

Generalized Contrastive Divergence: Joint Training of Energy-Based Model and Diffusion Model through Inverse Reinforcement Learning

Sangwoong Yoon¹, Dohyun Kwon², Himchan Hwang³, Yung-Kyun Noh^{1,4}, Frank C. Park^{3,5}

¹Korea Institute for Advanced Study, ²University of Seoul,

³Seoul National University, ⁴Hanyang University, ⁵Saige Research

swyoon@kias.re.kr, dhkwon@uos.ac.kr,

himchan@robotics.snu.ac.kr, nohyung@hanyang.ac.kr, fcp@snu.ac.kr

Abstract

We present Generalized Contrastive Divergence (GCD), a novel objective function for training an energy-based model (EBM) and a sampler simultaneously. GCD generalizes Contrastive Divergence [1], a celebrated algorithm for training EBM, by replacing Markov Chain Monte Carlo (MCMC) distribution with a trainable sampler, such as a diffusion model. In GCD, the joint training of EBM and a diffusion model is formulated as a minimax problem, which reaches an equilibrium when both models converge to the data distribution. The minimax learning with GCD bears interesting equivalence to inverse reinforcement learning, where the energy corresponds to a negative reward, the diffusion model is a policy, and the real data is expert demonstrations. We present preliminary yet promising results showing that joint training is beneficial for both EBM and a diffusion model. GCD enables EBM training without MCMC while improving the sample quality of a diffusion model.

1 Introduction

Diffusion models [2, 3] have achieved tremendous success in generating realistic high-dimensional samples. However, several challenges are remaining. First, a diffusion model does not directly minimize the deviation between the model and the data distributions. While recent theories show that diffusion models indirectly minimize a divergence between the model and data [4, 5, 6], a direct minimization can be more effective, if possible [7]. Second, generation is slow, requiring a large number of denoising steps. Shortening the number of steps often incurs a trade-off in sample quality [7, 8]. Third, the marginal likelihood of data can not be directly computed, and only its lower bound or importance-weighted approximation is available [2, 3].

In this paper, we propose Generalized Contrastive Divergence (GCD), a novel objective function motivated by Contrastive Divergence [1], that shows training a diffusion model in conjunction with an energy-based model (EBM) helps tackle these problems. First, learning with GCD is equivalent to the minimization of integral probability metric (IPM; e.g., Wasserstein distance) between a diffusion model and data, where the energy works as a critic in IPM. Second, therefore, a diffusion model with fewer steps can be fine-tuned with GCD to improve sample quality. Third, the jointly trained EBM can compute the marginal likelihood of data (up to constant).

This paper presents preliminary experimental results showing that GCD learning works successfully on a synthetic dataset and can be effective in improving the sample quality of the DDPM sampler, especially when the number of steps is small. We expect to obtain experimental results on larger real-world datasets in the near future. Furthermore, we discuss the connection between GCD learning

and inverse reinforcement learning. This connection opens up potential opportunities for advancing generative modeling by leveraging the recent ideas of reinforcement learning.

The contributions of this paper can be summarized as follows:

- We propose GCD learning, a novel method for jointly training EBM and a parametric sampler. The joint training is beneficial for both models.
 - GCD learning enables EBM training without MCMC, which is often computationally expensive and unstable.
 - GCD learning can be used to train a diffusion-style sampler with a small number of timesteps.
- We show that GCD learning has interesting equivalence to the integral probability metric minimization problems (e.g., WGAN) and inverse reinforcement learning.

2 Preliminaries

Diffusion for Sampling. In this paper, we focus on discrete-time diffusion-based samplers, such as DDPM [2]. A sampler first draws a Gaussian noise vector \mathbf{x}_0 and then iteratively samples from conditional distributions $\mathbf{x}_{t+1} \sim \pi(\mathbf{x}_{t+1}|\mathbf{x}_t, t)$ to obtain a final sample \mathbf{x}_T . We will denote the final samples as $\mathbf{x} = \mathbf{x}_T$ and their marginal distribution as $\pi(\mathbf{x})$. The initial samples \mathbf{x}_0 are drawn from $\mathcal{N}(0, I)$. The sampler $\pi(\mathbf{x})$ may be pre-trained through existing diffusion model techniques.

Energy-Based Models. An energy-based model (EBM) $q(\mathbf{x})$ uses a scalar function called an energy $E(\mathbf{x})$ to represent a probability distribution:

$$q(\mathbf{x}) = \frac{1}{Z} \exp(-E(\mathbf{x})/\tau), \quad E: \mathcal{X} \rightarrow \mathbb{R} \quad Z = \int_{\mathcal{X}} \exp(-E(\mathbf{x})/\tau) d\mathbf{x}, \quad (1)$$

where $\tau > 0$ is temperature, \mathcal{X} is the compact domain of data, and the normalization constant Z is assumed to be finite. The integral is replaced with a summation for discrete input data.

While multiple training algorithms for EBM have been proposed (see [9] for an overview), the standard method is maximum likelihood training which is equivalent to minimization of Kullback-Leibler (KL) divergence between data and a model: $\min_{q \in \mathcal{Q}} KL(p||q)$, where $p(\mathbf{x})$ is the data distribution, and \mathcal{Q} being the set of feasible $q(\mathbf{x})$'s. KL divergence is defined as $KL(p||q) = \int_{\mathcal{X}} p(\mathbf{x}) \log(p(\mathbf{x})/q(\mathbf{x})) d\mathbf{x}$. Meanwhile, the exact computation of the log-likelihood gradient requires convergent MCMC sampling.

Contrastive Divergence [1]. Hinton proposed an alternative objective function for training EBM which does not require convergent MCMC. Let us write \mathcal{T}_q as a single-step MCMC operator designed to draw sample from $q(\mathbf{x})$ and $\mathcal{T}_q^k(p)$ as the distribution of points after k steps of MCMC transition \mathcal{T}_q starting from $p(\mathbf{x})$. Then, Contrastive Divergence (CD) learning is defined as follows:

$$\min_{q \in \mathcal{Q}} KL(p||q) - KL(\mathcal{T}_q^k(p)||q). \quad (2)$$

CD learning is known to converge when $KL(p||q) \geq KL(\mathcal{T}_q^k(p)||q)$ is guaranteed [10]. However, CD still requires convergent MCMC when drawing a new sample from $q(\mathbf{x})$.

3 Generalized Contrastive Divergence

3.1 Generalized Contrastive Divergence Learning

We present Generalized Contrastive Divergence (GCD), an objective function for the joint training of EBM and a sampler. From CD (Eq. 2), we replace MCMC distribution $\mathcal{T}_q^k(p)(\mathbf{x})$ with an arbitrary sampler $\pi(\mathbf{x})$ which will later be set as a DDPM-like discrete-time diffusion model. We also introduce maximization with respect to $\pi(\mathbf{x})$, motivated by the fact that $\mathcal{T}_q^k(p)(\mathbf{x})$ converges towards $q(\mathbf{x})$.

Definition (GCD). Suppose the data distribution $p(\mathbf{x})$, the EBM $q(\mathbf{x})$, and the sampler $\pi(\mathbf{x})$ are supported on the domain \mathcal{X} . Let \mathcal{Q} and Π be the sets of feasible $q(\mathbf{x})$'s and $\pi(\mathbf{x})$'s, respectively. **GCD learning** is defined as the following minimax problem:

$$\inf_{q \in \mathcal{Q}} \sup_{\pi \in \Pi} KL(p||q) - KL(\pi||q), \quad (3)$$

where the changes made from CD are highlighted. This minimax problem has a desirable equilibrium.

Proposition 1 (Equilibrium). *Suppose $p(\mathbf{x}) \in \mathcal{Q} = \Pi$. Then, $p(\mathbf{x}) = q(\mathbf{x}) = \pi(\mathbf{x})$ is the equilibrium of GCD learning.*

Eq. 3 can be rewritten with respect to energy. Plugging EBM’s definition (Eq. 1), we obtain the following equivalent problem, which we solve in practice:

$$\inf_{E \in \mathcal{E}} \sup_{\pi \in \Pi} \mathbb{E}_p[E(\mathbf{x})] - \mathbb{E}_\pi[E(\mathbf{x})] + \tau \mathcal{H}(\pi), \quad (4)$$

where $\mathcal{H}(\pi) = -\int \pi(\mathbf{x}) \log \pi(\mathbf{x}) d\mathbf{x}$ is the differential entropy of $\pi(\mathbf{x})$. The set of energy functions \mathcal{E} is derived from \mathcal{Q} :

$$\mathcal{E} := \left\{ E : \mathcal{X} \rightarrow \mathbb{R} \mid \exp(-E(\mathbf{x})/\tau)/Z \in \mathcal{Q} \text{ where } Z = \int \exp(-E(\mathbf{x})/\tau) d\mathbf{x} < \infty \right\}. \quad (5)$$

In other words, $E \in \mathcal{E}$ if and only if $E(\mathbf{x}) = -\tau \log q(\mathbf{x}) + c$ for some $q \in \mathcal{Q}$ and a constant c . Similarly to Eq. 3, GCD learning in energy (Eq. 4) also reaches the equilibrium when $E(\mathbf{x}) = -\tau \log p(\mathbf{x}) + c$ and $\pi(\mathbf{x}) = p(\mathbf{x})$ (see Appendix A.2). GCD learning has interesting theoretical connections to other generative modeling problems and inverse reinforcement learning.

3.2 Equivalence to Entropy-Regularized Integral Probability Metric Minimization

Generative modeling is often formulated as the minimization of Integral Probability Metric (IPM; [11]), which measures the deviation between the data distribution and the model. WGAN [12] is a well-known example. GCD learning (Eq. 4) can also be viewed as the minimization of IPM between $p(\mathbf{x})$ and $\pi(\mathbf{x})$ but under entropy regularization for $\pi(\mathbf{x})$. The energy $E(\mathbf{x})$ works as a critic in IPM, and the feasible set \mathcal{E} characterizes IPM (hence we write IPM as $D^\mathcal{E}(p||\pi)$).

Proposition 2 (Entropy-regularized IPM minimization). *Assume that $p(\mathbf{x}) \in \mathcal{Q} \cap \Pi$. In addition, suppose that \mathcal{E} is closed under negation. Consider the following problem:*

$$\inf_{\pi \in \Pi} D^\mathcal{E}(p||\pi) - \tau \mathcal{H}(\pi), \text{ where } D^\mathcal{E}(p||\pi) := \sup_{E \in \mathcal{E}} |\mathbb{E}_p[E(\mathbf{x})] - \mathbb{E}_\pi[E(\mathbf{x})]|. \quad (6)$$

This problem has the same optimal value and the same equilibrium point ($E(\mathbf{x}) = -\tau \log p(\mathbf{x}) + c$ and $\pi(\mathbf{x}) = p(\mathbf{x})$) to the GCD learning in energy (Eq. 4).

Proposition 2 holds because Eq. 4 and Eq. 6 are the primal and the dual problems of the same objective function, where total duality holds. A detailed proof is in Appendix A.2.

In conventional IPM minimization problems, such as WGAN, where $\tau = 0$ and \mathcal{E} being the set of 1-Lipschitz functions, the critic is not guaranteed to converge to a quantity related to $\log p(\mathbf{x})$. The entropy regularization is critical for obtaining an accurate energy estimate and is the key difference from a prior work [7], where DDPM is optimized for only IPM.

3.3 Equivalence to Maximum Entropy Inverse Reinforcement Learning

GCD learning can also be interpreted through the lens of reinforcement learning (RL). Eq. 4 is a special case of maximum entropy inverse reinforcement learning (IRL; [13, 14, 15]). The connection between diffusion modeling and Markov Diffusion Processes (MDP) is well noted in previous works [7, 16]. The sampler $\pi(\mathbf{x})$ is an agent with a stochastic policy $\pi(\mathbf{x}_{t+1}|\mathbf{x}_t, t)$. The intermediate samples $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_T (= \mathbf{x})$, combined with time indices, form a trajectory of states $\{(\mathbf{x}_t, t)\}_{t=0}^T$. The action corresponds to the choice of the next state.

However, as no explicit reward signal is available in generative modeling, the reward must be inferred from the data, making the problem one of IRL. A training data \mathbf{x} serves as expert demonstrations for the terminal state \mathbf{x}_T . There is no demonstration for intermediate states, and thus the reward is only given for the terminal state. Then, a natural choice of the reward for a sampler is $\log p(\mathbf{x})$, which is unknown. In GCD learning, the energy function is responsible for learning $\log p(\mathbf{x})$, providing the reward signal $-E(\mathbf{x})$ for training $\pi(\mathbf{x})$.

One notable features of GCD is that it maximizes the terminal state entropy $\mathcal{H}(\pi(\mathbf{x}_T))$ instead of the causal entropy $\mathcal{H}(\pi(\mathbf{x}_{t+1}|\mathbf{x}_t, t))$, as is done in other max-ent IRL methods [13, 14, 15]. The connection to IRL allows us to analyze the training from RL perspective and to employ techniques proven to be effective in RL.

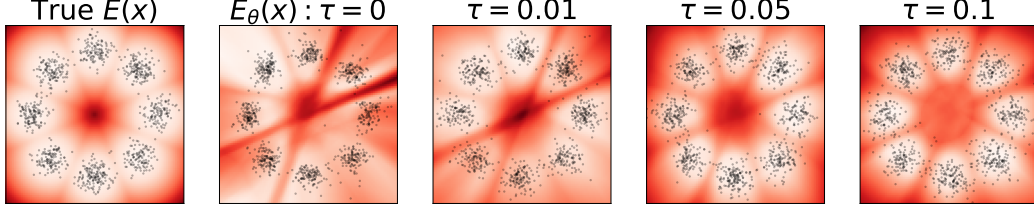


Figure 1: 2D Density Estimation on 8 Gaussians. The red shades indicate the energy (white is low), and the dots are samples. Without entropy regularization ($\tau = 0$), the energy does not reflect the data distribution, and the samples are collapsed to the mode.

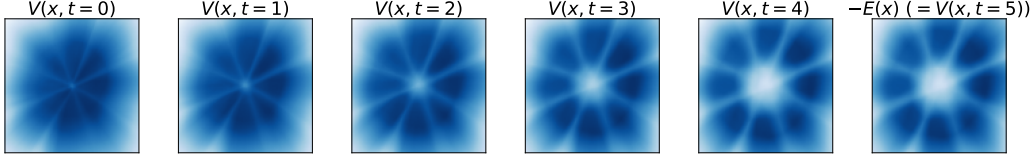


Figure 2: Value functions at each time step ($\tau = 0.1$ case). Blue indicates a high value.

4 Joint Training of EBM and Diffusion Model

Now we present an algorithm for GCD learning. While GCD learning is not restricted on the choice of $\pi(\mathbf{x})$, we shall focus on the case where $\pi(\mathbf{x})$ is a DDPM-style sampler. We write θ and ϕ as the parameters of EBM $q_\theta(\mathbf{x})$ and a sampler $\pi_\phi(\mathbf{x})$, respectively. GCD learning can be written as follows:

$$\min_{\theta} \max_{\phi} \mathcal{L}(\theta, \phi), \quad \mathcal{L} = \mathbb{E}_p[E_\theta(\mathbf{x})] - \mathbb{E}_{\pi_\phi}[E_\theta(\mathbf{x})] + \tau \mathcal{H}(\pi_\phi). \quad (7)$$

To solve this minimax problem, we alternatively update the EBM $E_\theta(\mathbf{x})$ and the diffusion model $\pi_\phi(\mathbf{x})$ as typically done in gradient descent-ascent algorithms, e.g., [17]. Temperature τ is treated as a hyperparameter.

EBM Update. The energy is updated to discriminate the training data and the samples from $\pi_\phi(\mathbf{x})$.

$$\nabla_{\theta} \mathcal{L} = \mathbb{E}_{p(\mathbf{x})}[\nabla_{\theta} E_{\theta}(\mathbf{x})] - \mathbb{E}_{\pi_{\phi}(\mathbf{x})}[\nabla_{\theta} E_{\theta}(\mathbf{x})]. \quad (8)$$

This update corresponds to the reward learning step in IRL [13, 15].

Diffusion Model Update. The sampler is updated to maximize the reward $-E_\theta(\mathbf{x})$ under entropy regularization. When a sampler admits a reparametrization trick, we may directly take the gradient of $-E_\theta(\mathbf{x})$. However, the reparametrization trick may not be applicable to a diffusion model due to the memory overhead. Instead, we can employ REINFORCE-style policy gradient algorithm [18] (Derivation in Appendix C):

$$\begin{aligned} \nabla_{\phi} \mathcal{L} &= \nabla_{\phi} (\mathbb{E}_{\pi_{\phi}(\mathbf{x})}[-E_{\theta}(\mathbf{x})] + \tau \cdot \mathcal{H}(\pi_{\phi})) = \mathbb{E}_{\pi_{\phi}(\mathbf{x})} [\nabla_{\phi} \log \pi_{\phi}(\mathbf{x}) (-E_{\theta}(\mathbf{x}) - \tau \log \pi_{\phi}(\mathbf{x}))] \\ &= \mathbb{E}_{\pi_{\phi}(\mathbf{x}_{0:T})} \left[\sum_{t=0}^{T-1} \nabla_{\phi} \log \pi_{\phi}(\mathbf{x}_{t+1} | \mathbf{x}_t, t) [-E_{\theta}(\mathbf{x}_T) - \tau \log \pi_{\phi}(\mathbf{x}_T) - B(\mathbf{x}_t, t)] \right], \end{aligned} \quad (9)$$

where the term $-\tau \log \pi_{\phi}(\mathbf{x}_T)$ can be seen as an additional reward for exploration. This update corresponds to the forward RL in typical IRL. For sample efficiency, we use proximal policy optimization [19] to update the parameters multiple times for training data. The baseline function $B(\mathbf{x}_t, t)$ is introduced to reduce the variance [20] and is described shortly.

Marginal Probability Estimation. The log probability $\log \pi_{\phi}(\mathbf{x}_T)$ in Eq. 9 can not be directly computed in DDPM. We estimate $\log \pi_{\phi}(\mathbf{x}_T)$ using mini-batch samples from $\pi_{\phi}(\mathbf{x})$. Following [21], we employ the nearest-neighbor density estimator [22], which gives $\log \pi_{\phi}(\mathbf{x}) \approx D \log d_k(\mathbf{x}) + c$ where D is the dimensionality of data, $d_k(\mathbf{x})$ is the distance to k -nearest neighbor, and c is constant with respect to \mathbf{x} .

Baseline Estimation. We model the baseline of $-E_{\theta}(\mathbf{x})$ and $-\tau \log \pi_{\phi}(\mathbf{x}_T)$ separately. For energy, we employ the state value function as the baseline $V(\mathbf{x}_t, t) = \mathbb{E}_{\pi_{\phi}}[-E_{\theta}(\mathbf{x}_T) | \mathbf{x}_t]$. We parametrize

$V(\mathbf{x}, t)$ as a separate model and estimate its parameter by minimizing temporal difference loss $\min_{V(\cdot, t)} (V(\mathbf{x}_{t+1}, t+1) - V(\mathbf{x}_t, t))^2$ for $t = 0, \dots, T-1$. Meanwhile, we simply use the estimated entropy as the baseline $\tau \mathcal{H}(\pi_\phi) = -\tau \mathbb{E}_{\pi_\phi} [\log \pi_\phi(\mathbf{x}_T)]$, where the expectation is approximated using exponential moving average. As a result, the baseline is given as $B(\mathbf{x}_t, t) = V(\mathbf{x}_t, t) + \tau \mathcal{H}(\pi_\phi)$.

5 Experiments

In our experiment, we show the synergistic benefit of jointly training EBM and a diffusion model on 2D 8 Gaussians data. We use time-conditioned MLP networks for both a value network and a policy network (i.e., sampler). The time step is encoded into a 128D vector using sinusoidal positional embedding and concatenated into the hidden neurons of MLP. The last time step ($t = T$) of the value network is treated as the negative energy. The conditional distributions $\pi(\mathbf{x}_{t+1}|\mathbf{x}_t, t)$ in the sampler are Gaussians where the mean is determined by the policy network. We use the linear variance schedule from 10^{-5} to 10^{-2} . Throughout the experiment, we use $T = 5$ for GCD.

We first train both EBM and a sampler from scratch (Figure 1). In the conventional IPM minimization ($\tau = 0$), the critic, i.e., the energy, deviates significantly from the true energy of $p(\mathbf{x})$. Furthermore, the samples tend to collapse to the mode of the data density, under-representing the variance of data. The entropy regularization makes the energy produce $q(\mathbf{x})$ that accurately reflects $p(\mathbf{x})$. Meanwhile, the value functions reflect the evolution of the energy from the initial Gaussian density (Figure 2), similar to Diffusion Recovery Likelihood [23].

Next, we show GCD can improve the sample quality of pre-trained DDPM samplers (Table 1). DDPM samplers with varying T are trained on 8 Gaussians until convergence, and Wasserstein distances from their samples to real data are measured. The quality of samples deteriorates significantly as T becomes small. However, if we further fine-tune DDPM using GCD learning (GCD-FT), we obtain a Wasserstein distance even smaller than that of $T = 1000$ DDPM. This result is possible because GCD learning directly minimizes IPM between the samples and data. Similar results were reported in [7]. If we do not fine-tune the sampler (GCD-Scratch), the distance slightly increases, showing there is a gain from DDPM pretraining. It is interesting to note that GCD learning with entropy regularization ($\tau = 0.05$), outperforms conventional IPM minimization ($\tau = 0$). One hypothesis is that the entropy term facilitates the exploration and thus leads to a better policy.

6 Related Work

GCD is an attempt to train generative samplers using RL. When applying RL, a key design choice is defining the reward signal. One source of reward is human feedback [24, 25], which became popular after being applied in large language models [26]. However, since human feedback is typically costly, a machine learning model, such as an image aesthetic quality estimator, can be used as a substitute for the reward function [16, 27]. Assuming the absence of human feedback or an external reward function, GCD employs the IRL approach where the reward function is inferred from demonstrations. In GCD, the reward is the log-likelihood of data, learned from training data through EBM. The previous work of Fan and Lee [7] can also be seen as IRL, but their reward function is not interpreted as the log-likelihood.

GCD can be seen as a novel method for EBM training that does not require MCMC. EBMs are a powerful class of generative models that shows promising performance in compositional generation [28] and out-of-distribution detection [29, 30]. However, EBM suffers from instability in training,

Method	T	τ	$SW (\downarrow)$	AUC (\uparrow)
DDPM	5	-	0.967 ± 0.005	-
DDPM	10	-	0.824 ± 0.002	-
DDPM	100	-	0.241 ± 0.003	-
DDPM	1000	-	0.123 ± 0.014	-
GCD-Scratch	5	0.05	0.114 ± 0.025	0.867
GCD-FT	5	0.05	0.086 ± 0.008	0.880
GCD-FT	5	0	0.152 ± 0.008	0.513

Table 1: Fine-tuning DDPM for 8 Gaussians. SW indicates sliced Wasserstein-2 distance (1,000 projections) between samples and data estimated from 10k samples. The standard deviation is computed from 5 independent samplings. AUC indicates how well the energy discriminates the uniform distribution over the domain from data. AUC is computed with 10k samples. The maximum value of AUC vs the uniform distribution is about 0.9059.

often requiring computationally expensive MCMC that is difficult to tune and not convergent in practice [31, 21, 32, 23, 33]. GCD enables the stable training of EBM by removing MCMC from training.

GCD bears a formal resemblance to algorithms that train EBM using minimax formulation. CoopFlow [34] jointly trains a normalizing flow sampler and EBM, but the flow is not directly optimized to minimize the divergence to EBM. Divergence Triangle [35] trains a sampler, an encoder, and EBM simultaneously, resulting in a more complex problem than a minimax problem. Due to the entropy regularization, GCD may have a connection to entropy-constrained optimal transport problems [36, 37] which regularizes the entropy of the joint distribution. However, the entropy of a single marginal distribution is regularized in GCD.

7 Conclusion

Bridging generative modeling and IRL, GCD opens up new opportunities to improve generative modeling by leveraging tools from reinforcement learning. We are planning to scale up the experiments to large-scale tasks, such as image generation. Besides generation tasks, GCD may also be useful for other tasks that require an accurate energy estimation, such as out-of-distribution detection [21, 29].

Acknowledgments and Disclosure of Funding

S. Yoon was supported by a KIAS Individual Grant (AP095701) via the Center for AI and Natural Sciences at Korea Institute for Advanced Study. Y.-K. Noh was partly supported by NRF/MSIT (No. 2018R1A5A7059549, 2021M3E5D2A01019545) and IITP/MSIT (IITP-2021-0-02068, 2020-0-01373, RS-2023-00220628). This work was supported in part by IITP-MSIT grant 2021-0-02068 (SNU AI Innovation Hub), IITP-MSIT grant 2022-0-00480 (Training and Inference Methods for Goal-Oriented AI Agents), KIAT grant P0020536 (HRD Program for Industrial Innovation), ATC+ MOTIE Technology Innovation Program grant 20008547, SRRC NRF grant RS-2023-00208052, SNU-AIIS, SNU-IAMD, SNU BK21+ Program in Mechanical Engineering, and SNU Institute for Engineering Research.

References

- [1] Geoffrey E Hinton. Training products of experts by minimizing contrastive divergence. *Neural computation*, 14(8):1771–1800, 2002.
- [2] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 6840–6851. Curran Associates, Inc., 2020.
- [3] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 2256–2265, Lille, France, 07–09 Jul 2015. PMLR.
- [4] Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of score-based diffusion models. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 1415–1428. Curran Associates, Inc., 2021.
- [5] Dohyun Kwon, Ying Fan, and Kangwook Lee. Score-based generative modeling secretly minimizes the wasserstein distance. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 20205–20217. Curran Associates, Inc., 2022.
- [6] Sitan Chen, Sinho Chewi, Jerry Li, Yuanzhi Li, Adil Salim, and Anru Zhang. Sampling is as easy as learning the score: theory for diffusion models with minimal data assumptions. In *The Eleventh International Conference on Learning Representations*, 2023.

- [7] Ying Fan and Kangwook Lee. Optimizing ddpm sampling with shortcut fine-tuning. *arXiv preprint arXiv:2301.13362*, 2023.
- [8] Zhifeng Kong and Wei Ping. On fast sampling of diffusion probabilistic models. *arXiv preprint arXiv:2106.00132*, 2021.
- [9] Yang Song and Diederik P Kingma. How to train your energy-based models. *arXiv preprint arXiv:2101.03288*, 2021.
- [10] Siwei Lyu. Unifying non-maximum likelihood learning objectives with minimum kl contraction. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc., 2011.
- [11] Alfred Müller. Integral probability metrics and their generating classes of functions. *Advances in Applied Probability*, 29(2):429–443, 1997.
- [12] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 214–223. PMLR, 06–11 Aug 2017.
- [13] Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, Anind K Dey, et al. Maximum entropy inverse reinforcement learning. In *AAAI*, volume 8, pages 1433–1438. Chicago, IL, USA, 2008.
- [14] Brian D. Ziebart, J. Andrew Bagnell, and Anind K. Dey. Modeling interaction via the principle of maximum causal entropy. In *Proceedings of the 27th International Conference on International Conference on Machine Learning, ICML’10*, page 1255–1262, Madison, WI, USA, 2010. Omnipress.
- [15] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. *Advances in neural information processing systems*, 29, 2016.
- [16] Ying Fan, Olivia Watkins, Yuqing Du, Hao Liu, Moonkyung Ryu, Craig Boutilier, Pieter Abbeel, Mohammad Ghavamzadeh, Kangwook Lee, and Kimin Lee. Reinforcement learning for fine-tuning text-to-image diffusion models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [17] Tianyi Lin, Chi Jin, and Michael Jordan. On gradient descent ascent for nonconvex-concave minimax problems. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 6083–6093. PMLR, 13–18 Jul 2020.
- [18] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Reinforcement learning*, pages 5–32, 1992.
- [19] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [20] Evan Greensmith, Peter L Bartlett, and Jonathan Baxter. Variance reduction techniques for gradient estimates in reinforcement learning. *Journal of Machine Learning Research*, 5(9), 2004.
- [21] Yilun Du, Shuang Li, B. Joshua Tenenbaum, and Igor Mordatch. Improved contrastive divergence training of energy based models. In *Proceedings of the 38th International Conference on Machine Learning (ICML-21)*, 2021.
- [22] Lyudmyla F Kozachenko and Nikolai N Leonenko. Sample estimate of the entropy of a random vector. *Problemy Peredachi Informatsii*, 23(2):9–16, 1987.
- [23] Ruiqi Gao, Yang Song, Ben Poole, Ying Nian Wu, and Diederik P Kingma. Learning energy-based models by diffusion recovery likelihood. In *International Conference on Learning Representations*, 2021.
- [24] Kevin Clark, Paul Vicol, Kevin Swersky, and David J Fleet. Directly fine-tuning diffusion models on differentiable rewards. *arXiv preprint arXiv:2309.17400*, 2023.

- [25] Shu Zhang, Xinyi Yang, Yihao Feng, Can Qin, Chia-Chih Chen, Ning Yu, Zeyuan Chen, Huan Wang, Silvio Savarese, Stefano Ermon, et al. Hive: Harnessing human feedback for instructional visual editing. *arXiv preprint arXiv:2303.09618*, 2023.
- [26] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 27730–27744. Curran Associates, Inc., 2022.
- [27] Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion models with reinforcement learning. *arXiv preprint arXiv:2305.13301*, 2023.
- [28] Yilun Du, Shuang Li, and Igor Mordatch. Compositional visual generation with energy based models. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 6637–6647. Curran Associates, Inc., 2020.
- [29] Sangwoong Yoon, Yung-Kyun Noh, and Frank Park. Autoencoding under normalization constraints. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 12087–12097. PMLR, 18–24 Jul 2021.
- [30] Sangwoong Yoon, Young-Uk Jin, Yung-Kyun Noh, and Frank C. Park. Energy-based models for anomaly detection: A manifold diffusion recovery approach. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [31] Yilun Du and Igor Mordatch. Implicit generation and modeling with energy based models. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d Alche-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 3608–3618. Curran Associates, Inc., 2019.
- [32] Zhisheng Xiao, Karsten Kreis, Jan Kautz, and Arash Vahdat. {VAEBM}: A symbiosis between variational autoencoders and energy-based models. In *International Conference on Learning Representations*, 2021.
- [33] Hankook Lee, Jongheon Jeong, Sejun Park, and Jinwoo Shin. Guiding energy-based models via contrastive latent variables. In *International Conference on Learning Representations*, 2023.
- [34] 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 *International Conference on Learning Representations*, 2022.
- [35] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [36] Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc., 2013.
- [37] Petr Mokrov, Alexander Korotin, and Evgeny Burnaev. Energy-guided entropic neural optimal transport. *arXiv preprint arXiv:2304.06094*, 2023.
- [38] Ernest K Ryu and Wotao Yin. *Large-scale convex optimization: algorithms & analyses via monotone operators*. Cambridge University Press, 2022.

A Proofs

A.1 Proof of Proposition 1

Solving the inner maximization yields $\pi^*(\mathbf{x}) = q(\mathbf{x})$, which is attainable since $\Pi = \mathcal{Q}$. Consequently, the outer minimization simplifies to $\min_{q \in \mathcal{Q}} KL(p||q)$, which reaches its minimum at $q^*(\mathbf{x}) = p(\mathbf{x})$, attainable as $p \in \mathcal{Q}$.

A.2 Proof of Proposition 2

We will prove Proposition 2 by showing that GCD learning in energy (Eq. 4) and the entropy-regularized IPM minimization (Eq. 6) are the primal and the dual problems of the same objective function, where total duality holds. Total duality implies strong duality, and thus the primal and the dual optimal values are the same.

Primal and Dual Problems. We consider a function $\Phi : \mathcal{E} \times \Pi \rightarrow \mathbb{R}$ defined as follows:

$$\Phi(E, \pi) = \mathbb{E}_p[E(\mathbf{x})] - \mathbb{E}_\pi[E(\mathbf{x})] + \tau \mathcal{H}(\pi). \quad (10)$$

The **primal** problem generated by $\Phi(E, \pi)$ is given as

$$\inf_{E \in \mathcal{E}} \sup_{\pi \in \Pi} \mathbb{E}_p[E(\mathbf{x})] - \mathbb{E}_\pi[E(\mathbf{x})] + \tau \mathcal{H}(\pi), \quad (11)$$

which is equivalent to GCD learning in energy (Eq. 4). The **dual** problem generated by $\Phi(E, \pi)$ is equivalent to the minimization of IPM under entropy regularization:

$$\inf_{\pi \in \Pi} \sup_{E \in \mathcal{E}} -\mathbb{E}_p[E(\mathbf{x})] + \mathbb{E}_\pi[E(\mathbf{x})] - \tau \mathcal{H}(\pi) = \inf_{\pi \in \Pi} \sup_{E \in \mathcal{E}} |\mathbb{E}_p[E(\mathbf{x})] - \mathbb{E}_\pi[E(\mathbf{x})]| - \tau \mathcal{H}(\pi) \quad (12)$$

$$= \inf_{\pi \in \Pi} D^\mathcal{E}(p||\pi) - \tau \mathcal{H}(\pi), \quad (13)$$

where we used the condition that \mathcal{E} is closed under negation to obtain Eq. 12.

Total Duality. We will show that total duality holds between these primal and dual problems. The function Φ is linear (and hence convex) for $E(\mathbf{x})$ and concave for $\pi(\mathbf{x})$. Then, total duality holds if and only if Φ has a saddle point ([38, Ch. 1, pp. 12-14]).

Let us define a function $E^p(\mathbf{x}) = -\tau \log p(\mathbf{x}) + c$, where c being a constant, which can be considered as the energy of $p(\mathbf{x})$ computed under the temperature τ . Let us show that $(E^*, \pi^*) = (E^p, p)$ is the saddle point of Φ for any constant c in E^p . Note that $(E^p, p) \in \mathcal{E} \times \Pi$ as $p \in \mathcal{Q} \cap \Pi$.

$$\Phi(E^p, \pi) = \mathbb{E}_p[E^p(\mathbf{x})] + \mathbb{E}_\pi[\tau \log p - \tau \log \pi] \quad (14)$$

$$= \mathbb{E}_p[E^p(\mathbf{x})] - \tau KL(\pi||p) \leq \Phi(E^p, p) \quad \forall \pi \in \Pi \quad (15)$$

Also,

$$\Phi(E, p) = \mathbb{E}_p[E(\mathbf{x})] - \mathbb{E}_p[E(\mathbf{x})] + \tau \mathcal{H}(p) \quad (16)$$

$$= \tau \mathcal{H}(p) = \Phi(E^p, p) \quad \forall E \in \mathcal{E} \quad (17)$$

Combining Eq. 15 and Eq. 17, we obtain $\Phi(E^p, \pi) \leq \Phi(E^p, p) \leq \Phi(E, p)$ for all $(E, \pi) \in \mathcal{E} \times \Pi$, and thus (E^p, p) is a saddle point of Φ . The saddle point is the equilibrium for both the primal and the dual problems.

B Algorithm

Algorithm 1 Generalized Contrastive Divergence

```

Initialize sample buffer  $\mathcal{B} = \emptyset$ .
Set  $i = 0$ . Set sampler update period  $n_{sampler}$  and the number of PPO update  $n_{PPO}$ .
for  $\mathbf{x}$  in dataset do
    Obtain  $m$  samples from  $\pi(\mathbf{x})$  and append them into  $\mathcal{B}$ .
    Update EBM  $E(\mathbf{x})$  using  $\mathcal{B}$  by Eq. 8.
    Update value functions  $V(\mathbf{x}_t, t)$ 's using  $\mathcal{B}$ .
    if  $i \% n_{sampler} = 0$  then
        for  $j = 1, \dots, n_{PPO}$  do
            Update  $\pi(\mathbf{x})$  using PPO using  $\mathcal{B}_{on}$ .
        end for
        Reset  $\mathcal{B}$  as  $\emptyset$ .
    end if
     $i \leftarrow i + 1$ 
end for

```

C Derivation of Policy Gradient

Here, we provide a detailed derivation on the policy gradient update for sample parameters (Eq. 9).

$$\nabla_{\phi} \mathcal{L} = \nabla_{\phi} \int \pi_{\phi}(\mathbf{x}_T) [-E_{\theta}(\mathbf{x}_T) - \tau \log \pi_{\phi}(\mathbf{x}_T)] d\mathbf{x}_T \quad (18)$$

$$= \int \nabla_{\phi} \pi_{\phi}(\mathbf{x}_T) [-E_{\theta}(\mathbf{x}_T) - \tau \log \pi_{\phi}(\mathbf{x}_T)] d\mathbf{x}_T \quad (19)$$

$$= \int \nabla_{\phi} \left(\int \pi_{\phi}(\mathbf{x}_{0:T}) d\mathbf{x}_{0:T-1} \right) [-E_{\theta}(\mathbf{x}_T) - \tau \log \pi_{\phi}(\mathbf{x}_T)] d\mathbf{x}_T \quad (20)$$

$$= \int \nabla_{\phi} \pi_{\phi}(\mathbf{x}_{0:T}) [-E_{\theta}(\mathbf{x}_T) - \tau \log \pi_{\phi}(\mathbf{x}_T)] d\mathbf{x}_{0:T} \quad (21)$$

$$= \int \pi_{\phi}(\mathbf{x}_{0:T}) \nabla_{\phi} \log \pi_{\phi}(\mathbf{x}_{0:T}) [-E_{\theta}(\mathbf{x}_T) - \tau \log \pi_{\phi}(\mathbf{x}_T)] d\mathbf{x}_{0:T} \quad (22)$$

$$= \mathbb{E}_{\pi_{\phi}(\mathbf{x}_{0:T})} [\nabla_{\phi} \log \pi_{\phi}(\mathbf{x}_{0:T}) [-E_{\theta}(\mathbf{x}_T) - \tau \log \pi_{\phi}(\mathbf{x}_T)]] \quad (23)$$

$$= \mathbb{E}_{\pi_{\phi}(\mathbf{x}_{0:T})} \left[\left(\sum_{t=0}^{t=T-1} \nabla_{\phi} \log \pi_{\phi}(\mathbf{x}_{t+1} | \mathbf{x}_t, t) \right) [-E_{\theta}(\mathbf{x}_T) - \tau \log \pi_{\phi}(\mathbf{x}_T)] \right]. \quad (24)$$