438–12448, 2020.
- [32] Yang Song, Sahaj Garg, Jiaxin Shi, and Stefano Ermon. Sliced score matching: A scalable approach to density and score estimation. *arXiv preprint arXiv:1905.07088*, 2019.
- [33] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
- [34] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. *arXiv preprint arXiv:2303.01469*, 2023.
- [35] Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674, 2011.
- [36] Daniel Watson, William Chan, Jonathan Ho, and Mohammad Norouzi. Learning fast samplers for diffusion models by differentiating through sample quality.
- [37] Zhisheng Xiao, Karsten Kreis, and Arash Vahdat. Tackling the generative learning trilemma with denoising diffusion gans. In *International Conference on Learning Representations*, 2021.
- [38] Serin Yang, Hyunmin Hwang, and Jong Chul Ye. Zero-shot contrastive loss for text-guided diffusion image style transfer. *arXiv preprint arXiv:2303.08622*, 2023.
- [39] Qinsheng Zhang and Yongxin Chen. Diffusion normalizing flow. *Advances in Neural Information Processing Systems*, 34:16280–16291, 2021.
- [40] Ye Zhu, Yu Wu, Kyle Olszewski, Jian Ren, Sergey Tulyakov, and Yan Yan. Discrete contrastive diffusion for cross-modal and conditional generation. *arXiv preprint arXiv:2206.07771*, 2022.

A Derivations and Proofs

A.1 Score Connection

We derive the following relation.

$$\nabla_x \log p_\alpha(\mathbf{x}) = -\frac{\hat{\epsilon}(\mathbf{x}, \alpha)}{\sqrt{\sigma(-\alpha)}} \quad (14)$$

To keep track of intermediate random variables and their associated distributions, we will use the cumbersome but more precise information theory notation where capitals represent a random variable and lowercase represents values. Define the channel that mixes the signal, X , with Gaussian noise as $Z_\alpha \equiv \sqrt{\sigma(\alpha)}X + \sqrt{\sigma(-\alpha)}\mathcal{E}$ with $\mathcal{E} \sim \mathcal{N}(0, \mathbb{I})$ and data distribution $p(X)$, α represents the log of the Signal-to-Noise Ratio (SNR), and $\sigma(\cdot)$ is the sigmoid function. In this notation, the probability that a mixture distribution takes a value, \mathbf{x} , would be written $p(Z_\alpha = \mathbf{x})$.

We start by re-writing the left-hand side in new notation, and expand the definition based on the noisy channel model, using $\delta(\cdot)$ for the Dirac delta.

$$\begin{aligned} \nabla_x \log p(Z_\alpha = \mathbf{x}) &= 1/p(Z_\alpha = \mathbf{x}) \nabla_x p(Z_\alpha = \mathbf{x}) \\ &= 1/p(Z_\alpha = \mathbf{x}) \nabla_x \int d\bar{\mathbf{x}} d\epsilon \delta(\mathbf{x} - \sqrt{\sigma(\alpha)}\bar{\mathbf{x}} - \sqrt{\sigma(-\alpha)}\epsilon) p(X = \bar{\mathbf{x}}) p(\mathcal{E} = \epsilon) \\ &= 1/p(Z_\alpha = \mathbf{x}) \int d\bar{\mathbf{x}} d\epsilon (\nabla_x \delta(\mathbf{x} - \sqrt{\sigma(\alpha)}\bar{\mathbf{x}} - \sqrt{\sigma(-\alpha)}\epsilon)) p(X = \bar{\mathbf{x}}) p(\mathcal{E} = \epsilon) \\ &= 1/p(Z_\alpha = \mathbf{x}) \int d\bar{\mathbf{x}} d\epsilon ((-1/\sqrt{\sigma(-\alpha)}) \nabla_\epsilon \delta(\mathbf{x} - \sqrt{\sigma(\alpha)}\bar{\mathbf{x}} - \sqrt{\sigma(-\alpha)}\epsilon)) p(X = \bar{\mathbf{x}}) p(\mathcal{E} = \epsilon) \\ &= 1/p(Z_\alpha = \mathbf{x}) \frac{1}{\sqrt{\sigma(-\alpha)}} \int d\bar{\mathbf{x}} d\epsilon \delta(\mathbf{x} - \sqrt{\sigma(\alpha)}\bar{\mathbf{x}} - \sqrt{\sigma(-\alpha)}\epsilon) p(X = \bar{\mathbf{x}}) \nabla_\epsilon p(\mathcal{E} = \epsilon) \\ &= -1/p(Z_\alpha = \mathbf{x}) \frac{1}{\sqrt{\sigma(-\alpha)}} \int d\bar{\mathbf{x}} d\epsilon \delta(\mathbf{x} - \sqrt{\sigma(\alpha)}\bar{\mathbf{x}} - \sqrt{\sigma(-\alpha)}\epsilon) p(X = \bar{\mathbf{x}}) \epsilon p(\mathcal{E} = \epsilon) \\ &= -1/p(Z_\alpha = \mathbf{x}) \frac{1}{\sqrt{\sigma(-\alpha)}} \int d\epsilon p(Z_\alpha = \mathbf{x}, \mathcal{E} = \epsilon) \epsilon \\ &= -\frac{1}{\sqrt{\sigma(-\alpha)}} \int d\epsilon p(\mathcal{E} = \epsilon | Z_\alpha = \mathbf{x}) \epsilon \\ &= -\frac{\mathbb{E}_{\epsilon \sim p(\mathcal{E} | Z_\alpha = \mathbf{x})}[\epsilon]}{\sqrt{\sigma(-\alpha)}} \\ &= -\frac{\hat{\epsilon}(\mathbf{x}, \alpha)}{\sqrt{\sigma(-\alpha)}} \end{aligned}$$

In the second line we expand, and in the third we just move the gradient inside the integral. In the fourth line we use the chain rule to relate the gradient over \mathbf{x} to the gradient over ϵ (introducing a sign flip). In the fifth line we use integration by parts to move the gradient (second sign flip). Taking the gradient of the Gaussian in the sixth line gives our third sign flip, and the factor of ϵ . We can conclude by writing the expression in terms of a conditional distribution, and relating that to the optimal denoiser in Eq. 1.

A.2 Mixture Distribution Density

In this section, we derive the expression that shows that the density of a continuum of Gaussian mixture distributions can be written in terms of the optimal denoiser, $\hat{\epsilon}$, for the data distribution.

$$\begin{aligned} -\log p_\alpha(\mathbf{x}) &= c + 1/2 \int_{-\infty}^{\infty} d\bar{\alpha} \mathbb{E}_{p(\epsilon)} [\|\epsilon - \sqrt{\frac{\sigma(-\bar{\alpha})}{\sigma(-\beta)}} \hat{\epsilon}(\sqrt{\sigma(\bar{\alpha})}\mathbf{x} + \sqrt{\sigma(-\bar{\alpha})}\epsilon, \beta)\|_2^2] \\ \beta &\equiv \sigma^{-1}(\sigma(\alpha)\sigma(\bar{\alpha})) \end{aligned} \quad (15)$$

As in the previous section, we will adopt information theory notation. If we define the optimal denoiser for the input distribution, $p_\alpha(\mathbf{x})$, with a subscript as $\hat{\epsilon}_\alpha(\cdot, \cdot)$, we can write the density analogously to Eq. 2.

$$-\log p_\alpha(\mathbf{x}) = c + 1/2 \int_{-\infty}^{\infty} d\bar{\alpha} \mathbb{E}_{p(\epsilon)} [\|\epsilon - \hat{\epsilon}_\alpha(\sqrt{\sigma(\bar{\alpha})}\mathbf{x} + \sqrt{\sigma(-\bar{\alpha})}\epsilon, \bar{\alpha})\|_2^2] \quad (16)$$

Note that we now have to keep track of two log SNR values. One indicates how much noise is added to the new “data” distribution, the other is how much noise we add and then try to remove with our denoiser. The goal is to relate $\hat{\epsilon}_\alpha$ to $\hat{\epsilon}$. We can formally write down the optimal solution using the relation in Eq. 1.

$$\hat{\epsilon}_\alpha(\mathbf{x}, \bar{\alpha}) = \mathbb{E}_{\epsilon \sim p(\mathcal{E}|Z=\mathbf{x})} [\epsilon]$$

Now, however, the noise channel is defined differently. The channel mixes the signal, $\bar{\mathbf{x}} \sim p_\alpha(\bar{X})$, with Gaussian noise, $\bar{\epsilon} \sim \mathcal{N}(0, \mathbb{I})$, as $Z \equiv \sqrt{\sigma(\bar{\alpha})}\bar{X} + \sqrt{\sigma(-\bar{\alpha})}\bar{\mathcal{E}}$. And the noisy variable, $\bar{X} = \sqrt{\sigma(\alpha)}X + \sqrt{\sigma(-\alpha)}\mathcal{E}$, where we must be careful to distinguish the two independent sources of Gaussian noise.

We start by expanding definitions.

$$\begin{aligned} \hat{\epsilon}_\alpha(\mathbf{x}_\alpha, \bar{\alpha}) &= \frac{1}{p(Z = \mathbf{x}_\alpha)} \int d\bar{\epsilon} p(\mathcal{E} = \bar{\epsilon}, Z = \mathbf{x}_\alpha) \bar{\epsilon} \\ &= \frac{1}{p(Z = \mathbf{x}_\alpha)} \int d\epsilon d\bar{\epsilon} d\bar{\mathbf{x}} d\mathbf{x} \delta(\mathbf{x}_\alpha - \sqrt{\sigma(\bar{\alpha})}\bar{\mathbf{x}} - \sqrt{\sigma(-\bar{\alpha})}\bar{\epsilon}) \delta(\bar{\mathbf{x}} - \sqrt{\sigma(\alpha)}\mathbf{x} - \sqrt{\sigma(-\alpha)}\epsilon) \bar{\epsilon} \\ &= \frac{1}{p(Z = \mathbf{x}_\alpha)} \int d\epsilon d\bar{\epsilon} d\mathbf{x} \delta(\mathbf{x}_\alpha - \sqrt{\sigma(\bar{\alpha})}(\sqrt{\sigma(\alpha)}\mathbf{x} + \sqrt{\sigma(-\alpha)}\epsilon) - \sqrt{\sigma(-\bar{\alpha})}\bar{\epsilon}) \bar{\epsilon} \end{aligned}$$

Now do a change of variables, a 2-d rotation with:

$$\begin{aligned} \bar{\epsilon}' &= a\bar{\epsilon} + b\epsilon, \epsilon' = -b\bar{\epsilon} + a\epsilon, \\ a &= \sqrt{\sigma(\bar{\alpha})\sigma(\alpha)/(1 - \sigma(\alpha)\sigma(\bar{\alpha}))}, b = \sqrt{\sigma(-\bar{\alpha})/(1 - \sigma(\alpha)\sigma(\bar{\alpha}))}. \end{aligned}$$

This change of variables leads to the following.

$$\begin{aligned} \hat{\epsilon}_\alpha(\mathbf{x}_\alpha, \bar{\alpha}) &= \frac{1}{p(Z = \mathbf{x}_\alpha)} \int d\epsilon' d\bar{\epsilon}' d\mathbf{x} \delta(\mathbf{x}_\alpha - \sqrt{\sigma(\alpha)\sigma(\bar{\alpha})}\mathbf{x} - \sqrt{1 - \sigma(\alpha)\sigma(\bar{\alpha})}\bar{\epsilon}') (b\bar{\epsilon}' + a\epsilon') \\ &= b \frac{1}{p(Z = \mathbf{x}_\alpha)} \int d\epsilon' d\bar{\epsilon}' d\mathbf{x} \delta(\mathbf{x}_\alpha - \sqrt{\sigma(\alpha)\sigma(\bar{\alpha})}\mathbf{x} - \sqrt{1 - \sigma(\alpha)\sigma(\bar{\alpha})}\bar{\epsilon}') \bar{\epsilon}' \\ &= b \frac{1}{p(Z = \mathbf{x}_\alpha)} \int d\epsilon' d\bar{\epsilon}' d\mathbf{x} \delta(\mathbf{x}_\alpha - \sqrt{\sigma(\beta)}\mathbf{x} - \sqrt{1 - \sigma(\beta)}\bar{\epsilon}') \bar{\epsilon}' \\ \hat{\epsilon}_\alpha(\mathbf{x}_\alpha, \bar{\alpha}) &= b\hat{\epsilon}(\mathbf{x}_\alpha, \beta), \quad \beta \equiv \sigma^{-1}(\sigma(\bar{\alpha})\sigma(\alpha)), b = \sqrt{\sigma(-\bar{\alpha})/(1 - \sigma(\alpha)\sigma(\bar{\alpha}))} \end{aligned}$$

Note in the second line that the expectation of ϵ' is zero, and we move the constant for the other term, b , outside the integral. In the third line, we define β which represents the log SNR of the two consecutive noisy channels with $\alpha, \bar{\alpha}$. Then we recognize the resulting integral as Eq. 1, the optimal denoiser for recovering samples from the original (non-noisy) data distribution in Gaussian noise.

A.3 Main Plot 1D Two-mode Gaussian’s Analytical Solution

In this section, we calculate the analytical solution to the 1D two-mode Gaussian in Fig. 1.

Fig. 1 is plotting the norm of difference between the ground-truth denoiser $\hat{\epsilon}_{gt}(\cdot, \cdot)$ and the estimated denoiser $\hat{\epsilon}(\cdot, \cdot)$:

$$\text{denoiser_err}(\mathbf{x}, \alpha) = \|\hat{\epsilon}(\mathbf{x}, \alpha) - \hat{\epsilon}_{gt}(\mathbf{x}, \alpha)\|_2^2$$

To get this error plot, we need to analytically calculate ground-truth denoiser $\hat{\epsilon}_{gt}(\mathbf{x}, \alpha)$. From score connection Eq. 9, for any intermediate noisy density $\log p_\alpha(\mathbf{x})$, denoiser function $\hat{\epsilon}(\mathbf{x}, \alpha)$ can be derived from score function $\nabla_x \log p_\alpha(\mathbf{x})$. Therefore, ultimately what we need to calculate here is the score function of any noisy distribution $p_\alpha(\mathbf{x})$.

The data we used is a mixture of two Gaussians, $\mathcal{N}(\mu = -5, \mathbb{I})$ and $\mathcal{N}(\mu = 5, \mathbb{I})$, and the noise distribution is consists of data plus noise, then the noisy distribution should also be a mixture of Gaussians. We just need to relate the parameters of the noisy mixture of Gaussians to the parameters of the mixture of Gaussians.

Start with one mode of the Gaussian mixture for the data. We could represent it in terms of the standard normal random variable, ϵ .

$$\mathbf{x}_d = \mu + \sigma\epsilon$$

Now call \mathbf{x}_α the random variable after applying a noisy channel with log-SNR α , here we present sigmoid function as $\sigma(\cdot)$.

$$\mathbf{x}_\alpha = \sqrt{\sigma(\alpha)}\mathbf{x}_d + \sqrt{\sigma(-\alpha)}\epsilon'$$

Note that we use a different ϵ' here. Now expand this, and then re-arrange.

$$\begin{aligned}\mathbf{x}_\alpha &= \sqrt{\sigma(\alpha)}(\mu + \sigma\epsilon) + \sqrt{\sigma(-\alpha)}\epsilon' \\ &= \sqrt{\sigma(\alpha)}\mu + \sqrt{\sigma(\alpha)}\sigma\epsilon + \sqrt{\sigma(-\alpha)}\epsilon'\end{aligned}$$

We want to represent this in a more canonical way to see what the variance and mean of this Gaussian is. Note that for two standard normal random variables, $a\epsilon + b\epsilon'$, we can represent them as a single random variable with the same variance, $\sqrt{a^2 + b^2}\epsilon''$ (reparameterization trick).

$$\mathbf{x}_\alpha = \sqrt{\sigma(\alpha)}\mu + \sqrt{\sigma(\alpha)\sigma + \sigma(-\alpha)}\epsilon''$$

Now we see that the noisy Gaussian (one component of a mixture) is just a modified version of the original. We have to change the mean (moving it towards zero when adding noise) and the variance.

In our example, we set $\sigma = 1$, so it simplifies further.

$$\mathbf{x}_\alpha = \sqrt{\sigma(\alpha)}\mu + \epsilon''$$

So the variance doesn't change, we just slowly shift the two mixtures together to the center.

Therefore, for one mode $\mathbf{x} \sim \mathcal{N}(\mu, I)$, the intermediate noisy log-density is.

$$\log p_\alpha(\mathbf{x}) = -1/2 \log(2\pi\sigma^2) - \frac{(\mathbf{x} - \sqrt{\sigma(\alpha)}\mu)^2}{2\sigma^2}$$

Take the gradient of the log-density via torch built-in function `torch.autograd.grad(log-density, samples)`, we have the ground-truth score function $\hat{\epsilon}_{gt}(\mathbf{x}, \alpha)$.

B Implementation Details

B.1 Synthetic Experiment – Maximum Mean Discrepancy Bandwidth Choice

The Maximum Mean Discrepancy (MMD) is a statistical test used to determine if two distributions are different. It works by comparing the mean embeddings of samples drawn from two distributions in high dimensional feature space. Specifically, if the distributions are the same, the means should be close; if they are different, the means should be far apart. The embeddings are typically constructed using a feature map associated with a kernel function, and here we select the Gaussian kernel:

$$K(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{y}\|^2}{2\sigma^2}\right)$$

The bandwidth parameter σ of the Gaussian kernel plays a critical role in the sensitivity and performance of MMD. The better bandwidth choice, the more effective MMD computation is. Intuitively, the bandwidth σ controls the scale at which differences between distributions are detected. A small σ makes the kernel sensitive to differences at small scales (fine details), while a large σ highlights differences at larger scales.

The choice of bandwidth is often related to the variance of the data, and the bandwidth should be on the order of the variance of the data. Through a small experiment where we calculate MMD score between our data and the standard Gaussian under various bandwidths σ s, we pick the one that maximizes the MMD score.

In the synthetic 2D Dino experiment, we plot the relationship between MMD scores and bandwidths (Fig. 8), setting $\sigma = 3e - 02$.

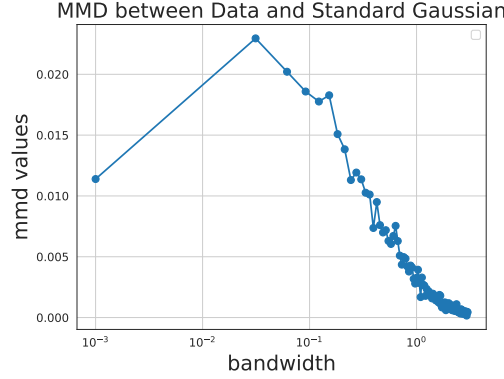


Figure 8: MMD between Dino data and the standard Gaussian. We plots the relationship between the MMD values and the bandwidth parameter used in the kernel function, and pick the bandwidth value with peak MMD score.

B.2 Details on Model Training

We adopt two models as baselines to fine-tune with CDL: DDPM model provided by Ho et al. 2020 and EDM model provided by Karras et al. 2019.

Model Checkpoints For DDPM model, the checkpoint ⁸ we used is a ema one pre-trained on unconditional CIFAR-10. The reason we are not using the most frequently used checkpoint (<https://huggingface.co/google/ddpm-cifar10-32>) is that, this is not EMA checkpoint and our calculation of the FID score of this model on 50,000 generated image via sequential DDPM sampler gives 12.43. This FID score is much higher than what reported on the original paper 3.17. Here we provide the FID scores of this non-EMA pre-trained model, results shown in Tab. 4.

	Parallel DDPM Sampler	Sequential DDPM Sampler
DDPM	10.69	12.43
CDL	7.83	10.06

Table 4: Evaluating sample quality from both parallel and sequential DDPM samplers.

For EDM model, in total eight checkpoints we used are “DDPM++ cont. (VP)” and “NCSN++ cont. (VE)” models pre-trained on three datasets (CIFAR-10, uncond-FFHQ, and uncond-AFHQv2 [14, 10, 3]) with two training settings (unconditional and conditional) ^{9 10 11 12 13 14 15 16}.

As for fine-tuning, we train all models with the same training setting in their original papers. For DDPM model, we train each model for 10 epochs and keep ‘learning rate / batch size’ ratio to be $10^{-4}/64$, and this training is on two A6000 GPUs. For EDM model, the following table list the exact our fine-tuning configurations, which is still of the same ‘learning rate/batch size’ ratio. This training is using eight V100 GPUs.

⁸https://github.com/pesser/pytorch_diffusion

⁹<https://nvlabs-fi-cdn.nvidia.com/edm/pretrained/edm-afhqv2-64x64-uncond-ve.pkl>

¹⁰<https://nvlabs-fi-cdn.nvidia.com/edm/pretrained/edm-afhqv2-64x64-uncond-vp.pkl>

¹¹<https://nvlabs-fi-cdn.nvidia.com/edm/pretrained/edm-cifar10-32x32-cond-ve.pkl>

¹²<https://nvlabs-fi-cdn.nvidia.com/edm/pretrained/edm-cifar10-32x32-cond-vp.pkl>

¹³<https://nvlabs-fi-cdn.nvidia.com/edm/pretrained/edm-cifar10-32x32-uncond-ve.pkl>

¹⁴<https://nvlabs-fi-cdn.nvidia.com/edm/pretrained/edm-cifar10-32x32-uncond-vp.pkl>

¹⁵<https://nvlabs-fi-cdn.nvidia.com/edm/pretrained/edm-ffhq-64x64-uncond-ve.pkl>

¹⁶<https://nvlabs-fi-cdn.nvidia.com/edm/pretrained/edm-ffhq-64x64-uncond-vp.pkl>

Dataset	Fine-tuning Configurations
uncond/cond CIFAR-10	-duration=0.5 -batch=128 -lr=2e-4
uncond AFHQ-64	-duration=0.5 -batch=32 -lr=5e-5 -cres=1,2,2,2 -dropout=0.25 -augment=0.15
uncond FFHQ-64	-batch=32 -lr=5e-5 -cres=1,2,2,2 -dropout=0.05 -augment=0.15

Table 5: Fine-tuning configurations for different datasets

Network and preconditioning	
Architecture of denoising function	(any)
Skip scaling $c_{skip}(\sigma)$	$\sigma_{data}^2 / (\sigma^2 + \sigma_{data}^2)$
Output scaling $c_{out}(\sigma)$	$\sigma \cdot \sigma_{data} / \sqrt{\sigma_{data}^2 + \sigma^2}$
Input scaling $c_{in}(\sigma)$	$1 / \sqrt{\sigma^2 + \sigma_{data}^2}$
Noise Cond. $c_{noise}(\sigma)$	$\frac{1}{4} \ln(\sigma)$
Training	
Noise distribution	$\ln(\sigma) \sim \mathcal{N}(P_{mean}, P_{std}^2)$
Loss weighting $\lambda(\sigma)$	$(\sigma^2 + \sigma_{data}^2) / (\sigma \cdot \sigma_{data})^2$

Table 6: CDL finetune on EDM experiment – fine-tuning design choices.

More Details on EDM Fine-tuning As the design choices of EDM model is very comprehensive and complicate, here we list the training noise distribution, loss weighting, network and preconditioning choices we make during CDL fine-tuning in Tab. 6.

B.3 Relationship Among Log-SNR, Timesteps, and Noise variance Sigma

To use the pre-trained models in the literature with our CDL loss, we need to translate " t ", a parameter representing time in a Markov chain that progressively adds noise to data in Ho et al. 2020 and Song et al. 2020, or " σ ", the variance scale of the Gaussian noise in Karras et al. 2022, to a log-SNR " α ".

Translation between Timesteps and Log-SNR For time-step t in DDPM and stochastic diffusion notation, we recommend readers check Kong et al. 2023 Appendix B.2 about the mapping between α and t .

Translation between Noise Variance Sigma and Log-SNR For variance scale of the Gaussian noise σ in EDM, referring to Eq.(7) and (8) in Karras et al. 2022, it's easy to translate the preconditioning:

$$\mathbf{x}_\alpha \equiv c_{in}(\sigma) \cdot (\mathbf{x} + \sigma \epsilon) \quad (17)$$

$$\mathbf{x}_\alpha \equiv \sqrt{\sigma(\alpha)} \mathbf{x} + \sqrt{\sigma(-\alpha)} \epsilon \equiv \sqrt{\sigma(\alpha)} \left(\mathbf{x} + \sqrt{\frac{\sigma(-\alpha)}{\sigma(\alpha)}} \epsilon \right) \quad (18)$$

From Eq. 18 and Eq. 17, we see that $\sigma \equiv \sqrt{\frac{\sigma(-\alpha)}{\sigma(\alpha)}}$, therefore, the relationship between α and σ should be:

$$\sigma \equiv \exp(-\alpha/2), \alpha \equiv -2 \ln(\sigma)$$

B.4 Contrastive Loss Implementation

To implement contrastive loss, we follow the definition in Sec. 3. First, we generate a random binary label y . Next, conditioned on y , we sample from either data distribution $p(\mathbf{x})$ or the noisy data distribution $p_\zeta(\mathbf{x})$. We calculate the point-wise log-likelihood, then the contrastive loss in Eq. 7.

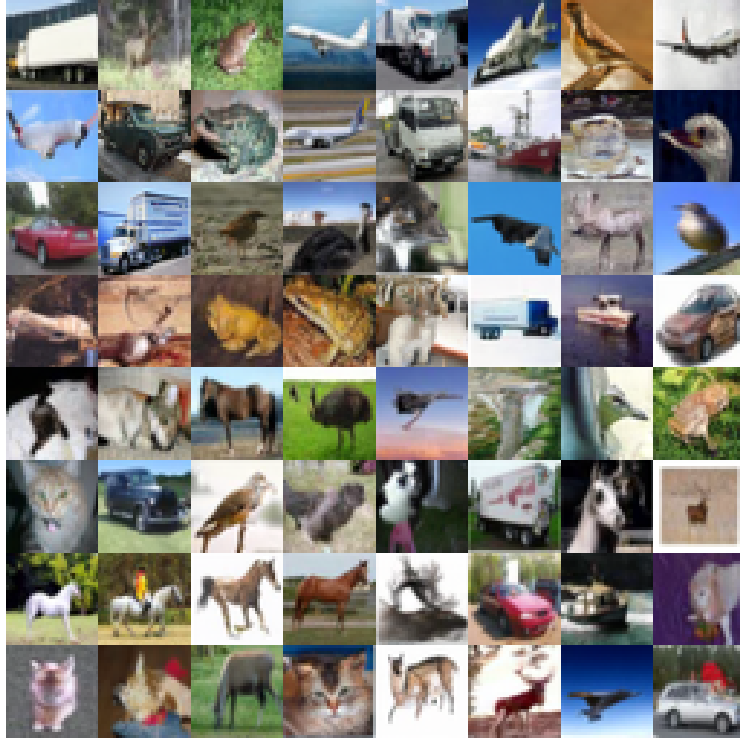


Figure 10: The CDL-loss fine-tuned EDM checkpoint generated examples from Unconditional CIFAR-10, via parallel DDPM sampler.



Figure 11: The CDL-loss fine-tuned EDM checkpoint generated examples from Unconditional AFHQ, via parallel DDPM sampler.

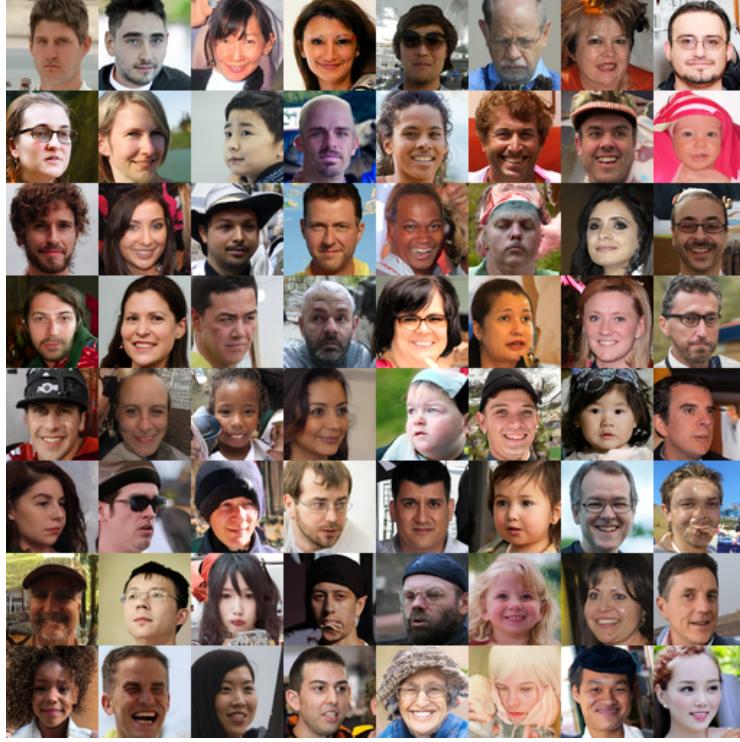


Figure 12: The CDL-loss fine-tuned EDM checkpoint generated examples from Unconditional FFHQ, via parallel DDPM sampler.

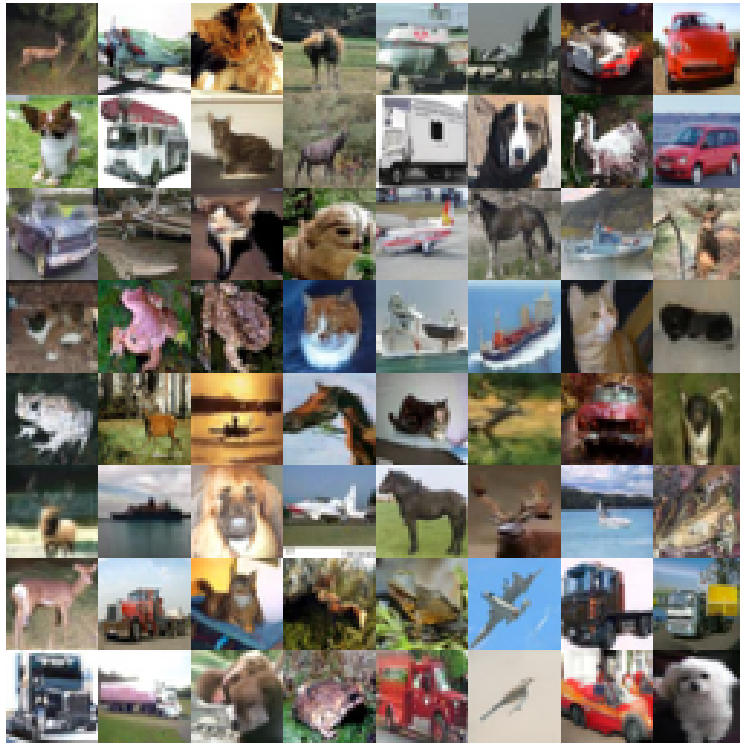


Figure 13: The CDL-loss fine-tuned EDM checkpoint generated examples from Conditional CIFAR-10, via sequential EDM sampler.

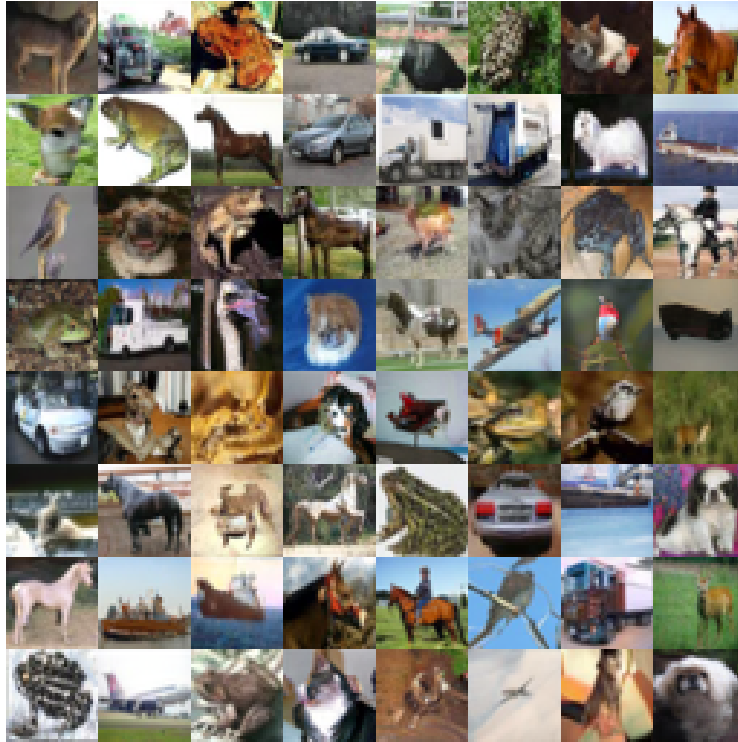


Figure 14: The CDL-loss fine-tuned EDM checkpoint generated examples from Unconditional CIFAR-10, via sequential EDM sampler.



Figure 15: The CDL-loss fine-tuned EDM checkpoint generated examples from Unconditional AFHQ, via sequential EDM sampler.



Figure 16: The CDL-loss fine-tuned EDM checkpoint generated examples from Unconditional FFHQ, via sequential EDM sampler.