

LEARNING ENERGY-BASED MODELS BY COOPERATIVE DIFFUSION RECOVERY LIKELIHOOD

Yaxuan Zhu

UCLA

yaxuanzhu@g.ucla.edu

Jianwen Xie

Akool Research

jianwen@ucla.edu

Ying Nian Wu

UCLA

ywu@stat.ucla.edu

Ruiqi Gao

Google DeepMind

ruiqig@google.com

ABSTRACT

Training energy-based models (EBMs) on high-dimensional data can be both challenging and time-consuming, and there exists a noticeable gap in sample quality between EBMs and other generative frameworks like GANs and diffusion models. To close this gap, inspired by the recent efforts of learning EBMs by maximizing diffusion recovery likelihood (DRL), we propose cooperative diffusion recovery likelihood (CDRL), an effective approach to tractably learn and sample from a series of EBMs defined on increasingly noisy versions of a dataset, paired with an initializer model for each EBM. At each noise level, the two models are jointly estimated within a cooperative training framework: Samples from the initializer serve as starting points that are refined by a few MCMC sampling steps from the EBM. The EBM is then optimized by maximizing recovery likelihood, while the initializer model is optimized by learning from the difference between the refined samples and the initial samples. In addition, we made several practical designs for EBM training to further improve the sample quality. Combining these advances, we significantly boost the generation performance compared to existing EBM methods on CIFAR-10 and ImageNet 32x32. And we have shown that CDRL has great potential to largely reduce the sampling time. We also demonstrate the effectiveness of our models for several downstream tasks, including classifier-free guided generation, compositional generation, image inpainting and out-of-distribution detection.

1 INTRODUCTION

Energy-based models (EBMs), as a class of probabilistic generative models, have exhibited their flexibility and practicality in a variety of application scenarios, such as realistic image synthesis (Xie et al., 2016; Nijkamp et al., 2019; Du & Mordatch, 2019; Arbel et al., 2021; Hill et al., 2022; Xiao et al., 2021; Lee et al., 2023; Grathwohl et al., 2021b), graph generation (Liu et al., 2021), compositional generation (Du et al., 2020; 2023), video generation (Xie et al., 2021c), 3D generation (Xie et al., 2021a; 2018b; Zhu et al., 2023), simulation-based inference (Glaser et al., 2022), stochastic optimization (Kong et al., 2022), out-of-distribution detection (Grathwohl et al., 2020; Liu et al., 2020), continue learning (Wang et al., 2023), internal learning (Zheng et al., 2021), learning set function (Ou et al., 2022), image style transfer (Zhao et al., 2021), continuous inverse optimal control (Xu et al., 2022), and latent space modeling (Pang et al., 2020). Despite these successes of EBMs, training and sampling from EBMs remains challenging, mainly because of the intractability of the partition function in the distribution.

Recently, Diffusion Recovery Likelihood (DRL) (Gao et al., 2021) has emerged as a powerful framework for estimating EBMs. Inspired by diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song & Ermon, 2019), DRL assumes a sequence of EBMs for the marginal distributions of a diffusion process, where each EBM is trained with recovery likelihood that maximizes the conditional probability of the data at the current noise level given their noisy versions at a higher noise level. Maximizing recovery likelihood is more tractable, as sampling from the conditional distribution is much easier than sampling from the marginal distribution. DRL achieves exceptional generation performance among EBM-based generative models. However, a noticeable performance gap still exists between the sample quality of EBMs and other generative frameworks like GANs or diffusion

models. Moreover, DRL requires around 30 MCMC sampling steps at each noise level to generate valid samples, which can be time-consuming during both training and sampling processes.

To further close the performance gap and expedite EBM training and sampling with fewer MCMC steps, we introduce Cooperative Diffusion Recovery Likelihood (CDRL), that jointly estimates a sequence of EBMs and MCMC initializers defined on data perturbed by a diffusion process. At each noise level, the initializer and EBM are updated by cooperative training: The initializer model proposes initial samples by predicting the samples at the current noise level given their noisy versions at a higher noise level. The initial samples are then refined by a few MCMC sampling steps from the conditional distribution defined by the EBM. Given the refined samples, the EBM is updated by maximizing recovery likelihood, and the initializer is updated to absorb the difference between the initial samples and the refined samples. The introduced initializer models learn to accumulate the MCMC transitions of the EBMs, and reproduce them by direct ancestral sampling. Combining with a new noise schedule and a variance reduction technique, we achieve significantly better performance than existing methods of estimating EBMs. We further incorporate classifier-free guidance (CFG) (Ho & Salimans, 2022) to boost the performance of conditional generation, and we observe similar trade-offs between sample quality and sample diversity as CFG for diffusion models when adjusting the guidance strength. In addition, we showcase that our approach can be applied to perform several useful downstream tasks, including compositional generation, image inpainting and out-of-distribution detection.

Our main contributions are as follows: (1) We propose cooperative diffusion recovery likelihood (CDRL) that tractably and efficiently learns and samples from a sequence of EBMs and MCMC initializers; (2) We make several practical design choices related to noise scheduling, MCMC sampling, noise variance reduction for EBM training; (3) Empirically we demonstrate that CDRL achieves significant improvements on sample quality compared to existing EBM approaches, on CIFAR-10 and ImageNet 32×32 datasets; (4) We show that CDRL has great potential to enable more efficient sampling with sampling adjustment techniques; (5) We demonstrate CDRL’s ability in compositional generation, image inpainting and out-of-distribution (OOD) detection, as well as its compatibility with classifier-free guidance for conditional generation.

2 PRELIMINARIES ON ENERGY-BASED MODELS

Let $\mathbf{x} \sim p_{\text{data}}(\mathbf{x})$ be a training example from an underlying data distribution. An energy-based model defines the density of \mathbf{x} by

$$p_{\theta}(\mathbf{x}) = \frac{1}{Z_{\theta}} \exp(f_{\theta}(\mathbf{x})), \quad (1)$$

where f_{θ} is the unnormalized log density, or negative energy, parametrized by a neural network with a scalar output. Z_{θ} is the normalizing constant or partition function. The derivative of the log-likelihood function of an EBM can be approximately written as

$$\mathcal{L}'(\theta) = \mathbb{E}_{p_{\text{data}}} \left[\frac{\partial}{\partial \theta} f_{\theta}(\mathbf{x}) \right] - \mathbb{E}_{p_{\theta}} \left[\frac{\partial}{\partial \theta} f_{\theta}(\mathbf{x}) \right], \quad (2)$$

where the second term is analytically intractable and has to be estimated by Monte Carlo samples from the current model p_{θ} . Therefore, applying gradient-based optimization for an EBM usually involves an inner loop of MCMC sampling, which can be time-consuming for high-dimensional data.

3 COOPERATIVE DIFFUSION RECOVERY LIKELIHOOD

3.1 DIFFUSION RECOVERY LIKELIHOOD

Given the difficulty of sampling from the marginal distribution $p(\mathbf{x})$ defined by an EBM, we could instead estimate a sequence of EBMs defined on increasingly noisy versions of the data and jointly estimate them by maximizing *recovery likelihood*. Specifically, assume a sequence of noisy training examples perturbed by a Gaussian diffusion process: $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_T$ such that $\mathbf{x}_0 \sim p_{\text{data}}$; $\mathbf{x}_{t+1} = \alpha_{t+1}\mathbf{x}_t + \sigma_{t+1}\epsilon$. Denote $\mathbf{y}_t = \alpha_{t+1}\mathbf{x}_t$ for notation simplicity. The marginal distributions of $\{\mathbf{y}_t; t = 1, \dots, T\}$ are modeled by a sequence of EBMs: $p_{\theta}(\mathbf{y}_t) = \frac{1}{Z_{\theta,t}} \exp(f_{\theta}(\mathbf{y}_t; t))$. Then the

conditional EBM of \mathbf{y}_t given the sample \mathbf{x}_{t+1} at a higher noise level can be derived as

$$p_\theta(\mathbf{y}_t|\mathbf{x}_{t+1}) = \frac{1}{\tilde{Z}_{\theta,t}(\mathbf{x}_{t+1})} \exp \left(f_\theta(\mathbf{y}_t; t) - \frac{1}{2\sigma_{t+1}^2} \|\mathbf{y}_t - \mathbf{x}_{t+1}\|^2 \right), \quad (3)$$

where $\tilde{Z}_{\theta,t}(\mathbf{x}_{t+1})$ is the partition function of the conditional EBM dependent on \mathbf{x}_{t+1} . Compared with the marginal EBM $p_\theta(\mathbf{y}_t)$, when σ_{t+1} is small, the extra quadratic term in $p_\theta(\mathbf{y}_t|\mathbf{x}_{t+1})$ constrains the conditional energy landscape to be localized around \mathbf{x}_{t+1} , making the latter less multi-modal and easier to sample from with MCMC. In the extreme case when σ_{t+1} is infinitesimal, $p_\theta(\mathbf{y}_t|\mathbf{x}_{t+1})$ is approximately a Gaussian distribution that can be tractably sampled from and has a close connection to diffusion models (Gao et al., 2021). In the other extreme case when $\sigma_{t+1} \rightarrow \infty$, the conditional distribution falls back to the marginal distribution, and we lose the advantage of being more MCMC friendly for the conditional distribution. Therefore, we need to maintain a small σ_{t+1} between adjacent time steps, and to equip the model with the ability of generating new samples from white noises, we end up with estimating a sequence of EBMs defined on the diffusion process. We use the variance-preserving noise schedule (Song et al., 2021), under which case we have $\mathbf{x}_t = \bar{\alpha}_t \mathbf{x}_0 + \bar{\sigma}_t \epsilon$, where $\bar{\alpha}_t = \prod_{i=1}^T \alpha_i$ and $\bar{\sigma}_t = \sqrt{1 - \bar{\alpha}_t^2}$.

We estimate each EBM by maximizing the following recovery log-likelihood function at each noise level (Bengio et al., 2013):

$$\mathcal{J}_t(\theta) = \frac{1}{n} \sum_{i=1}^n \log p_\theta(\mathbf{y}_{t,i}|\mathbf{x}_{t+1,i}), \quad (4)$$

where $\{\mathbf{y}_{t,i}, \mathbf{x}_{t+1,i}\}$ are pair of samples at time steps t and $t+1$. Sampling from $p_\theta(\mathbf{y}_t|\mathbf{x}_{t+1})$ can be achieved by running K steps of Langevin dynamics from the initialization point $\tilde{\mathbf{y}}_t^0 = \mathbf{x}_{t+1,i}$ and iterating

$$\tilde{\mathbf{y}}_t^{\tau+1} = \tilde{\mathbf{y}}_t^\tau + \frac{s_t^2}{2} \left(\nabla_{\mathbf{y}} f_\theta(\tilde{\mathbf{y}}_t^\tau; t) - \frac{1}{\sigma_{t+1}^2} (\tilde{\mathbf{y}}_t^\tau - \mathbf{x}_{t+1}) \right) + s_t \epsilon^\tau, \quad (5)$$

where s_t is the step size and τ is the index of the current sampling step. With the samples, the updating of EBMs then follows the same learning gradients as MLE (Equation 2), as the extra quadratic term $-\frac{1}{2\sigma_{t+1}^2} \|\mathbf{y}_t - \mathbf{x}_{t+1}\|^2$ in $p_\theta(\mathbf{y}_t|\mathbf{x}_{t+1})$ does not involve learnable parameters. It is worth noting that maximizing recovery likelihood still guarantees an unbiased estimator of the true parameters of the *marginal distribution* of the data.

3.2 AMORTIZING MCMC SAMPLING BY INITIALIZER MODEL

Although $p_\theta(\mathbf{y}_t|\mathbf{x}_{t+1})$ is easier to sample from than $p_\theta(\mathbf{y}_t)$, when σ_{t+1} is not infinitesimal, the initialization of MCMC sampling, \mathbf{x}_{t+1} , may still be far from the data manifold of \mathbf{y}_t . This necessitates a certain amount of MCMC sampling steps at each noise level (e.g., 30 steps of Langevin dynamics in Gao et al. (2021)). Naively reducing the number of sampling steps would lead to training divergence or performance degradation.

To address this issue, we propose to learn an initializer model jointly with the EBM at each noise level, which maps \mathbf{x}_{t+1} closer to the manifold of \mathbf{y}_t . Our work is inspired by the CoopNets work Xie et al. (2018a; 2021b; 2022), which shows that jointly training a top-down generator via MCMC teaching will help the training of a single EBM model. We take this idea and generalize it to the recovery-likelihood model. More discussions are included in Appendix B. Specifically, the initializer model at noise level t is defined as

$$q_\phi(\mathbf{y}_t|\mathbf{x}_{t+1}) \sim \mathcal{N}(\mathbf{g}_\phi(\mathbf{x}_{t+1}; t), \tilde{\sigma}_t^2 \mathbf{I}). \quad (6)$$

It serves as a coarse approximation to $p_\theta(\mathbf{y}_t|\mathbf{x}_{t+1})$, as the former is a single-mode Gaussian distribution while the latter can be multi-modal. A more general formulation would be to involve latent variables \mathbf{z}_t following a certain simple prior $p(\mathbf{z}_t)$ into \mathbf{g}_ϕ . Then $q_\phi(\mathbf{y}_t, t|\mathbf{x}_{t+1}) = \mathbb{E}_{p(\mathbf{z}_t)} [q_\phi(\mathbf{y}_t, \mathbf{z}_t, t|\mathbf{x}_{t+1})]$ can be non-Gaussian (Xiao et al., 2022). However, we empirically find that the simple initializer in Equation 6 works well. Compared with the more general formulation, the simple initializer avoids the inference of \mathbf{z}_t which may again require MCMC sampling, and leads to more stable training. Different from (Xiao et al., 2022), samples from the initializer just serves as

the starting points and are refined by sampling from the EBM, instead of being treated as the final samples. We follow (Ho et al., 2020) to set $\tilde{\sigma}_t = \sqrt{\frac{1-\bar{\alpha}_t^2}{1-\bar{\alpha}_{t+1}^2}}\sigma_t$. If we treat the sequence of initializers as the reverse process, such choice of $\tilde{\sigma}_t$ corresponds to the lower bound of the standard deviation given by p_{data} being a delta function (Sohl-Dickstein et al., 2015).

3.3 COOPERATIVE TRAINING

We jointly train the sequence of EBMs and initializers in a cooperative fashion. Specifically, at each iteration, for a randomly sampled noise level t , we obtain an initial sample $\hat{\mathbf{y}}_t$ from the initializer model. Then a synthesized sample $\tilde{\mathbf{y}}_t$ from $p(\mathbf{y}_t|\mathbf{x}_{t+1})$ is generated by initializing from $\hat{\mathbf{y}}_t$ and running a few steps of Langevin dynamics (Equation 5). The parameters of EBM are then updated by maximizing the recovery log-likelihood function (Equation 4). The learning gradient of EBM is

$$\nabla_{\theta}\mathcal{J}_t(\theta) = \nabla_{\theta} \left[\frac{1}{n} \sum_{i=1}^n f_{\theta}(\mathbf{y}_{t,i}; t) - \frac{1}{n} \sum_{i=1}^n f_{\theta}(\tilde{\mathbf{y}}_{t,i}; t) \right]. \quad (7)$$

To train the initializer model that amortizes the MCMC sampling process, we treat the revised sample $\tilde{\mathbf{y}}_t$ by the EBM as the observed data of the initializer model, and estimate the parameters of the initializer by maximizing log-likelihood:

$$\mathcal{L}_t(\phi) = \frac{1}{n} \sum_{i=1}^n \left[-\frac{1}{2\tilde{\sigma}_t^2} \|\tilde{\mathbf{y}}_{t,i} - \mathbf{g}_{\phi}(\mathbf{x}_{t+1,i}; t)\|^2 \right]. \quad (8)$$

That is, the initializer model learns to absorb the difference between $\hat{\mathbf{y}}_t$ and $\tilde{\mathbf{y}}_t$ at each iteration so that $\hat{\mathbf{y}}_t$ is getting closer to the samples from $p_{\theta}(\mathbf{y}_t|\mathbf{x}_{t+1})$. In practice, we re-weight $\mathcal{L}_t(\phi)$ across different noise levels by removing the coefficient $\frac{1}{2\tilde{\sigma}_t^2}$, similar to the ‘‘simple loss’’ in diffusion models. The training algorithm is summarized in Algorithm 1.

After training, we generate new samples by starting from Gaussian white noise and progressively samples $p_{\theta}(\mathbf{y}_t|\mathbf{x}_{t+1})$ at decreasingly lower noise levels. For each noise level, an initial proposal is generated from the initializer model, followed by a few steps of Langevin dynamics from the EBM. See Algorithm 2 for a summary.

3.4 NOISE VARIANCE REDUCTION

We further propose a simple way to reduce the variance of training gradients. In principle, the pair of \mathbf{x}_t (or \mathbf{y}_t) and \mathbf{x}_{t+1} is generated by $\mathbf{x}_t \sim \mathcal{N}(\bar{\alpha}_t\mathbf{x}_0, \bar{\sigma}_t^2\mathbf{I})$ and $\mathbf{x}_{t+1} \sim \mathcal{N}(\alpha_{t+1}\mathbf{x}_t, \sigma_{t+1}^2\mathbf{I})$. Alternatively, we can fix the Gaussian white noise $\mathbf{e} \sim \mathcal{N}(0, \mathbf{I})$, and sample pair $(\mathbf{x}'_t, \mathbf{x}'_{t+1})$ by

$$\begin{aligned} \mathbf{x}'_t &= \bar{\alpha}_t\mathbf{x}_0 + \bar{\sigma}_t\mathbf{e} \\ \mathbf{x}'_{t+1} &= \bar{\alpha}_{t+1}\mathbf{x}'_t + \bar{\sigma}_{t+1}\mathbf{e}. \end{aligned} \quad (9)$$

In other words, both \mathbf{x}'_t and \mathbf{x}'_{t+1} are linear interpolation between the clean sample \mathbf{x}_0 and a sampled white noise image \mathbf{e} . \mathbf{x}'_t and \mathbf{x}'_{t+1} have the same marginal distributions as \mathbf{x}_t and \mathbf{x}_{t+1} . But \mathbf{x}'_t is deterministic given \mathbf{x}_0 and \mathbf{x}'_{t+1} , while there’s still variance for \mathbf{x}_t given \mathbf{x}_0 and \mathbf{x}_{t+1} . This schedule is related to the ODE forward process used in flow matching (Lipman et al., 2022) and rectified flow (Liu et al., 2022).

3.5 CONDITIONAL GENERATION AND CLASSIFIER-FREE GUIDANCE

(Ho & Salimans, 2022) proposed classifier-free guidance that greatly improves the sample quality of conditional diffusion models, and trades-off between sample quality and sample diversity by adjusting the guidance strength. Given the close connection between EBMs and diffusion models, we show that it is possible to apply classifier-free guidance in CDRL as well. Specifically, suppose c is the context (e.g., a label or a text description). At each noise level we jointly estimate an unconditional EBM $p_{\theta}(\mathbf{y}_t) \propto \exp(f_{\theta}(\mathbf{y}_t; t))$ and a conditional EBM $p_{\theta}(\mathbf{y}_t|c) \propto \exp(f_{\theta}(\mathbf{y}_t; c, t))$. By Bayes rule:

$$p_{\theta}(c|\mathbf{y}_t) = \frac{p_{\theta}(c, \mathbf{y}_t)}{p_{\theta}(\mathbf{y}_t)} = \frac{p_{\theta}(\mathbf{y}_t|c)p(c)}{p_{\theta}(\mathbf{y}_t)}. \quad (10)$$

Algorithm 1 CDRL Training

Input: (1) observed data $\mathbf{x}_0 \sim p_{\text{data}}(\mathbf{x})$; (2) Number of noise levels T ; (3) Number of Langevin sampling steps K per noise level; (4) Langevin step size at each noise level s_t ; (5) Learning rate η_θ for EBM f_θ ; (6) Learning rate η_ϕ for initializer g_ϕ ;

Output: Parameters θ, ϕ

Randomly initialize θ and ϕ .

repeat

Sample noise level t from $\{0, 1, \dots, T-1\}$.

Sample $\epsilon \sim \mathcal{N}(0, \mathbf{I})$. Let $\mathbf{x}_{t+1} = \bar{\alpha}_{t+1}\mathbf{x}_0 + \bar{\sigma}_{t+1}\epsilon$, $\mathbf{y}_t = \alpha_{t+1}(\bar{\alpha}_t\mathbf{x}_0 + \bar{\sigma}_t\epsilon)$.

Generate the initial sample $\hat{\mathbf{y}}_t$ following Equation 6.

Generate the refined sample \mathbf{y}_t by running K steps of Langevin dynamics starting from $\hat{\mathbf{y}}_t$ following Equation 5.

Update EBM parameter θ following the gradients in Equation 7.

Update initializer parameter ϕ by maximizing Equation 8.

until converged

Algorithm 2 CDRL Sampling

Input: (1) Number of noise levels T ; (2) Number of Langevin sampling steps K at each noise level; (3) Langevin step size at each noise level δ_t ; (4) Trained EBM f_θ ; (5) Trained initializer g_ϕ ;

Output: Samples $\tilde{\mathbf{x}}_0$

Randomly initialize $\mathbf{x}_T \sim \mathcal{N}(0, \mathbf{I})$.

for $t = T-1$ **to** 0 **do**

Generate initial proposal $\hat{\mathbf{y}}_t$ following Equation 6.

Update $\hat{\mathbf{y}}_t$ to $\tilde{\mathbf{y}}_t$ by K iterations of Equation 5.

Let $\tilde{\mathbf{x}}_t = \tilde{\mathbf{y}}_t / \alpha_{t+1}$.

end for

With classifier-free guidance, we assume that the log-density of \mathbf{y}_t is scaled to

$$\log \tilde{p}_\theta(\mathbf{y}_t | c) = \log [p_\theta(\mathbf{y}_t | c) p_\theta(c | \mathbf{y}_t)^w] + \text{const.} = (w+1)f_\theta(\mathbf{y}_t; c, t) - w f_\theta(\mathbf{y}_t; t) + \text{const.}, \quad (11)$$

where w controls the guidance strength. Similarly, for the initializer model, we jointly estimate an unconditional model $q_\phi(\mathbf{y}_t | \mathbf{x}_{t+1}) \sim \mathcal{N}(\mathbf{g}_\phi(\mathbf{x}_{t+1}; t), \tilde{\sigma}_t^2 \mathbf{I})$ and a conditional model $q_\phi(\mathbf{y}_t | c, \mathbf{x}_{t+1}) \sim \mathcal{N}(\mathbf{g}_\phi(\mathbf{x}_{t+1}; c, t), \tilde{\sigma}_t^2 \mathbf{I})$. Since both models follow Gaussian distributions, the scaled conditional distribution with classifier-free guidance is still a Gaussian (Dhariwal & Nichol, 2021):

$$\tilde{q}_\phi(\mathbf{y}_t | c, \mathbf{x}_{t+1}) \propto q_\phi(\mathbf{y}_t | c, \mathbf{x}_{t+1}) q_\phi(c | \mathbf{y}_t, \mathbf{x}_{t+1})^w \sim \mathcal{N}((w+1)\mathbf{g}_\phi(\mathbf{x}_{t+1}; c, t) - w\mathbf{g}_\phi(\mathbf{x}_{t+1}; t), \tilde{\sigma}_t^2 \mathbf{I}). \quad (12)$$

3.6 COMPOSITIONALITY IN ENERGY-BASED MODEL

One attractive property of EBMs is compositionality: one can combine multiple EBMs conditioned on individual concepts, and re-normalize it to create a new distribution conditioned on the intersection of those concepts. Specifically, given two EBMs $p_\theta(\mathbf{x} | c_1) \propto \exp(f_\theta(\mathbf{x}; c_1))$ and $p_\theta(\mathbf{x} | c_2) \propto \exp(f_\theta(\mathbf{x}; c_2))$ that are conditional on two separate concepts, (Du et al., 2020; Lee et al., 2023) constructs a new EBM conditional on both concepts as $p_\theta(\mathbf{x} | c_1, c_2) \propto \exp(f_\theta(\mathbf{x}; c_1) + f_\theta(\mathbf{x}; c_2))$ based on the production of expert (Hinton, 2002). Here we show that a negative energy term is missing in that formulation, which can be easily added back with classifier-free guidance and leads to better compositional generation empirically. Specifically, suppose the two concepts c_1 and c_2 are conditionally independent given the observed data \mathbf{x} . Then we have

$$\begin{aligned} \log p_\theta(\mathbf{x} | c_1, c_2) &= \log p_\theta(c_1, c_2 | \mathbf{x}) + \log p_\theta(\mathbf{x}) + \text{const.} \\ &= \log p_\theta(c_1 | \mathbf{x}) + \log p_\theta(c_2 | \mathbf{x}) + \log p_\theta(\mathbf{x}) + \text{const.} \\ &= \log p_\theta(\mathbf{x} | c_1) + \log p_\theta(\mathbf{x} | c_2) - \log p_\theta(\mathbf{x}) + \text{const.} \end{aligned}$$

The composition can be generalized to include arbitrary number of concepts. Suppose we have M conditionally independent concepts $c_i, i = 1, \dots, M$, then

$$\log p_\theta(\mathbf{x} | c_i, i = 1, \dots, M) = \sum_{i=1}^M \log p_\theta(\mathbf{x} | c_i) - (M-1) \log p_\theta(\mathbf{x}) + \text{const.} \quad (13)$$

We can combine the compositional log-density (Equation 13) with classifier-free guidance (Equation 11) to further improve the alignment of generated samples with given concepts. The scaled log-density function is given by

$$\begin{aligned} & \log [p(\mathbf{x}|c_i, i = 1, \dots, M)p(c_i, i = 1, \dots, M|\mathbf{x})^w] \\ &= (w + 1) \sum_{i=1}^M \log p_{\theta}(\mathbf{x}|c_i) - (Mw + M - 1) \log p(\mathbf{x}) + \text{const.} \end{aligned} \quad (14)$$

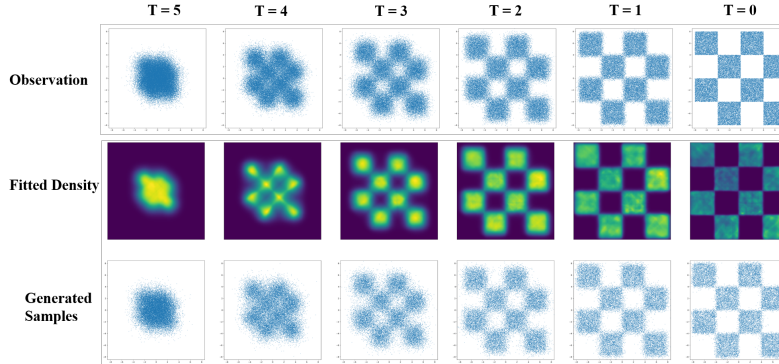


Figure 1: Density estimation results using CDRL on 2D checkerboard distribution. Top: observed samples at each noise level. Middle: fitted density by CDRL at each noise level. Bottom: generated samples at each noise level.

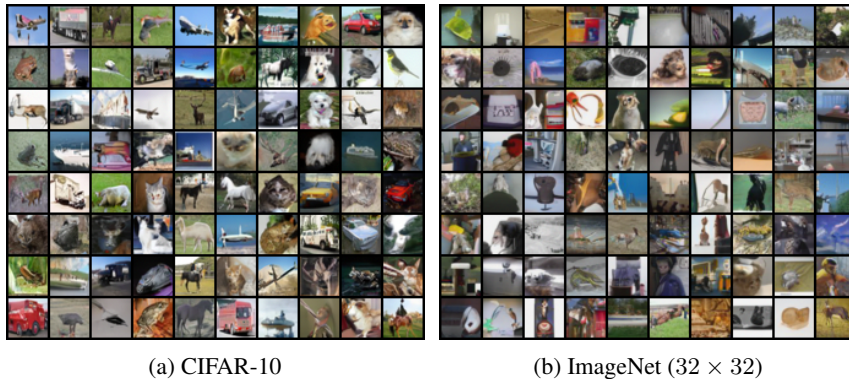


Figure 2: Unconditional generated examples on CIFAR-10 and ImageNet (32 × 32) datasets.

4 EXPERIMENTS

We evaluate our model’s performance across various scenarios. Section 4.1 demonstrates unconditional generation. Section 4.2 highlights our model’s potential to further optimize sampling efficiency. The focus shifts to conditional generation and classifier-free guidance in Section 4.3. Section 4.4 elucidates our model’s prowess in likelihood estimation and OOD detection, and Section 4.5 showcases compositional generation. Please refer to the Appendix A for implementation details, Appendix C for image inpainting with our trained models, Appendix E for comparing the sampling time between CDRL and other EBM models, Appendix F for understanding the role of EBM and initializer in the generation process and Appendix G for ablation study. We name our approach ‘CDRL’ in the following sections.

Our experiments primarily involve three datasets: (i) CIFAR-10 (Krizhevsky & Hinton, 2009) comprises images from 10 categories, with 50k training samples and 10k test samples at a resolution of 32 × 32. We use its training set for unconditional generation. (ii) ImageNet (Deng et al.,

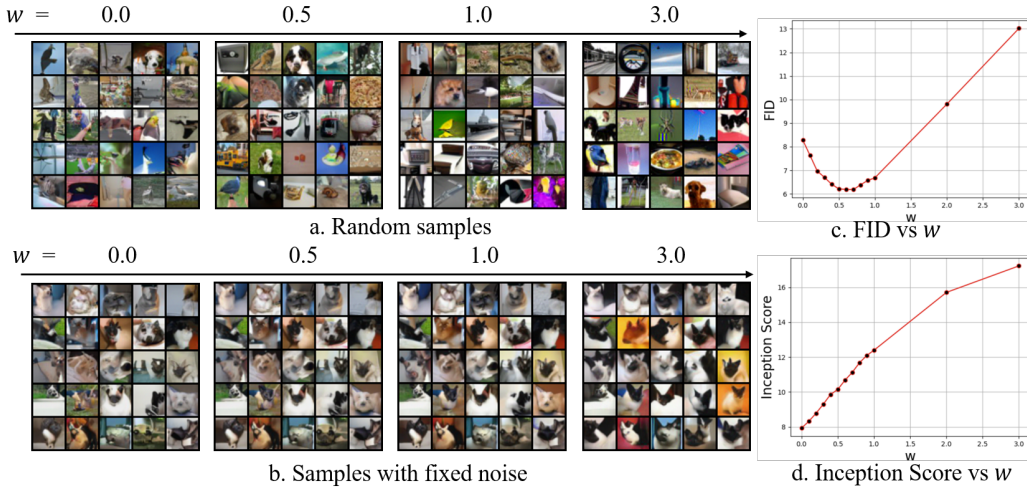


Figure 3: Conditional generation on ImageNet (32×32) with classifier free guidance. (a) Random samples generated under different guided weights w ; (b) Samples generated with fixed noise under different guided weights. The class label is set to be Siamese Cat. Sub-images at the same position are samples with the same random noise as well as class label and only differ in the guided weight; (c) FID Scores at different guided weight w ; (d) Inception Scores at different guided weight w .

2009) contains approximately 1.28M images from 1000 categories. We use its training set for both conditional and unconditional generation, focusing on a downsampled version (32×32) of the dataset. (iii) CelebA (Liu et al., 2015) consists of around 200k human face images, each annotated with attributes. We downsample the dataset to 64×64 and use it for compositionality and image inpainting tasks.

4.1 UNCONDITIONAL IMAGE GENERATION

We first showcase our model’s capabilities in unconditional generation on CIFAR-10 and ImageNet 32×32 datasets. FID scores (Heusel et al., 2017) are reported in Tables 1 and 3 respectively, with generated examples displayed in Figure 2. We use the EBM architecture proposed in Gao et al. (2021). We also employ a larger version called CDRL-large, featuring twice as many channels in each layer. For the initializer network, we follow (Nichol & Dhariwal, 2021)’s structure using a U-Net (Ronneberger et al., 2015) but halving the number of channels. Compared to Gao et al. (2021), CDRL achieves significant improvements in FID scores. Additionally, CDRL uses the same number of noise levels (6 in total) as DRL but requires only half the MCMC steps at each noise level (reduced from 30 to 15), significantly lowering computational costs. With the large architecture, CDRL results in 3.68 FID score on CIFAR-10 and 9.35 on ImageNet 32×32 , which, to the best of our knowledge, are the state-of-the-art among existing EBM approaches, and competitive to other strong generative model classes such as GANs and diffusion models.

4.2 SAMPLING EFFICIENCY

Similar to sampling acceleration of the diffusion model (Song et al.; Liu et al.; Lu et al., 2022b), we foresee the development of post-training techniques to further accelerate CDRL sampling. Although designing an advanced MCMC sampling algorithm could be a standalone project, we present a straightforward yet effective sampling adjustment technique to demonstrate CDRL’s potential in further reducing sampling time. Specifically, we propose to decrease the number of sampling steps, and meanwhile adjust the MCMC sampling step size to be inversely proportional to the square root of the number of sampling steps. As shown in Table 2, while we trained CDRL with 15 MCMC sampling steps at each noise level, we can reduce the sampling steps to 8, 5, and 3 during inference, without losing much perceptual quality.

Table 1: FID for CIFAR-10 unconditional generation.

Models	FID ↓	Models	FID ↓
EBM based method		Other likelihood based method	
NT-EBM (Nijkamp et al., 2022)	78.12	VAE (Kingma & Welling, 2014)	78.41
LP-EBM (Pang et al., 2020)	70.15	PixelCNN (Salimans et al., 2017)	65.93
Adaptive CE (Xiao & Han, 2022)	65.01	PixelIQN (Ostrovski et al., 2018)	49.46
EBM-SR (Nijkamp et al., 2019)	44.50	Residual Flow (Chen et al., 2019)	47.37
JEM (Grathwohl et al., 2020)	38.40	Glow (Kingma & Dhariwal, 2018)	45.99
EBM-IG (Du & Mordatch, 2019)	38.20	DC-VAE (Parmar et al., 2021)	17.90
EBM-FCE (Gao et al., 2020)	37.30	GAN based method	
CoopVAEBM (Xie et al., 2021b)	36.20	WGAN-GP (Gulrajani et al., 2017)	36.40
CoopNets (Xie et al., 2018a)	33.61	SN-GAN (Miyato et al., 2018)	21.70
Divergence Triangle (Han et al., 2020)	30.10	BigGAN (Brock et al., 2019)	14.80
VARA (Grathwohl et al., 2021b)	27.50	StyleGAN2-DiffAugment (Zhao et al., 2020)	5.79
EBM-CD (Du et al., 2021)	25.10	Diffusion-GAN (Xiao et al., 2022)	3.75
GEBM (Arbel et al., 2021)	19.31	StyleGAN2-ADA (Karras et al., 2020)	2.92
HAT-EBM (Hill et al., 2022)	19.30	Score based and Diffusion method	
CF-EBM (Zhao et al., 2021)	16.71	NCSN (Song & Ermon, 2019)	25.32
CoopFlow (Xie et al., 2022)	15.80	NCSN-v2 (Song & Ermon, 2020)	10.87
CLEL-base (Lee et al., 2023)	15.27	NCSN++ (Song et al., 2021)	2.20
VAEBM (Xiao et al., 2021)	12.16	DDPM Distillation (Luhman & Luhman, 2021)	9.36
DRL (Gao et al., 2021)	9.58	DDPM++(VP, NLL) (Kim et al., 2021)	3.45
CLEL-large (Lee et al., 2023)	8.61	DDPM (Ho et al., 2020)	3.17
EGC (Unsupervised) (Guo et al., 2023)	5.36	DDPM++(VP, FID) (Kim et al., 2021)	2.47
CDRL (Ours)	4.31		
CDRL-large (Ours)	3.68		

4.3 CONDITIONAL SYNTHESIS WITH CLASSIFIER-FREE GUIDANCE

We evaluate our model for conditional generation on the ImageNet32 dataset, employing classifier-free guidance as outlined in Section 3.5. Generation results for varying guided weights w are displayed in Figure 3. As w increases, sample quality improves, and the conditioned class features become more prominent, though diversity may decrease. This is also evident from the FID and Inception Score (Salimans et al., 2016) curves in Figures 3c and 3d. While the Inception Score consistently improves (increasing), the FID score first improves (decreasing) and then worsens (increasing), reaching the optimal value of 6.18 at a guidance weight of 0.7. Additional generation results can be found in the Appendix D.1.

Table 2: FID for CIFAR-10 with sampling adjust- Table 3: FID for ImageNet (32×32) unconditional generation.

Models	Number of noise level \times Number of MCMC steps	FID ↓	Models	FID ↓
DRL (Gao et al., 2021)	$6 \times 30 = 180$	9.58	EBM-IG (Du & Mordatch, 2019)	60.23
CDRL	$6 \times 15 = 90$	4.31	PixelCNN (Salimans et al., 2017)	40.51
CDRL (step 8)	$6 \times 8 = 48$	4.58	EBM-CD (Du et al., 2021)	32.48
CDRL (step 5)	$6 \times 5 = 30$	5.37	CF-EBM (Zhao et al., 2021)	26.31
CDRL (step 3)	$6 \times 3 = 18$	9.67	CLEL-base (Lee et al., 2023)	22.16
			DRL (Gao et al., 2021)	- (not converge)
			DDPM++(VP, NLL) (Kim et al., 2021)	8.42
			CDRL (Ours)	9.35

4.4 LIKELIHOOD ESTIMATION AND OUT-OF-DISTRIBUTION DETECTION

A distinctive feature of the EBM is its ability to model the unnormalized log-likelihood directly using the energy function. This capability allows it to undertake tasks beyond generation. We initially showcase CDRL’s prowess in estimating the density of a 2D checkerboard distribution, with results presented in Figure 1. These findings confirm CDRL’s capacity to accurately determine log-likelihood while concurrently generating valid samples.

Moreover, we demonstrate CDRL’s utility in out-of-distribution (OOD) detection tasks. For this endeavor, we employ the model trained on CIFAR-10 as a detector and use the energy at the lowest

noise level to serve as the OOD prediction score. The AUROC score of our CDRL model, with CIFAR-10 interpolation, CIFAR-100, and CelebA data as OOD samples, is provided in Table 4. CDRL achieves strong results in OOD detection comparing with the baseline approaches. More results can be found in Table 8 in the appendix.

4.5 COMPOSITIONALITY



Figure 4: Attribute-compositional generation on CelebA (64×64) with guided weight $w = 3$. Left: generated samples under different attribute compositions. Right: control attributes ('√', 'x' and '-' indicate true, false and no control respectively).

To evaluate the compositionality of EBMs, we conduct experiments on CelebA (64×64) datasets with *Male*, *Smile*, and *Young* as the three conditional concepts. We estimate EBMs conditional on each single concept separately, and assume simple unconditional initializer models. Classifier-free guidance is adopted when conducting compositional generation (Equation 14). Specifically, we treat images with a certain attribute value as individual classes. We randomly assign each image in a training batch to a class based on the controlled attribute value. For example, an image with Male=True and Smile=True may be assigned to class 0 if the Male attribute is picked or class 2 if the Smile attribute is picked. For the conditional network structure, we make EBM f_θ conditional on attributes c_i and use an unconditional initializer model g_ϕ to propose the initial distribution. We focus on showcasing the compositionality ability of EBM itself, although it is also possible to use a conditional initializer model similar to Section 3.5. Our results are displayed in Figure 4, with images generated at a guided weight of $w = 3.0$. More generation results with different guidance weights can be found in the Appendix D.1. Images generated with composed attributes following Equation 14 contain features of both attributes, and increasing the guided weight makes the corresponding attribute more prominent. This demonstrates CDRL’s ability and the effectiveness of Equation 14.

Table 4: AUROC scores in OOD detection using CDRL and other explicit density models on CIFAR-10

	Cifar-10 interpolation	Cifar-100	CelebA
PixelCNN (Salimans et al., 2017)	0.71	0.63	-
GLOW (Kingma & Dhariwal, 2018)	0.51	0.55	0.57
NVAE (Vahdat & Kautz, 2020)	0.64	0.56	0.68
EBM-IG (Du & Mordatch, 2019)	0.70	0.50	0.70
VAEBM (Xiao et al., 2021)	0.70	0.62	0.77
EBM-CD (Du et al., 2021)	0.65	0.83	-
CLEL-Base (Lee et al., 2023)	0.72	0.72	0.77
CDRL (ours)	0.75	0.78	0.84

5 CONCLUSION AND FUTURE WORK

We propose CDRL, an EBM learning method employing cooperative diffusion recovery likelihood, significantly improving generation performance of EBMs. CDRL excels in compositional generation, out-of-distribution detection, image inpainting, and compatibility with classifier-free guidance for conditional generation. The present limitation might be the fact that a certain number of MCMC steps is still needed during generation. Besides, we also seek to scale our model for high-resolution image generation in the future. Our work aims to promote further research on EBMs as generative models. However, powerful generative models may pose negative social impacts, such as deepfakes, misinformation, privacy invasion, and eroding public trust, necessitating exploration of effective preventive measures.

REFERENCES

- Michael Arbel, Liang Zhou, and Arthur Gretton. Generalized energy based models. In *The Ninth International Conference on Learning Representations, ICLR*, 2021.
- Yoshua Bengio, Li Yao, Guillaume Alain, and Pascal Vincent. Generalized denoising auto-encoders as generative models. *Advances in neural information processing systems*, 2013.
- Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. In *The Seventh International Conference on Learning Representations, ICLR*, 2019.
- Tong Che, Ruixiang Zhang, Jascha Sohl-Dickstein, Hugo Larochelle, Liam Paull, Yuan Cao, and Yoshua Bengio. Your gan is secretly an energy-based model and you should use discriminator driven latent sampling. *Advances in Neural Information Processing Systems*, 2020.
- Ricky TQ Chen, Jens Behrmann, David Duvenaud, and Jörn-Henrik Jacobsen. Residual flows for invertible generative modeling. In *Advances in Neural Information Processing Systems*, 2019.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, CVPR*, 2009.
- Prafulla Dhariwal and Alexander Quinn Nichol. Diffusion models beat gans on image synthesis. In *Advances in Neural Information Processing Systems*, 2021.
- Yilun Du and Igor Mordatch. Implicit generation and generalization in energy-based models. *arXiv preprint arXiv:1903.08689*, 2019.
- Yilun Du, Shuang Li, and Igor Mordatch. Compositional visual generation with energy based models. *Advances in Neural Information Processing Systems*, 2020.
- Yilun Du, Shuang Li, Joshua B. Tenenbaum, and Igor Mordatch. Improved contrastive divergence training of energy-based models. In *Proceedings of the 38th International Conference on Machine Learning, ICML*, 2021.
- Yilun Du, Conor Durkan, Robin Strudel, Joshua B Tenenbaum, Sander Dieleman, Rob Fergus, Jascha Sohl-Dickstein, Arnaud Doucet, and Will Grathwohl. Reduce, reuse, recycle: Compositional generation with energy-based diffusion models and mcmc. In *Proceedings of the Fortieth International Conference on Machine Learning, ICML*, 2023.
- Bradley Efron. Tweedie’s formula and selection bias. *Journal of the American Statistical Association*, 106(496):1602–1614, 2011.
- Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 12873–12883, 2021.
- Ruiqi Gao, Yang Lu, Junpei Zhou, Song-Chun Zhu, and Ying Nian Wu. Learning generative convnets via multi-grid modeling and sampling. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, pp. 9155–9164, 2018.
- Ruiqi Gao, Erik Nijkamp, Diederik P Kingma, Zhen Xu, Andrew M Dai, and Ying Nian Wu. Flow contrastive estimation of energy-based models. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2020.
- Ruiqi Gao, Yang Song, Ben Poole, Ying Nian Wu, and Diederik P. Kingma. Learning energy-based models by diffusion recovery likelihood. In *The Ninth International Conference on Learning Representations, ICLR*, 2021.
- Pierre Glaser, Michael Arbel, Arnaud Doucet, and Arthur Gretton. Maximum likelihood learning of energy-based models for simulation-based inference. *arXiv preprint arXiv:2210.14756*, 2022.

- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.
- Will Grathwohl, Kuan-Chieh Wang, Jörn-Henrik Jacobsen, David Duvenaud, Mohammad Norouzi, and Kevin Swersky. Your classifier is secretly an energy based model and you should treat it like one. In *The Eighth International Conference on Learning Representations, ICLR*, 2020.
- Will Sussman Grathwohl, Jacob Jin Kelly, Milad Hashemi, Mohammad Norouzi, Kevin Swersky, and David Duvenaud. No MCMC for me: Amortized sampling for fast and stable training of energy-based models. In *The Ninth International Conference on Learning Representations, ICLR*, 2021a.
- Will Sussman Grathwohl, Jacob Jin Kelly, Milad Hashemi, Mohammad Norouzi, Kevin Swersky, and David Duvenaud. No MCMC for me: Amortized sampling for fast and stable training of energy-based models. In *The Ninth International Conference on Learning Representations, ICLR*, 2021b.
- Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron Courville. Improved training of wasserstein gans. In *The Sixth International Conference on Learning Representations, ICLR*, 2017.
- Qiushan Guo, Chuofan Ma, Yi Jiang, Zehuan Yuan, Yizhou Yu, and Ping Luo. Eg3: Image generation and classification via a diffusion energy-based model. *arXiv preprint arXiv:2304.02012*, 2023.
- Tian Han, Erik Nijkamp, Xiaolin Fang, Mitch Hill, Song-Chun Zhu, and Ying Nian Wu. Divergence triangle for joint training of generator model, energy-based model, and inferential model. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2019.
- Tian Han, Erik Nijkamp, Linqi Zhou, Bo Pang, Song-Chun Zhu, and Ying Nian Wu. Joint training of variational auto-encoder and latent energy-based model. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2020.
- 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. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems*, 2017.
- Mitch Hill, Erik Nijkamp, Jonathan Mitchell, Bo Pang, and Song-Chun Zhu. Learning probabilistic models from generator latent spaces with hat ebm. In *NeurIPS*, 2022.
- Geoffrey E Hinton. Training products of experts by minimizing contrastive divergence. *Neural computation*, 14(8):1771–1800, 2002.
- Geoffrey E. Hinton. A practical guide to training restricted boltzmann machines. In *Neural Networks: Tricks of the Trade - Second Edition*, pp. 599–619. Springer, 2012.
- Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, 2020.
- Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Training generative adversarial networks with limited data. In *Advances in Neural Information Processing Systems*, 2020.
- Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *Advances in Neural Information Processing Systems*, 2022.
- Dongjun Kim, Seungjae Shin, Kyungwoo Song, Wanmo Kang, and Il-Chul Moon. Soft truncation: A universal training technique of score-based diffusion model for high precision score estimation. In *International Conference on Machine Learning, ICML*, 2021.

- Taesup Kim and Yoshua Bengio. Deep directed generative models with energy-based probability estimation. *arXiv preprint arXiv:1606.03439*, 2016.
- Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *Advances in neural information processing systems*, 2021.
- Diederik P. Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. In *Advances in Neural Information Processing Systems*, 2018.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *The Second International Conference on Learning Representations, ICLR*, 2014.
- Lingkai Kong, Jiaming Cui, Yuchen Zhuang, Rui Feng, B. Aditya Prakash, and Chao Zhang. End-to-end stochastic optimization with energy-based model. In *NeurIPS*, 2022.
- Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. 2009.
- Rithesh Kumar, Sherjil Ozair, Anirudh Goyal, Aaron Courville, and Yoshua Bengio. Maximum entropy generators for energy-based models. *arXiv preprint arXiv:1901.08508*, 2019.
- Yann LeCun, Sumit Chopra, Raia Hadsell, M Ranzato, and Fufie Huang. A tutorial on energy-based learning. *Predicting structured data*, 1(0), 2006.
- Hankook Lee, Jongheon Jeong, Sejun Park, and Jinwoo Shin. Guiding energy-based models via contrastive latent variables. In *The Eleventh International Conference on Learning Representations, ICLR*, 2023.
- Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. Pseudo numerical methods for diffusion models on manifolds. In *10th International Conference on Learning Representations, ICLR 2022*.
- Meng Liu, Keqiang Yan, Bora Oztekin, and Shuiwang Ji. Graphebm: Molecular graph generation with energy-based models. In *EBM Workshop at ICLR*, 2021.
- Weitang Liu, Xiaoyun Wang, John D. Owens, and Yixuan Li. Energy-based out-of-distribution detection. In *Advances in Neural Information Processing Systems*, 2020.
- Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of the 2015 International Conference on Computer Vision (ICCV)*, Santiago, Chile, 2015.
- Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ODE solver for diffusion probabilistic model sampling in around 10 steps. In *NeurIPS*, 2022a.
- Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ODE solver for diffusion probabilistic model sampling in around 10 steps. In *NeurIPS*, 2022b.
- Eric Luhman and Troy Luhman. Knowledge distillation in iterative generative models for improved sampling speed. *arXiv preprint arXiv:2101.02388*, 2021.
- Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. In *The Sixth International Conference on Learning Representations, ICLR*, 2018.
- Jiquan Ngiam, Zhenghao Chen, Pang Wei Koh, and Andrew Y. Ng. Learning deep energy models. In *Proceedings of the 28th international conference on machine learning, ICML*, 2011.
- Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *Proceedings of the 38th International Conference on Machine Learning, ICML*, 2021.

- Erik Nijkamp, Mitch Hill, Song-Chun Zhu, and Ying Nian Wu. Learning non-convergent non-persistent short-run MCMC toward energy-based model. In *Advances in Neural Information Processing Systems*, 2019.
- Erik Nijkamp, Ruiqi Gao, Pavel Sountsov, Srinivas Vasudevan, Bo Pang, Song-Chun Zhu, and Ying Nian Wu. Learning energy-based model with flow-based backbone by neural transport mcmc. In *The Tenth International Conference on Learning Representations, ICLR*, 2022.
- Georg Ostrovski, Will Dabney, and Rémi Munos. Autoregressive quantile networks for generative modeling. In *Proceedings of the 35th International Conference on Machine Learning, ICML*, 2018.
- Zijing Ou, Tingyang Xu, Qinliang Su, Yingzhen Li, Peilin Zhao, and Yatao Bian. Learning neural set functions under the optimal subset oracle. *NeurIPS*, 2022.
- Bo Pang, Tian Han, Erik Nijkamp, Song-Chun Zhu, and Ying Nian Wu. Learning latent space energy-based prior model. In *Advances in Neural Information Processing Systems*, 2020.
- Gaurav Parmar, Dacheng Li, Kwonjoon Lee, and Zhuowen Tu. Dual contradistinctive generative autoencoder. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2021.
- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.
- Herbert E Robbins. *An empirical Bayes approach to statistics*. Springer, 1992.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022*. IEEE, 2022.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In Nassir Navab, Joachim Hornegger, William M. Wells III, and Alejandro F. Frangi (eds.), *Medical Image Computing and Computer-Assisted Intervention - MICCAI*, 2015.
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L. Denton, Seyed Kamyar Seyed Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, Jonathan Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. In *NeurIPS*, 2022.
- Tim Salimans and Jonathan Ho. Should ebms model the energy or the score? In *Energy Based Models Workshop-ICLR 2021*, 2021.
- Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. In *The Tenth International Conference on Learning Representations, ICLR*, 2022.
- Tim Salimans, Ian J. Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. In Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems*, 2016.
- Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P Kingma. Pixelcnn++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications. In *The Fifth International Conference on Learning Representations, ICLR*, 2017.
- Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *Proceedings of the 32nd International Conference on Machine Learning, ICML*, 2015.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *9th International Conference on Learning Representations, ICLR 2021*.
- Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In *Advances in Neural Information Processing Systems*, 2019.
- Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. In *Advances in Neural Information Processing Systems*, 2020.

- Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *The Ninth International Conference on Learning Representations, ICLR 2021*, 2021.
- Arash Vahdat and Jan Kautz. NVAE: A deep hierarchical variational autoencoder. In *Advances in Neural Information Processing Systems*, 2020.
- Fu-Yun Wang, Da-Wei Zhou, Liu Liu, Han-Jia Ye, Yatao Bian, De-Chuan Zhan, and Peilin Zhao. Beef: Bi-compatible class-incremental learning via energy-based expansion and fusion. In *The Eleventh International Conference on Learning Representations, ICLR, 2023*.
- Zhisheng Xiao and Tian Han. Adaptive multi-stage density ratio estimation for learning latent space energy-based model. In *NeurIPS*, 2022.
- Zhisheng Xiao, Karsten Kreis, Jan Kautz, and Arash Vahdat. Vaebm: A symbiosis between variational autoencoders and energy-based models. In *The Ninth International Conference on Learning Representations, ICLR, 2021*.
- Zhisheng Xiao, Karsten Kreis, and Arash Vahdat. Tackling the generative learning trilemma with denoising diffusion gans. In *The Tenth International Conference on Learning Representations, ICLR, 2022*.
- Jianwen Xie, Yang Lu, Song-Chun Zhu, and Ying Nian Wu. A theory of generative convnet. In *Proceedings of the 33rd International Conference on Machine Learning, ICML, 2016*.
- Jianwen Xie, Yang Lu, Ruiqi Gao, Song-Chun Zhu, and Ying Nian Wu. Cooperative training of descriptor and generator networks. *IEEE transactions on pattern analysis and machine intelligence*, 42(1):27–45, 2018a.
- Jianwen Xie, Zilong Zheng, Ruiqi Gao, Wenguan Wang, Song-Chun Zhu, and Ying Nian Wu. Learning descriptor networks for 3d shape synthesis and analysis. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2018b.
- Jianwen Xie, Yang Lu, Ruiqi Gao, Song-Chun Zhu, and Ying Nian Wu. Cooperative training of descriptor and generator networks. *IEEE Trans. Pattern Anal. Mach. Intell.*, 42(1):27–45, 2020.
- Jianwen Xie, Yifei Xu, Zilong Zheng, Song-Chun Zhu, and Ying Nian Wu. Generative pointnet: Deep energy-based learning on unordered point sets for 3d generation, reconstruction and classification. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2021a.
- Jianwen Xie, Zilong Zheng, and Ping Li. Learning energy-based model with variational auto-encoder as amortized sampler. In *AAAI Conference on Artificial Intelligence (AAAI)*, pp. 10441–10451, 2021b.
- Jianwen Xie, Song-Chun Zhu, and Ying Nian Wu. Learning energy-based spatial-temporal generative convnets for dynamic patterns. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(2):516–531, 2021c.
- Jianwen Xie, Yaxuan Zhu, Jun Li, and Ping Li. A tale of two flows: Cooperative learning of langevin flow and normalizing flow toward energy-based model. In *The Tenth International Conference on Learning Representations, ICLR, 2022*.
- Yifei Xu, Jianwen Xie, Tianyang Zhao, Chris Baker, Yibiao Zhao, and Ying Nian Wu. Energy-based continuous inverse optimal control. *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- Xuwang Yin, Shiyang Li, and Gustavo K Rohde. Learning energy-based models with adversarial training. In *European Conference on Computer Vision*, pp. 209–226. Springer, 2022.
- Sangwoong Yoon, Young-Uk Jin, Yung-Kyun Noh, and Frank C Park. Energy-based models for anomaly detection: A manifold diffusion recovery approach. *arXiv preprint arXiv:2310.18677*, 2023.

- Shengyu Zhao, Zhijian Liu, Ji Lin, Jun-Yan Zhu, and Song Han. Differentiable augmentation for data-efficient GAN training. In *Advances in Neural Information Processing Systems*, 2020.
- Yang Zhao, Jianwen Xie, and Ping Li. Learning energy-based generative models via coarse-to-fine expanding and sampling. In *The Ninth International Conference on Learning Representations, ICLR*, 2021.
- Zilong Zheng, Jianwen Xie, and Ping Li. Patchwise generative convnet: Training energy-based models from a single natural image for internal learning. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, pp. 2961–2970, 2021.
- Song Chun Zhu, Yingnian Wu, and David Mumford. Filters, random fields and maximum entropy (frame): Towards a unified theory for texture modeling. *International Journal of Computer Vision*, 27:107–126, 1998.
- Yaxuan Zhu, Jianwen Xie, and Ping Li. Likelihood-based generative radiance field with latent space energy-based model for 3d-aware disentangled image representation. In *International Conference on Artificial Intelligence and Statistics, AISTATS*, 2023.

CONTENTS

1	Introduction	1
2	Preliminaries on Energy-Based Models	2
3	Cooperative Diffusion Recovery Likelihood	2
3.1	Diffusion recovery likelihood	2
3.2	Amortizing MCMC sampling by initializer model	3
3.3	Cooperative training	4
3.4	Noise variance reduction	4
3.5	Conditional generation and classifier-free guidance	4
3.6	Compositionality in energy-based model	5
4	Experiments	6
4.1	Unconditional image generation	7
4.2	Sampling efficiency	7
4.3	Conditional synthesis with classifier-free guidance	8
4.4	Likelihood Estimation and Out-Of-Distribution Detection	8
4.5	Compositionality	9
5	Conclusion and Future Work	9
A	Training Details	17
A.1	Network Structure	17
A.2	Hyper-parameters	18
A.3	Noise schedule and conditioning input	18
A.4	Overall Illustration of CDRL	19
B	Related Work	19
C	Image inpainting	20
D	More experimental Results	21
D.1	More generation results	21
D.2	Generate high resolution image	21
D.3	More OOD results	21
E	Sampling time	21
F	Understand the role of initializer and EBM	22
G	Ablation Study	22

G.1	What are the contributions of each individual design component?	22
G.2	Should the initializer be learned using the cooperative training algorithm or should it directly regress on the data?	23
G.3	Whether adding more diffusion levels to the origin DRL schedule to cover the high-noise region works better than the new schedule used in CDRL?	25
G.4	Can we further reduce the number of noise levels?	25
G.5	What is the effect of using different number of Langevin steps?	25

A TRAINING DETAILS

A.1 NETWORK STRUCTURE

We adopt the EBM structure from (Gao et al., 2021), starting with a 3×3 convolution layer with 128 channels (doubled to 256 in the CDRL(fat) setting). We use several downsample blocks for resolution adjustments, each containing multiple residual blocks. A 2×2 average pooling layer is applied to all downsample blocks except the last. Spectral normalization is applied to all convolution layers for stability, while ReLU activation is applied to the final feature map. Spatial and channel dimensions are summed to obtain the energy output. EBM building block structures are shown in Table 5, and network architecture hyperparameters in Table 6.

For the initializer network, we follow (Nichol & Dhariwal, 2021)’s structure using a Unet (Ronneberger et al., 2015) but halving the number of channels, reducing the initializer model size. For 32×32 images, we have feature map resolutions at 32×32 , 16×16 , and 4×4 . For 64×64 images, we add a 64×64 resolution. We set all feature map channel numbers to 64 and apply attention to resolutions 16×16 and 8×8 . Our initializer directly outputs the noised image \tilde{y}_t at predicted level t , while (Ho et al., 2020) outputs total injected noise ϵ .

For the class-conditioned generation task, we map class labels to one-hot vectors and use a fully-connected layer to map these vectors to class embedding vectors with the same dimensions as time embedding vectors. The class embedding is then added to the time embedding. We set the time embedding dimension to 512 for EBM and 256 for the initializer in the CDRL setting. In the CDRL(fat) setting, the time embedding dimension increases to 1024 for EBM, while the initializer’s dimension remains unchanged.

Table 5: Building block of the EBM.

(a) ResBlock	(b) Downsample Block	(c) Time Embedding
leakyReLU, 3×3 Conv2D	N ResBlocks	Sinusoidal Embedding
+ Dense(leakyReLU(temb))	Downsample 2×2	Dense, leakyReLU
leakyReLU, 3×3 Conv2D		Dense
+ Input		

Table 6: Structure hyper-parameter of EBM at different setting.

Model	M Downsample Blocks	N Resblocks Per Downsample Block	Number of Channels In Each Resolution
CDRL Image Size 32×32	4	8	(128, 256, 256, 256)
CDRL(fat) Image Size 32×32	4	8	(256, 512, 512, 512)
Compositionality Experiment	5	2	(128, 256, 256, 256, 256)
Inpainting Experiment	5	8	(128, 256, 256, 256, 256)

A.2 HYPER-PARAMETERS

We set the learning rate of EBM to be $\eta_\theta = 1e-4$ and the learning rate of initializer to be $\eta_\phi = 1e-5$. We use linear warm up for both EBM and initializer and let the initializer to start earlier than EBM. More specifically, given training iteration iter, we have:

$$\begin{aligned}\eta_\theta &= \min(1.0, \frac{\text{iter}}{10000}) \times 1e-4 \\ \eta_\phi &= \min(1.0, \frac{\text{iter} + 500}{10000}) \times 1e-5\end{aligned}\tag{15}$$

We use Adam as optimizer for both EBM and initializer. We set the $\beta s = (0.9, 0.999)$ and weight decay equals to 0.0. We also applied exponential moving average which has decay equals to 0.9999 to both the EBM and initializer. We use 8 A100 GPU for training. The training process usually needs around 400k iterations, which takes up around 6 days.

Following (Gao et al., 2021), we use a re-parameterization trick for the calculating the energy term. We build our EBM on $t = 0, 1, 2, 3, 4, 5$ and we assume the distribution of $t = 6$ to be simple Normal distribution during sampling. Given \mathbf{y}_t under noise level t , suppose we denote the output of the EBM network as $\hat{f}_\theta(\mathbf{y}_t, t)$, then we let the true energy term to be $f_\theta(\mathbf{y}_t, t) = \frac{\hat{f}_\theta(\mathbf{y}_t, t)}{s_t^2}$, where s_t is the Langevin step size at noise level t , i.e. we parameterize the energy as the multiplication of EBM network output and a noise-level dependent coefficient and we set this coefficient to equal to the square of Langevin step size. We use 15 step Langevin updates at each noise level and we let the Langevin step size of noise level t to follow:

$$s_t^2 = 0.054 \times \bar{\sigma}_t \times \sigma_{t+1}^2\tag{16}$$

where σ_{t+1} is the variance of the noise adding between noise level t and $t + 1$ and $\bar{\sigma}_t$ is the accumulative noise level at noise level t . During generation, we start from randomly sample $x_6 \sim \mathcal{N}(0, \mathbf{I})$ and do denosing using both initializer and Langevin Dynamics following Algorithm 2. During image generation, after getting samples \mathbf{x}_0 of the lowest noise level $t = 0$, we do one more denoising step without adding noise to further improve its quality. More specifically, we follow Tweedie’s formula (Efron, 2011; Robbins, 1992), which says that if we have $\mathbf{x} \sim p_{data}(\mathbf{x})$ and noise version image \mathbf{x}' which has conditional distribution $p(\mathbf{x}'|\mathbf{x}) = \mathcal{N}(\mathbf{x}, \sigma^2 \mathbf{I})$. The marginal distribution can be defined as $p(\mathbf{x}') = \int p_{data}(\mathbf{x})p(\mathbf{x}'|\mathbf{x})d\mathbf{x}$. Then we have

$$\mathbb{E}(\mathbf{x}|\mathbf{x}') = \mathbf{x}' + \sigma^2 \nabla_{\mathbf{x}'} \log p(\mathbf{x}')\tag{17}$$

In our case, we have $p(\mathbf{x}_t|\bar{\alpha}_t \mathbf{x}_0) = \mathcal{N}(\bar{\alpha}_t \mathbf{x}_0, \bar{\sigma}_t^2 \mathbf{I})$ and we use EBM to model the marginal distribution of \mathbf{x}_t as $p_{\theta,t}(\mathbf{x}_t)$, thus:

$$\begin{aligned}\mathbb{E}(\bar{\alpha}_t \mathbf{x}_0|\mathbf{x}_t) &= \mathbf{x}_t + \bar{\sigma}_t^2 \nabla_{\mathbf{x}_t} \log p_{\theta,t}(\mathbf{x}_t) \\ \mathbb{E}(\mathbf{x}_0|\mathbf{x}_t) &= \frac{\mathbf{x}_t + \bar{\sigma}_t^2 \nabla_{\mathbf{x}_t} \log p_{\theta,t}(\mathbf{x}_t)}{\bar{\alpha}_t}\end{aligned}\tag{18}$$

Suppose the samples we get at $t = 0$ is x_0 . This samples actually contains a small magnitude of noise corresponding to $\bar{\alpha}_0$, thus, we may use Equation 18 to further denoise it. In practice, we find that enlarging the denoising step by multiplying gradient term $\nabla_{\mathbf{x}_t} \log p_{\theta,t}(\mathbf{x}_t)$ by a coefficient larger than 1.0 gives us better results. We set this coefficient to be 2.0 in our experiments.

A.3 NOISE SCHEDULE AND CONDITIONING INPUT

We improve upon the noise schedule and conditioning input of DRL. Let $\lambda_t = \log \frac{\bar{\alpha}_t^2}{\bar{\sigma}_t^2}$ denote the logarithm of signal-to-noise ratio at noise level t . Inspired by (Kingma et al., 2021), we send λ_t as the conditioning input to f_θ and \mathbf{g}_ϕ instead of t .

For the noise schedule, we keep the design of using 6 noise levels as in DRL. Inspired by (Nichol & Dhariwal, 2021), we construct a cosine schedule such that λ_t is defined as $\lambda_t = -2 \log(\tan(at + b))$.

a and b are calculated from the maximum log SNR (denoted as λ_{\max}) and minimum log SNR (denoted as λ_{\min}) using:

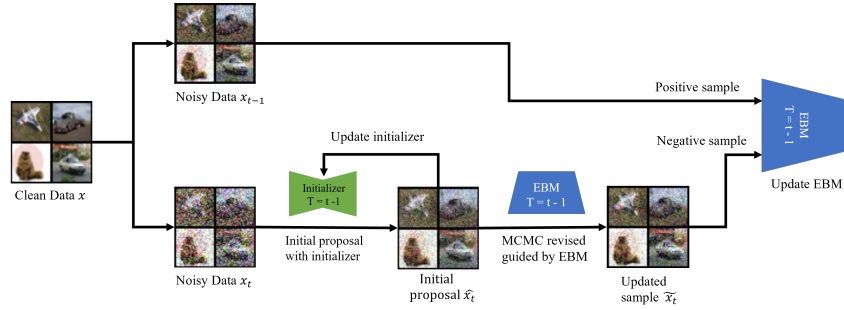
$$b = \arctan(\exp(-0.5\lambda_{\max})), \quad (19)$$

$$a = \arctan(\exp(-0.5\lambda_{\min})) - b. \quad (20)$$

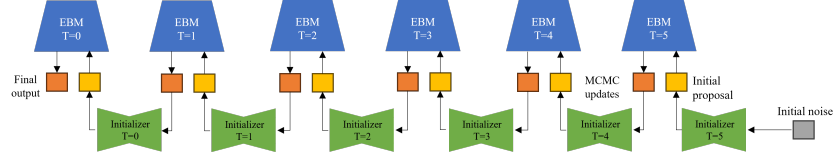
We set $\lambda_{\max} = 9.8$ and $\lambda_{\min} = -5.1$ to match the accumulative noise level $\bar{\alpha}_t$ of the original Recovery Likelihood model T6 at the highest and lowest noise levels. Figure 6 illustrates the noise schedule of DRL and our proposed schedule. Compared to DRL’s original schedule, our new schedule focuses more on regions with lower signal-to-noise ratios, crucial for generating low-frequency, high-level concepts in samples.

A.4 OVERALL ILLUSTRATION OF CDRL

In figure 5, we give an overview illustration of the training and sampling process of CDRL.



(a) An overview of the CDRL training process: In the training phase, we commence by selecting a pair of images at noise levels t and $t - 1$. The image at noise level t is then input into the initializer to produce an initial proposal. Subsequently, this initial proposal undergoes refinement through MCMC process guided by the underlying energy function. The enhanced sample derived from this process is utilized to update both the energy function and the initializer.



(b) An Overview of the CDRL Sampling process: The sampling phase starts from Gaussian noise. Beginning at the highest noise level, an initial proposal is generated by the initializer specific to that noise level. This is followed by refinement of the samples through MCMC sampling. This procedure is iteratively repeated, descending through progressively lower noise levels, until the lowest noise level is reached

Figure 5: Overview of CDRL

B RELATED WORK

Energy-Based Learning Energy-based models (EBMs) (Zhu et al., 1998; LeCun et al., 2006; Ngiam et al., 2011; Hinton, 2012) define unnormalized probabilistic distributions and are typically trained through maximum likelihood estimation. Methods such as contrastive divergence (Hinton, 2002; Du et al., 2021), persistent chain (Xie et al., 2016), replay buffer (Du & Mordatch, 2019) or short-run MCMC sampling (Nijkamp et al., 2019) approximate the analytically intractable learning gradient. To scale up and stabilize EBM training for high-fidelity data generation, strategies like multi-grid sampling (Gao et al., 2018), progressive training (Zhao et al., 2021), and diffusion (Gao et al., 2021) have been adopted.

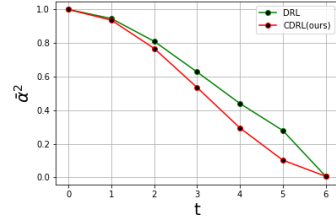


Figure 6: Noise level schedule. The green line is the noise levels used by DRL (Gao et al., 2021) while the red line is the noise levels used by our CDRL.

EBMs have also been connected to other models, such as adversarial training (Arbel et al., 2021; Che et al., 2020), variational autoencoders (Xiao et al., 2021), contrastive guidance (Lee et al., 2023), and noise contrastive estimation (Gao et al., 2020). To alleviate MCMC burden, various methods have been proposed, including amortizing MCMC sampling with learned networks (Kim & Bengio, 2016; Xie et al., 2018a; Kumar et al., 2019; Xiao et al., 2021; Han et al., 2019; Grathwohl et al., 2021a). Among them, cooperative networks (CoopNets) (Xie et al., 2018a) jointly train a top-down generator and an EBM via MCMC teaching, using the generator as a fast initializer for Langevin sampling. CoopNets variants have also been studied (Xie et al., 2021b; 2022). Our work improves the recovery likelihood learning algorithm of EBMs (Gao et al., 2021) by learning a fast MCMC initializer for EBM sampling, leveraging the cooperative learning scheme. Compared to (Xie et al., 2020) that applied cooperative training to an initializer and an EBM for the *marginal* distribution of the clean data, our approach only requires learning *conditional* initializers and sampling from *conditional* EBMs, which are much more tractable than their marginal counterparts.

Diffusion Model Diffusion models, originating from Sohl-Dickstein et al. (2015) and further developed in works such as Song & Ermon (2020); Ho et al. (2020), generate samples by progressively denoising them from a high noise level to clean data. These models have achieved remarkable success in generating high-quality samples from complex distributions, thanks to various architectural and framework innovations Ho et al. (2020); Song et al.; Kim et al. (2021); Song et al. (2021); Dhariwal & Nichol (2021); Karras et al. (2022); Ho & Salimans (2022). Notably, Dhariwal & Nichol (2021) emphasizes that the generative performance of diffusion models can be enhanced with the aid of a classifier, while Ho & Salimans (2022) further demonstrates that this guided scoring can be estimated by the differential scores of a conditional model versus an unconditional model. Enhancements in sampling speed have been realized through distillation techniques Salimans & Ho (2022) and the development of fast SDE/ODE samplers Song et al.; Karras et al. (2022); Lu et al. (2022a). Recent advancements Rombach et al. (2022); Saharia et al. (2022); Ramesh et al. (2022) have successfully applied conditional diffusion models to the task of text-to-image generation, achieving significant breakthroughs.

EBM shares a close relationship with diffusion models, as both frameworks can provide a score to guide the generation process, whether through Langevin dynamics or SDE/ODE solvers. As Salimans & Ho (2021) discusses, the distinction between the two lies in their implementation approach: EBMs model the log-likelihood directly, while diffusion models concentrate on the gradient of the log-likelihood. This distinction enables EBM to be used in some potential applications. These include utilizing advanced sampling techniques Du et al. (2023), transformed into classifiers Guo et al. (2023), or employed in the detection of abnormal samples through estimated likelihood Grathwohl et al. (2020); Liu et al. (2020).

The focus of this work is to push the development of EBM. And our work connects to diffusion models (Ho et al., 2020; Xiao et al., 2022) by learning a sequence of EBMs and MCMC initializers to reverse the diffusion process. Contrasting (Ho et al., 2020), our framework employs more expressive conditional EBMs instead of normal distributions. (Xiao et al., 2022) also suggests multimodal distributions, trained by generative adversarial networks (Goodfellow et al., 2020), for the reverse process.

C IMAGE INPAINTING

We demonstrate our learned model’s inpainting ability on the 64×64 CelebA dataset. For each image, we mask a portion and let the model fill in the masked area. We add noise to the masked image up to the final noise level and allow the model to gradually denoise the image, similar to the standard generation process. During inpainting, we only update the masked area, retaining the unmasked area’s values. This is achieved by resetting the unmasked area values to the current noisy version after each Langevin update step of the EBM or initializer proposal step. Our results are shown in Figure 7. We test two masking types: a regular square mask and an irregularly shaped mask. CDRL successfully inpaints valid and diverse values in the masked area, with inpainted results differing from the observations. This indicates that our model does not merely memorize data points but fills in meaningful unobserved areas based on the dataset’s statistical features.

D MORE EXPERIMENTAL RESULTS

D.1 MORE GENERATION RESULTS

In this section we show more generation results. Figure 10 shows more compositionality results with different guidance weight on CelebA 64×64 dataset. Here, $w = 0.0$ equals to the original setting in Equation 13 in the main paper without guidance. In Figure 11, 12, 13, 14 and 15 provide more results for conditional generation on ImageNet32 (32×32) with different guidance weight. Figure 11 gives random samples while each figure in Figure 12, 13, 14 and 15 contains samples from a certain class under different guidance weight w .

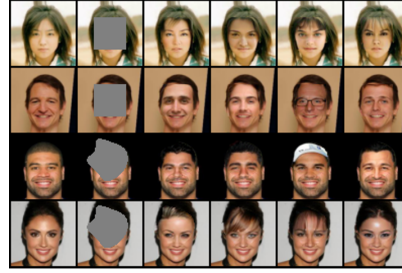


Figure 7: Image inpainting results on CelebA (64×64). The first two rows use square masks and the last two rows uses irregular masks. Columns one displays original images. Column two shows masked images. Columns three to six exhibit inpainted images using different initialization seeds.

D.2 GENERATE HIGH RESOLUTION IMAGE

The recent trend in generative modeling involves either utilizing the latent space of a VAE, as in latent Diffusion Rombach et al. (2022), or initially generating a low-resolution image and then scaling up, as demonstrated by techniques like Imagen Sahraria et al. (2022). This process might reduce the modeled space to dimensions such as 32×32 or 64×64 , aligning with the resolution on which we conducted our experiments in the main paper. Here, we undertook additional experiments by learning CDRL following Rombach et al. (2022). We conduct experiments on the CelebA-HQ dataset. The samples are shown in Figure 9. We report the FID score in Table 7.

Table 7: FID score for CelebA-HQ (256 x 256) dataset

Model	FID score
GLOW (Kingma & Dhariwal, 2018)	68.93
VAEBM (Xiao et al., 2021)	20.38
ATEBM (Yin et al., 2022)	17.31
VQGAN+Transformer (Esser et al., 2021)	10.2
LDM Rombach et al. (2022)	5.11
CDRL(ours)	10.74

D.3 MORE OOD RESULTS

Here we include the results for OOD task on more datasets. We also include more recent baselines. The results are included in Table 8.

E SAMPLING TIME

In this section, we measure the sampling time of CDRL and compare it with the following models: 1) CoopFlow (Xie et al., 2022), which composes a EBM and a Normalizing Flow model; 2) VAEBM (Xiao et al., 2021), which composes a VAE with an EBM and achieves strong generation performance; 3) The original DRL (Gao et al., 2021) model with 30 step MCMC steps at each noise level. We run the sampling process of each model individually on a single A6000 GPU to generate a batch of 100 samples on the Cifar10 dataset. Our CDRL model generates samples with better quality with relatively less time. And after applying the sampling adjustment techniques, the sampling time can be further reduced without hurting much sampling quality.

Table 8: AUROC scores in OOD detection using CDRL and other explicit density models on CIFAR-10. The score for DRL Gao et al. (2021) is reported by Yoon et al. (2023). Also for EBM-CD Du et al. (2021), we find different numbers from different sources, one from the Du et al. (2021) and the other from a recent work Yoon et al. (2023), we include both scores here and put the scores from Yoon et al. (2023) into brackets.

	Cifar-10 interpolation	Cifar-100	CelebA	SVHN	Texture
PixelCNN (Salimans et al., 2017)	0.71	0.63	-	0.32	0.33
GLOW (Kingma & Dhariwal, 2018)	0.51	0.55	0.57	0.24	0.27
NVAE (Vahdat & Kautz, 2020)	0.64	0.56	0.68	0.42	-
EBM-IG (Du & Mordatch, 2019)	0.70	0.50	0.70	0.63	0.48
VAEBM (Xiao et al., 2021)	0.70	0.62	0.77	0.83	-
EBM-CD (Du et al., 2021)	0.65 (-)	0.83 (0.53)	- (0.54)	0.91 (0.78)	0.88 (0.73)
CLEL (Lee et al., 2023)	0.72	0.72	0.77	0.98	0.94
DRL (Gao et al., 2021)	-	0.44	0.64	0.88	0.45
MPDR-S (Yoon et al., 2023)	-	0.56	0.73	0.99	0.66
MPDR_R (Yoon et al., 2023)	-	0.64	0.83	0.98	0.80
CDRL	0.75	0.78	0.84	0.82	0.65

Table 9: Comparison of different EBMs in terms of sampling time and number of MCMC steps. The sampling time are measured in second.

Method	Number of MCMC steps	Sampling Time	FID ↓
CoopFlow (Xie et al., 2022)	30	2.5	15.80
VAEBM (Xiao et al., 2021)	16	21.3	12.16
DRL (Gao et al., 2021)	$6 \times 30 = 180$	23.9	9.58
CDRL	$6 \times 15 = 90$	12.2	4.31
CDRL (8 steps)	$6 \times 8 = 48$	6.5	4.58
CDRL (5 steps)	$6 \times 5 = 30$	4.2	5.37
CDRL (3 steps)	$6 \times 3 = 18$	2.6	9.67

F UNDERSTAND THE ROLE OF INITIALIZER AND EBM

To further understand the roles of our initializer and EBM in image generation, we conduct two additional experiments using a pretrained CDRL model on the ImageNet Dataset (32×32). We evaluate two generation options: (a) images generated using only the initializer’s proposal, without the EBM’s Langevin Dynamics at each noise level, and (b) images generated with the full CDRL model, which includes the initializer’s proposal and 15-step Langevin updates at each time interval. As shown in Figure 8a and 8b, the initializer captures the object’s rough outline, while the Langevin updates on EBM fill in meaningful details. Furthermore, in Figure 8c, we display samples generated by fixing the initial noise image and sample noise of each initializer proposal step. The results reveal that images generated with the same initialization noises share basic elements while differing in details, highlighting the impact of both the initializer and Langevin sampling. The initializer provides a starting point, and the Langevin sampling process adds details.

G ABLATION STUDY

In this section, we carry out ablation studies to justify the choice of each component of our CDRL model. We use several experiments to answer the questions listed below.

G.1 WHAT ARE THE CONTRIBUTIONS OF EACH INDIVIDUAL DESIGN COMPONENT?

In our main paper, we have described the three main techniques that contribute most to our CDRL model. They are the new noise schedule design, cooperative training algorithm and noise variance

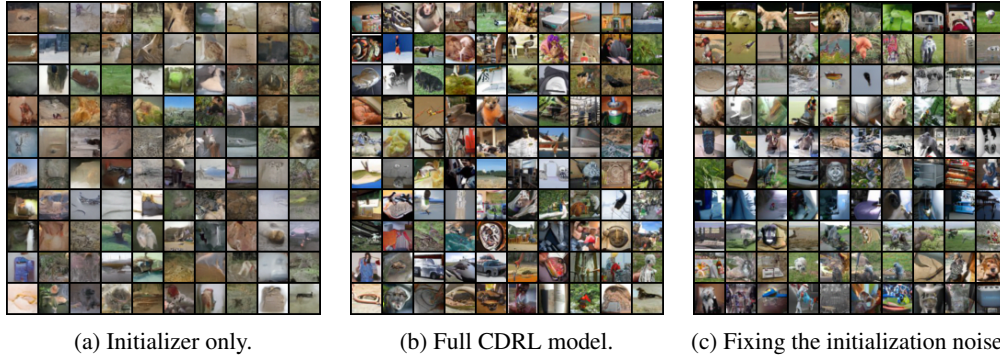


Figure 8: Illustration of the effect of initializer and EBM in the image generation process with a CDRL model pretrained on ImageNet Dataset (32×32). (a) Samples generated using only the proposal of the initializer; (b) Samples generated by the full CDRL model; (c) Samples generated by fixing the initial noise image and the sample noise of each initialization proposal step. Each row of images shared the same initial noise image and the sample noise of each initialization proposal step, but differed in the noises of Langevin sampling process at each noise level.

Table 10: Ablation study on CIFAR-10 dataset.

Models	FID ↓
DRL (Gao et al., 2021)	9.58
No cooperative training	6.47
No noise reduction	5.51
Initializer predicts $\hat{\mathbf{x}}_0$	5.17
Initializer predicts $\hat{\epsilon}$	4.95
CDRL	4.31

reduction. Here, we demonstrate the effect of each of the three techniques by comparing our CDRL with the following models (1) The original DRL model used in (Gao et al., 2021); (2) Model trained without cooperatively training framework: This corresponds to the original recovery-likelihood EBM using the same noise schedule and conditioning input as CDRL; (3) CDRL without noise reduction; (4) While CDRL use initializer to directly output $\hat{\mathbf{y}}_t$, here similar to (Xiao et al., 2022), we use the initializer to output the prediction of clean image $\hat{\mathbf{x}}_0$ and then transformed it to $\hat{\mathbf{y}}_t$. (5) Similarly, we can also use initializer to directly output the prediction of total added noise $\hat{\epsilon}$ similar to (Ho et al., 2020) and then transformed it to $\hat{\mathbf{y}}_t$.

We let all the models to share the same network structure and training setting on CIFAR-10 dataset and differ only in the way mentioned above. Shown in Table 10, our full model works the best among these settings, which justifies our design choice.

G.2 SHOULD THE INITIALIZER BE LEARNED USING THE COOPERATIVE TRAINING ALGORITHM OR SHOULD IT DIRECTLY REGRESS ON THE DATA?

Shown in equation 8 in the main paper, in our cooperative training algorithm, the initializer learns from the revised sample $\tilde{\mathbf{y}}_t$ at each step. A natural question to be asked here is whether we should instead regress it directly on data \mathbf{y}_t . To answer this, we try the option that the initializer directly learns from \mathbf{y}_t at each step. We denote this option as CDRL(data) to distinguish it from the original CDRL model. The results shown in Table 11 suggest that CDRL works better than CDRL(data), which supports our choice of the cooperative training algorithm.

To understand this, we may first dive into a deeper understanding of the learning behavior of the cooperative learning algorithm. We follow the analysis framework of (Nijkamp et al., 2019; Xie et al., 2022). Let $K_\theta(\mathbf{y}_t|\mathbf{y}_t', \mathbf{x}_{t+1})$ be the transition kernel of the K -step Langevin

sampling that refines the initial output \mathbf{y}_t' to the refined output \mathbf{y}_t . Let $(K_{\theta}q_{\phi})(\mathbf{y}_t|\mathbf{x}_{t+1}) = \int K_{\theta}(\mathbf{y}_t|\mathbf{y}_t', \mathbf{x}_{t+1})q_{\phi}(\mathbf{y}_t'|\mathbf{x}_{t+1})d\mathbf{y}_t'$ be the conditional distribution of \mathbf{y}_t , which is obtained by K steps of Langevin sampling starting from the output of the initializer $q_{\phi}(\mathbf{y}_t|\mathbf{x}_{t+1})$. Let $\pi(\mathbf{y}_t|\mathbf{x}_{t+1})$ be the true conditional distribution for denoising \mathbf{x}_{t+1} to retrieve \mathbf{y}_t . The maximum recovery likelihood for the EBM in equation 4 in the main paper is equivalent to minimizing the KL divergence $KL(\pi(\mathbf{x}_t|\mathbf{y}_{t+1})||p_{\theta}(\mathbf{x}_t|\mathbf{y}_{t+1}))$. Let us use j to index the learning iteration for model parameters. Given the current initializer model $q_{\phi}(\mathbf{y}_t|\mathbf{x}_{t+1})$, the EBM updates its parameters θ by minimizing

$$\theta_{j+1} = \arg \min_{\theta} KL(\pi(\mathbf{y}_t|\mathbf{x}_{t+1})||p_{\theta}(\mathbf{y}_t|\mathbf{x}_{t+1})) - KL((K_{\theta_j}q_{\phi})(\mathbf{y}_t|\mathbf{x}_{t+1})||p_{\theta}(\mathbf{y}_t|\mathbf{x}_{t+1})) \quad (21)$$

which is modified contrastive divergence. It is worth noting that, in the original contrastive divergence, the above $(K_{\theta_j}q_{\phi})(\mathbf{y}_t|\mathbf{x}_{t+1})$ is replaced by $(K_{\theta_j}\pi)(\mathbf{y}_t|\mathbf{x}_{t+1})$. That is, the MCMC chains are initialized by the true data. The learning shifts $p_{\theta}(\mathbf{y}_t|\mathbf{x}_{t+1})$ toward the true distribution $\pi(\mathbf{y}_t|\mathbf{x}_{t+1})$.

On the other hand, given the current EBM, the initializer model learns from the output distribution of the EBM’s MCMC. (That is, we train the initializer with $\tilde{\mathbf{y}}$ in equation 8 of the main paper.) The update of the parameters of the initializer at learning iteration $j + 1$ approximately follows the gradient of

$$\phi_{j+1} = \arg \min_{\phi} KL(K_{\theta}q_{\phi_j}(\mathbf{y}_t|\mathbf{x}_{t+1})||q_{\phi}(\mathbf{y}_t|\mathbf{x}_{t+1})) \quad (22)$$

The initializer $q_{\phi}(\mathbf{y}_t|\mathbf{x}_{t+1})$ learns to be the stationary distribution of the MCMC transition $K_{\theta}(\mathbf{x}_t|\mathbf{y}_{t+1})$ by shifting its mapping to the low energy regions of $p_{\theta}(\mathbf{x}_t|\mathbf{y}_{t+1})$. In a limit, the initializer $q_{\phi}(\mathbf{y}_t|\mathbf{x}_{t+1})$ minimizes $KL(K_{\theta}q_{\phi_j}(\mathbf{y}_t|\mathbf{x}_{t+1})||q_{\phi}(\mathbf{y}_t|\mathbf{x}_{t+1}))$ and gets close to the EBM $p_{\theta}(\mathbf{x}_t|\mathbf{y}_{t+1})$. The whole learning algorithm is a chasing game. That is, the initializer model $q_{\phi}(\mathbf{y}_t|\mathbf{x}_{t+1})$ chases the EBM $p_{\theta}(\mathbf{y}_t|\mathbf{x}_{t+1})$ toward the true distribution $\pi(\mathbf{y}_t|\mathbf{x}_{t+1})$.

Then we can begin the discussion of the benefits of the cooperative learning. Comparing with training the initializer directly with observed \mathbf{y} , the proposed training with $\tilde{\mathbf{y}}$ has the following benefits.

Firstly, in our algorithm, Equation 22 above shows that the MCMC of the EBM drives the evolution of the initializer that seeks to amortize the MCMC. At each learning iteration, in order to provide good initial examples for the current EBM’s MCMC, the initializer needs to be close enough to the EBM. Therefore, at each learning algorithm, training the initializer with the MCMC outputs $\tilde{\mathbf{y}}$ is a good strategy to maintain a proper distance between EBM and initializer model. If the initializer directly learns from the true distribution, even though it can move toward the true distribution quickly, it might not provide a good starting point for the MCMC. A good initializer should be helpful for finding the modes of the EBM. Let us imagine a situation in which the initializer model has shifted toward the true distribution by firstly learning directly with \mathbf{y} , but the EBM is still far from that. Due to a large divergence between EBM and initializer, the initializer is not helpful for EBM to draw fair samples, especially with a finite-step Langevin dynamics. A far-away initializer might lead to unstable training of the EBM.

Secondly, to consider a more generic case, in which we model our initializer via a non-Gaussian generator $\mathbf{y}_t = g_{\phi}(\mathbf{x}_{t+1}, \mathbf{z}, t)$, $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$. The randomness comes from the latent vector \mathbf{z} . In this case $q_{\phi}(\mathbf{y}_t|\mathbf{x}_{t+1}) = \int p(\mathbf{y}_t|\mathbf{x}_{t+1}, \mathbf{z})p(\mathbf{z})d\mathbf{z}$ is analytically intractable. Learning $q_{\phi}(\mathbf{y}_t|\mathbf{x}_{t+1})$ directly from \mathbf{y} independently requires MCMC inference for the posterior distribution $q_{\phi}(\mathbf{z}|\mathbf{x}_{t+1}, \mathbf{y}_t)$. However, the cooperative learning can get around the difficulty of inference of the latent variables \mathbf{z} . That is, at each learning iteration, we generate examples $\tilde{\mathbf{y}}$ from $q_{\phi}(\mathbf{y}_t|\mathbf{x}_{t+1})$ by first sampling $\hat{\mathbf{z}} \sim p(\mathbf{z})$ and then mapping it to $\tilde{\mathbf{y}}_t = g_{\phi}(\mathbf{x}_{t+1}, \hat{\mathbf{z}}, t)$. The $\tilde{\mathbf{y}}$ is used to initialize the EBM’s MCMC that produces $\tilde{\mathbf{y}}$. The learning equation of ϕ is $\frac{1}{n} \sum_{i=1}^n -\frac{1}{2\hat{\sigma}_t^2} \|\tilde{\mathbf{y}}_{t,i} - g_{\phi}(\mathbf{x}_{t+1,i}, \hat{\mathbf{z}}, t)\|^2$, where the latent variables $\hat{\mathbf{z}}$ is used. That is, we shift the mapping from $\hat{\mathbf{z}} \rightarrow \tilde{\mathbf{y}}$ to $\hat{\mathbf{z}} \rightarrow \tilde{\mathbf{y}}$ for accumulate the MCMC transition. Even though in our

Table 11: Comparison between the model whose initializer learns from the samples given by EBM and the one whose initializer learns directly from the data. Scores are reported on Cifar-10 dataset.

Model	FID ↓
CDRL(data)	5.95
CDRL	4.31

Table 12: Comparison between our CDRL model with the model using DRL noise level schedule but adding 2 more noise levels to the high noise region. Scores are reported on Cifar-10 dataset.

Model	FID ↓
CDRL(DRL-T8)	4.94
CDRL	4.31

paper, we currently use a Gaussian initializer, if we use a non-Gaussian initializer in the future, the current cooperative learning (i.e., training q_ϕ with \tilde{y}) can be much more beneficial and feasible.

G.3 WHETHER ADDING MORE DIFFUSION LEVELS TO THE ORIGIN DRL SCHEDULE TO COVER THE HIGH-NOISE REGION WORKS BETTER THAN THE NEW SCHEDULE USED IN CDRL?

In section A.3, we introduce the new noise schedule we used for CDRL. Comparing with the one used in the original DRL(Gao et al., 2021) paper, the noise schedule used in CDRL puts more attention on the high-noise area where $\bar{\alpha}$ is close to 0. We carry out a comparison that trains the CDRL model using the original DRL schedule but with 2 more extra noise levels in the high-noise region. We refer to this setting as CDRL(DRL-T8). As shown in Table 12, CDRL(DRL-T8) performs slightly worse than CDRL. Also, in terms of sampling time, CDRL(DRL-T8) required 30% more steps during sampling. That is, increasing the number of noise levels might not necessarily improve the performance, but it must increase the computational cost. Our new schedule is very important because it keeps the algorithm efficient.

G.4 CAN WE FURTHER REDUCE THE NUMBER OF NOISE LEVELS?

We test whether the noise level can be further reduced. The results in Table 13a show that further reducing noise level to 4 can make model more unstable even if we increase the number of the Langevin sample steps K . On the other hand reducing T to 5 gives reasonable but slightly worse results.

G.5 WHAT IS THE EFFECT OF USING DIFFERENT NUMBER OF LANGEVIN STEPS?

In Table 13b, we show the effect of changing the number of Langevin steps K . The results show that, on one hand, decreasing K to 10 gives us comparable but slightly worse results. On the other hand, increasing K to 30 doesn't give better results. This agrees with the observations by DRL(Gao et al., 2021). The observation of changing K implies that simply increasing the number of Langevin steps doesn't bring significant increase in the sample quality, which verifies the effectiveness of the initializer in our model.

Table 13: CDRL with different number of noise levels T and number of Langevin steps K . Scores are reported on Cifar-10 dataset.

(a) Results for reducing T		(b) Results for changing K	
Model	FID ↓	Model	FID ↓
T=4 (K=15,20,30)	- (not converge)	T=6 (K=10)	4.50
T=5 (K=15)	5.08	T=6 (K=15)	4.31
T=6 (K=15)	4.31	T=6 (K=30)	5.08



Figure 9: Samples on CelebAHQ (256×256)

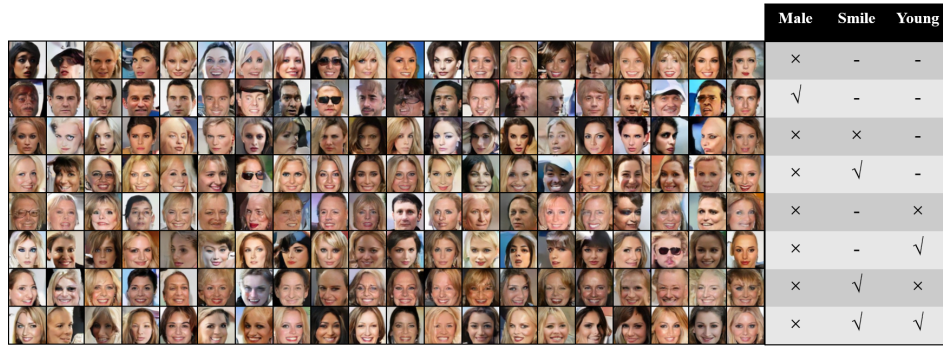
(a) $w = 0.0$ (b) $w = 0.5$ (c) $w = 1.0$

Figure 10: Attribute compositional samples on CelebA (64×64). Here we use guided weight $w = 0.0, 0.5, 1.0$. We let images at different guidance share the same random noise. Results can also be compared with Figure 4 which use $w = 3.0$

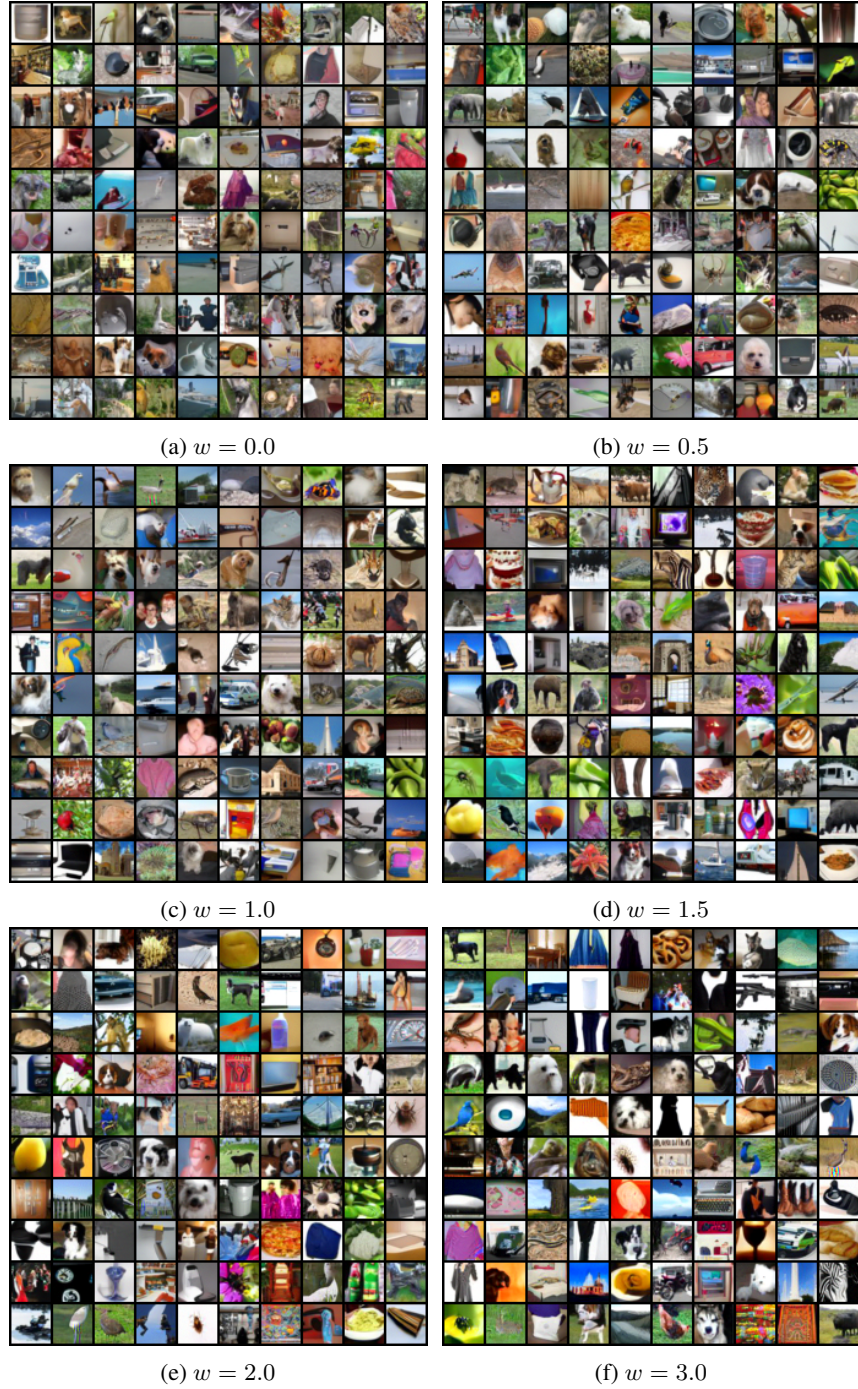


Figure 11: Conditional generated examples with different classifier free guidance weight on ImageNet32 (32×32). Samples are generated with randomly chosen class label.



Figure 12: Conditional generated examples with different classifier free guidance weight on ImageNet32 (32×32) with class label Tench.

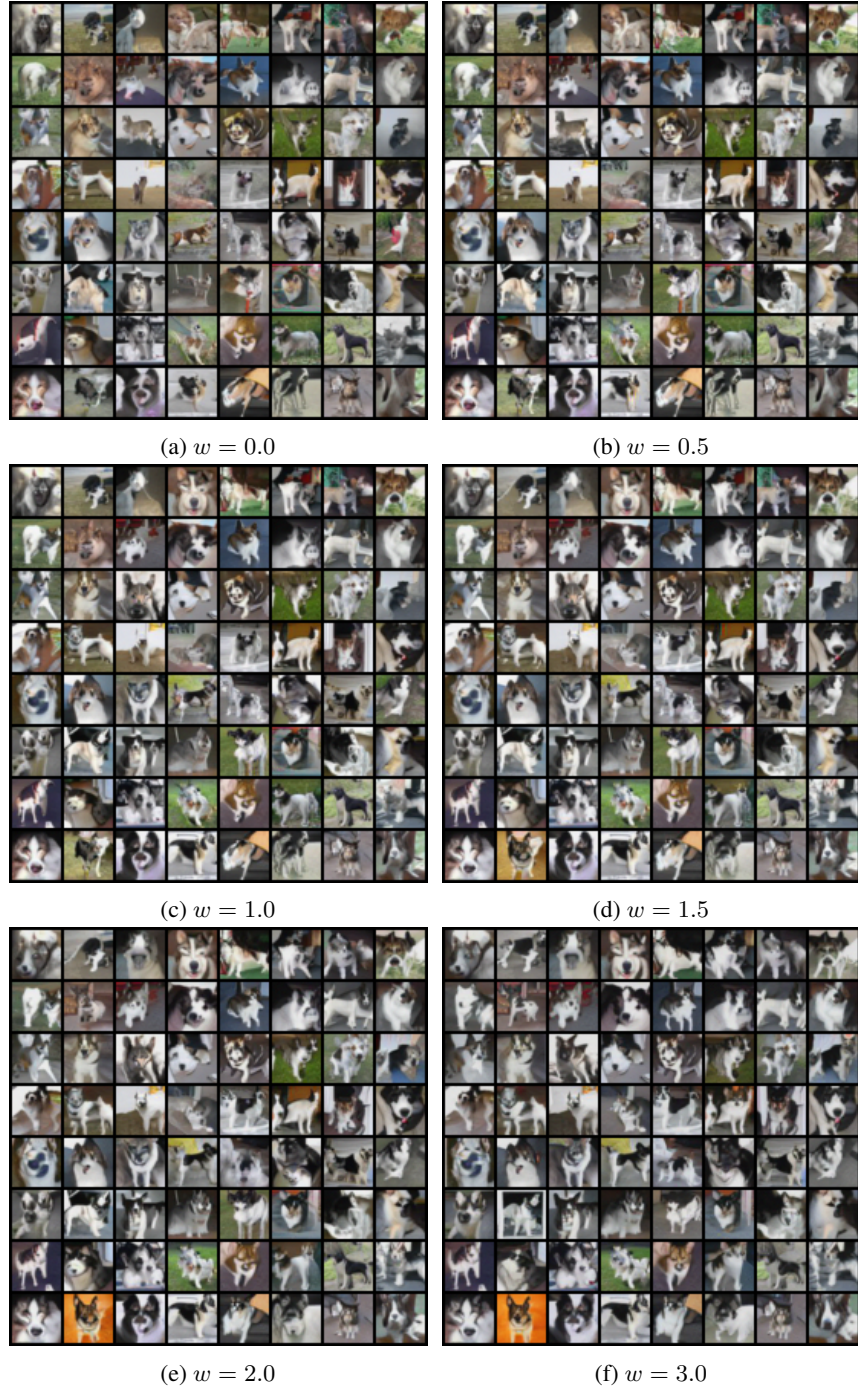


Figure 13: Conditional generated examples with different classifier free guidance weight on ImageNet32 (32×32) with class label Siberian Husky.



Figure 14: Conditional generated examples with different classifier free guidance weight on ImageNet32 (32×32) with class label Tow Truck.

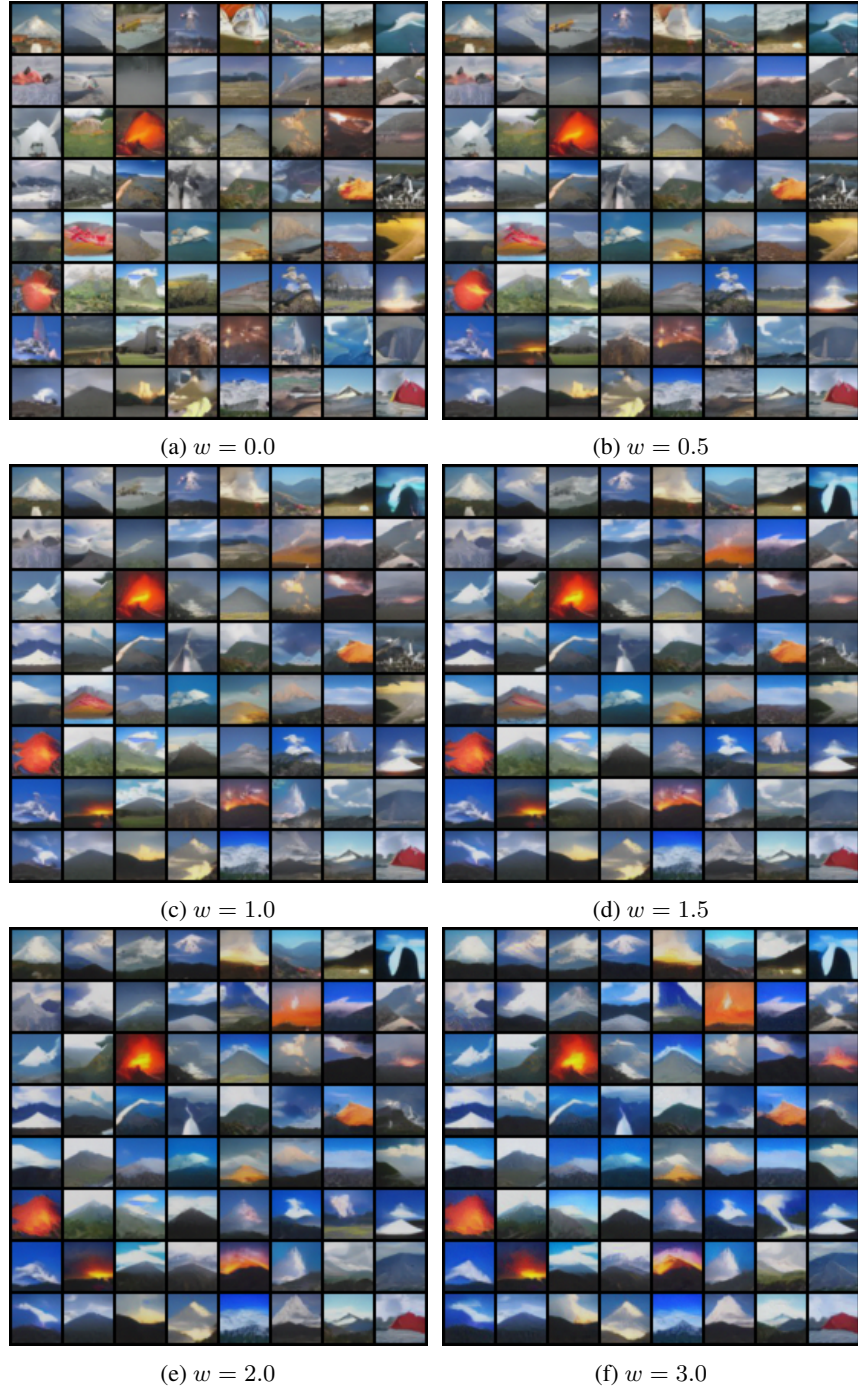


Figure 15: Conditional generated examples with different classifier free guidance weight on ImageNet32 (32×32) with class label Volcano.