

---

# Discriminator Contrastive Divergence: Semi-Amortized Generative Modeling by Exploring Energy of the Discriminator

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Generative Adversarial Networks (GANs) have shown great promise in model-  
2 ing high dimensional data. The learning objective of GANs usually minimizes  
3 some measure discrepancy, *e.g.*,  $f$ -divergence ( $f$ -GANs) or Integral Probability  
4 Metric (Wasserstein GANs). With  $f$ -divergence as the objective function, the  
5 discriminator essentially estimates the density ratio, and the estimated ratio proves  
6 useful in further improving the sample quality of the generator. However, how  
7 to leverage the information contained in the discriminator of Wasserstein GANs  
8 (WGAN) is less explored. In this paper, we introduce the Discriminator Contrastive  
9 Divergence, which is well motivated by the property of WGAN’s discriminator and  
10 the relationship between WGAN and energy-based model. Compared to standard  
11 GANs, where the generator is directly utilized to obtain new samples, our method  
12 proposes a semi-amortized generation procedure where the samples are produced  
13 with the generator’s output as an initial state. Then several steps of Langevin  
14 dynamics are conducted using the gradient of the discriminator. We demonstrate  
15 the benefits of significantly improved generation on both synthetic data and several  
16 real-world image generation benchmarks.

## 17 1 Introduction

18 Generative Adversarial Networks (GANs) [10] proposes a popular way to learn likelihood-free gener-  
19 ative models, which have shown promising results on various challenging tasks. Specifically, GANs  
20 are learned by finding the equilibrium of a min-max game between a generator and a discriminator or  
21 a critic. Assuming the optimal discriminator can be obtained, the generator substantially minimizes  
22 some discrepancy between the generated distribution and the target distribution.

23 Improving training GANs by exploring the discrepancy measure with the excellent property has stimu-  
24 lated fruitful lines of research works and is still an active area. Two well-known discrepancy measures  
25 for training GANs are  $f$ -divergence and Integral Probability Metric (IPM) [26].  $f$ -divergence is  
26 severe for directly minimization due to the intractable integral,  $f$ -GANs provide minimization instead  
27 of a variational approximation of  $f$ -divergence between the generated distribution  $p_{G_\theta}$  and the target  
28 distribution  $p_{\text{data}}$ . The discriminator in  $f$ -GANs serves as a density ratio estimator [36]. The other  
29 families of GANs are based on the minimization of an Integral Probability Metric (IPM). According  
30 to the definition of IPM, the critic needs to be constrained into a specific function class. When  
31 the critic is restricted to be 1-Lipschitz function, the corresponding IPM turns to the Wasserstein-1  
32 distance, which inspires the approaches of Wasserstein GANs (WGANs) [25, 1, 13].

33 No matter what kind of discrepancy is evaluated and minimized, the discriminator is usually discarded  
34 at the end of the training, and only the generator is kept to generate samples. A natural question to

ask is whether, and how we can leverage the remaining information in the discriminator to construct a more superior distribution than simply sampling from a generator.

Recent work [2, 35] has shown that a density ratio can be obtained through the output of discriminator, and a more superior distribution can be acquired by conducting rejection sampling or Metropolis-Hastings sampling with the estimated density ratio based on the original GAN [10].

However, the critical limitation of previous methods lies in that they can not be adapted to WGANs, which enjoy superior empirical performance over other variants. How to leverage the information of a WGAN’s critic model to improve image generation remains an open problem. In this paper, we do the following to address this:

- We provide a generalized view to unify different families of GANs by investigating the informativeness of the discriminators.
- We propose a semi-amortized generative modeling procedure so-called discriminator contrastive divergence (DCD), which achieves an intermediate between implicit and explicit generation and hence allows a trade-off between generation quality and speed.

Extensive experiments are conducted to demonstrate the efficacy of our proposed method on both synthetic setting and real-world generation scenarios, which achieves state-of-the-art performance on several standard evaluation benchmarks of image generation.

## 2 Methodology

We first introduce the Fenchel dual of the intractable partition function  $Z_\theta$  in Eq. 8:

**Theorem 1.** [38] *With  $H(q) = -\int q(x) \log q(x) dx$ , the Fenchel dual of log-partition  $Z_\theta$  is as follows:*

$$A(E_\theta) = \max_{q \in \mathcal{P}} \langle q(x), E_\theta(x) \rangle + H(q), \quad (1)$$

where  $\mathcal{P}$  denotes the space of distributions, and  $\langle q(x), E_\theta(x) \rangle = \int E_\theta(x) q(x) dx$ .

We put the Fenchel dual of  $A(E_\theta)$  back into the MLE objective in Eq. 9, we achieve the following min-max game formalization for training energy-based model based on MLE:

$$\min_{q \in \mathcal{P}} \max_{E_\theta \in \mathcal{E}} \underbrace{\mathbb{E}_{\mathbf{x} \sim P_{\text{data}}} [E_\theta(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim q} [E_\theta(\mathbf{x})]}_{\text{WGAN's objective for critic}} - \underbrace{H(q)}_{\text{entropy regularization}}. \quad (2)$$

The Fenchel dual view of MLE training in the energy-based model explicitly illustrates the gap and connection between the WGAN and Energy based model. If we consider the dual distribution  $q$  as the generated distribution  $p_{G_\theta}$ , and the  $D_\phi$  as the energy function  $E_\theta$ . The duality form for training energy-based models is essentially the WGAN’s objective with the entropy of the generator is regularized.

Hence to turn the discriminator in WGAN into an energy function, we may conduct several fine-tuning steps, as illustrated in Eq. 2. Note that maximizing the entropy of the  $p_{G_\theta}$  is indeed a challenging task, which needs to either use a tractable density generator, e.g., normalizing Flows [7], or maximize the mutual information between the latent variable  $\mathcal{Z}$  and the corresponding  $G_\theta(\mathcal{Z})$  when the  $G_\theta$  is a deterministic mapping. However, instead of maximizing the entropy of the generated distribution  $p_{G_\theta}$  directly, we derive our method based on the following fact:

**Proposition 1.** [19] *Update the generated distribution  $p_{G_\theta}$  according to the gradient estimated through Equation. 2, essentially minimized the Kullback–Leibler (KL) divergence between  $p_{G_\theta}$  and the distribution  $p_{D_\phi}$ , which refers to the distribution implied by using  $D_\phi$  as the energy function, as illustrated in Eq. 8, i.e.  $D_{\text{KL}}(p_{G_\theta} || p_{D_\phi})$ .*

To avoid the computation of  $H(p_{G_\theta})$ , motivated by the monotonic property of MCMC, as illustrated in Eq. 11, we propose Discriminator Contrastive Divergence (DCD), which replaces the gradient-based optimization on  $q(p_{G_\theta})$  in Eq. 2 with several steps of MCMC for finetuning the critic in WGAN into an energy function. To be more specific, we use Langevin dynamics[33] which leverages the gradient of the discriminator to conduct sampling:

$$x_k = x_{k-1} - \frac{\epsilon}{2} \nabla_x D_\phi(x_{k-1}) + \sqrt{\epsilon} \omega, \omega \sim \mathcal{N}(0, I), \quad (3)$$

79 where  $\epsilon$  refers to the step size. The GAN-based approaches are implicitly constrained by the dimension  
 80 of the latent noise, which is based on a widely applied assumption that the high dimensional data,  
 81 *e.g.*, images, actually distribute on a relatively low-dimensional manifold. Apart from searching the  
 82 reasonable point in the data space, we could also find the lower energy part of the latent manifold by  
 83 conducting Langevin dynamics in the latent space which are more stable in practice, *i.e.*:

$$z_t^l = z_t^{l-1} - \frac{\epsilon}{2} \nabla_z D_\phi (G_\theta(z_t)^{l-1}) + \sqrt{\epsilon} \omega, \omega \sim \mathcal{N}(0, \mathcal{I}). \quad (4)$$

84 Ideally, the proposal should be accepted or rejected according to the Metropolis–Hastings algorithm:

$$\alpha := \min \left\{ 1, \frac{D_\phi(x_k) q(x_{k-1}|x_k)}{D_\phi(x_{k-1}) q(x_k|x_{k-1})} \right\}, \quad (5)$$

85 where  $q$  refers to the proposal which is defined as:

$$q(x'|x) \propto \exp \left( -\frac{1}{4\tau} \|x' - x - \tau \nabla \log \pi(x)\|_2^2 \right). \quad (6)$$

86 In practice, we find the rejection steps described in Eq. 5 do not boost performance. For simplicity,  
 87 following [31, 8], we apply Eq. 3 in experiments as an approximate version. The whole tuning  
 88 procedure is illustrated in Algorithm 1.

89 After fine-tuning, the discriminator function can be approximated seen as an unnormalized probability  
 90 function, which implies a unique distribution  $p_{D_\phi}$ . And similar to the  $p_*$  implied in the rejection  
 91 sampling-based method, it is reasonable to assume that  $p_{D_\phi}$  is a superior distribution of  $p_{G_\theta}$ . Sampling  
 92 from  $p_{D_\phi}$  can be implemented through the Langevin dynamics, as illustrated in Eq. 3 with  $p_{G_\theta}$  serves  
 93 as the initial distribution.

## 94 3 Experiments

### 95 3.1 Synthetic Density Modeling

96 Displaying the level sets is a meaningful way to  
 97 study learned critic. Following the [2, 13], we  
 98 investigate the impacts of our method on two  
 99 challenging low-dimensional synthetic settings:  
 100 twenty-five isotropic Gaussian distributions ar-  
 101 ranged in a grid and eight Gaussian distributions  
 102 arranged in a ring (Fig. 1a). For all different set-  
 103 tings, both the generator and the discriminator  
 104 of the WGAN model are implemented as neural  
 105 networks with four fully connected layers and  
 106 Relu activations. The Lipschitz constraint is  
 107 restricted through spectral normalization [25],  
 108 while the prior is a two-dimensional multivari-  
 109 ate Gaussian with a mean of 0 and a standard  
 110 deviation of 1.

111 To investigate whether the proposed Discrimina-  
 112 tor Contrastive Divergence is capable of tuning  
 113 the distribution induced by the discriminator as  
 114 desired energy function, *i.e.*  $p_{D_\phi}$ , we visual-  
 115 ize both the value surface of the critic and the  
 116 samples obtained from  $p_{D_\phi}$  with Langevin dy-  
 117 namics. The results are shown in Figure. 1. As can be observed, the original WGAN (Fig. 1b) is strong  
 118 enough to cover most modes, but there are still some spurious links between two different modes. The  
 119 enhanced distribution  $p_{D_\phi}$  (Fig. 1c), however, has the ability to reduce spurious links and recovers the  
 120 modes with underestimated density. More precisely, after the MCMC fine-tuning procedure (Fig. 1c),  
 121 the gradients of the value surface become more meaningful so that all the regions with high density in  
 122 data distribution  $p_{\text{data}}$  are assigned with high  $D_\phi$  value, *i.e.*, lower energy ( $\exp(-D_\phi)$ ). By contrast,  
 123 in the original discriminator (Fig. 1b), the lower energy regions in  $p_{D_\phi}$  are not necessarily consistent  
 124 with the high-density region of  $p_{\text{data}}$ .

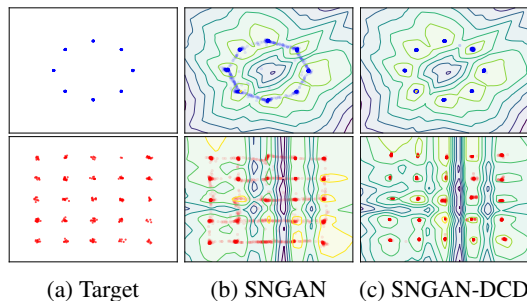


Table 1: Density modeling on synthetic distributions. **Top:** 8 Gaussian distribution. **Bottom:** 25 Gaussian distribution. **Left:** Distribution of real data. **Middle:** Distribution defined by the generator of SNGAN. The surface is the level set of the critic. Yellow corresponds to higher value while purple corresponds to lower. **Right:** Distribution defined by the SNGAN-DCD. The surface is the level set of the proposed energy function.

Model	Inception	FID
<b>CIFAR-10 Unconditional</b>		
PixelCNN [37]	4.60	65.93
PixelIQN [28]	5.29	49.46
EBM [8]	6.02	40.58
WGAN-GP [13]	$7.86 \pm .07$	18.12
MoLM [29]	$7.90 \pm .10$	18.9
SNGAN [25]	$8.22 \pm .05$	21.7
ProgressiveGAN [18]	$8.80 \pm .05$	-
NCSN [31]	$8.87 \pm .12$	25.32
DCGAN w/ DRS [2]	3.073	-
DCGAN w/ MH-GAN [35]	3.379	-
ResNet-SAGAN w/ DOT [32]	$8.50 \pm .12$	19.71
<b>SNGAN-DCD (Pixel)</b>	$8.54 \pm .11$	21.67
<b>SNGAN-DCD (Latent)</b>	<b><math>9.11 \pm .04</math></b>	<b>16.24</b>
<b>CIFAR-10 Conditional</b>		
EBM [8]	8.30	37.9
SNGAN [25]	$8.43 \pm .09$	15.43
<b>SNGAN-DCD (Pixel)</b>	$8.73 \pm .13$	22.84
<b>SNGAN-DCD (Latent)</b>	$8.81 \pm .11$	15.05
BigGAN [3]	<b>9.22</b>	<b>14.73</b>

Table 2: Inception and FID scores for CIFAR-10.

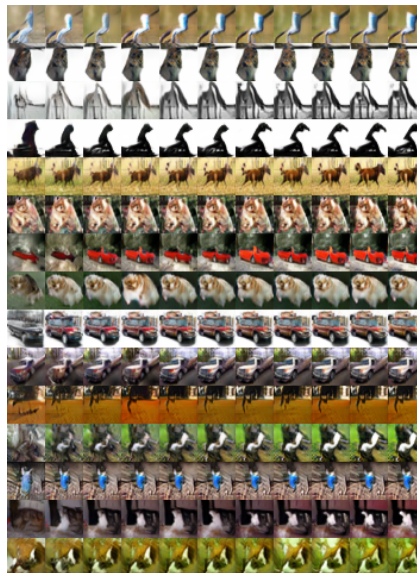


Figure 1: Unconditional CIFAR-10 Langevin dynamics visualization.

### 125 3.2 Real-World Image Generation

126 For quantitative evaluation, we report the inception score [30] and FID [15] scores on CIFAR-10  
127 in Tab. 2. As shown in the Tab. 2, in pixel space, by introducing the proposed DCD algorithm,  
128 we achieve a significant improvement of inception score over the SNGAN. The reported inception  
129 score is even higher than most values achieved by class-conditional generative models. Our FID  
130 score of 21.67 on CIFAR-10 is competitive with other top generative models. When the DCD is  
131 conducted in the latent space, we further achieve a 9.11 inception score and a 16.24 FID, which is a  
132 new state-of-the-art performance of IS. When combined with label information to perform conditional  
133 generation, we further improve the FID to 15.05, which is comparable with current state-of-the-art  
134 large-scale trained models [3]. Some visualization of generated examples can be found in Fig 1,  
135 which demonstrates that the Markov chain is able to generate more realistic samples, suggesting that  
136 the MCMC process is meaningful and effective. Tab. 4 and Tab. 5 shows the performance on STL-10  
137 and ImageNet respectively, which demonstrate that as a generalized method, DCD is not over-fitted to  
138 the specific dataset. More experiment details and the generated samples can be found in Appendix. I.

### 139 4 Conclusion and Future Work

140 Based on the density ratio estimation perspective, the discriminator in  $f$ -GANs could be adapted to a  
141 wide range of applications, *e.g.*, mutual information estimation [17] and bias correction of generative  
142 models [12]. However, as another important branch in GANs, the available information in WGANs’  
143 discriminator is less explored. In this paper, we narrow down the scope and focus on how to leverage  
144 the discriminator of WGANs to further improve the sample quality. We first present a comprehensive  
145 theoretical analysis on the informativeness of WGANs’ discriminator. Motivated by the theoretical  
146 understanding, we investigate the possibility of turning the discriminator of WGANs into an energy  
147 function and propose a tuning and sampling procedure named “Discriminator Contrastive Divergence”.  
148 The final generation process is semi-amortized, where we take the generator as the initial state and  
149 then conduct several MCMC steps. Empirical results demonstrate the effectiveness of the proposed  
150 method on several tasks. We hope our work can shed some light on a generalized view to a method  
151 of connecting different GANs and energy-based models, which will stimulate more exploration into  
152 the potential of current deep generative models. One potential direction for future work is to conduct  
153 DCD in each layer of the generator. This can be seen as a compromise between the latent and the  
154 pixel space, which may lead to further sampling quality improvements.

## References

- 155 [1] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein gan. *arXiv preprint*  
156 *arXiv:1701.07875*, 2017.
- 158 [2] Samaneh Azadi, Catherine Olsson, Trevor Darrell, Ian Goodfellow, and Augustus Odena.  
159 Discriminator rejection sampling. *arXiv preprint arXiv:1810.06758*, 2018.
- 160 [3] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity  
161 natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018.
- 162 [4] Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsuper-  
163 vised feature learning. In *Proceedings of the fourteenth international conference on artificial*  
164 *intelligence and statistics*, pages 215–223, 2011.
- 165 [5] Thomas M Cover and Joy A Thomas. *Elements of information theory*. John Wiley & Sons,  
166 2012.
- 167 [6] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-  
168 scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern*  
169 *recognition*, pages 248–255. Ieee, 2009.
- 170 [7] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp.  
171 *arXiv preprint arXiv:1605.08803*, 2016.
- 172 [8] Yilun Du and Igor Mordatch. Implicit generation and generalization in energy-based models.  
173 *arXiv preprint arXiv:1903.08689*, 2019.
- 174 [9] Chelsea Finn, Paul Christiano, Pieter Abbeel, and Sergey Levine. A connection between  
175 generative adversarial networks, inverse reinforcement learning, and energy-based models.  
176 *arXiv preprint arXiv:1611.03852*, 2016.
- 177 [10] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil  
178 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural*  
179 *information processing systems*, pages 2672–2680, 2014.
- 180 [11] Will Grathwohl, Kuan-Chieh Wang, Jörn-Henrik Jacobsen, David Duvenaud, Mohammad  
181 Norouzi, and Kevin Swersky. Your classifier is secretly an energy based model and you should  
182 treat it like one. *arXiv preprint arXiv:1912.03263*, 2019.
- 183 [12] Aditya Grover, Jiaming Song, Alekh Agarwal, Kenneth Tran, Ashish Kapoor, Eric Horvitz, and  
184 Stefano Ermon. Bias correction of learned generative models using likelihood-free importance  
185 weighting. *arXiv preprint arXiv:1906.09531*, 2019.
- 186 [13] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville.  
187 Improved training of wasserstein gans. In *Advances in neural information processing systems*,  
188 pages 5767–5777, 2017.
- 189 [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
190 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
191 pages 770–778, 2016.
- 192 [15] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.  
193 Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances*  
194 *in Neural Information Processing Systems*, pages 6626–6637, 2017.
- 195 [16] Jean-Baptiste Hiriart-Urruty and Claude Lemaréchal. *Fundamentals of convex analysis*. Springer  
196 Science & Business Media, 2012.
- 197 [17] R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman,  
198 Adam Trischler, and Yoshua Bengio. Learning deep representations by mutual information  
199 estimation and maximization. *arXiv preprint arXiv:1808.06670*, 2018.
- 200 [18] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for  
201 improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017.

- 202 [19] Taesup Kim and Yoshua Bengio. Deep directed generative models with energy-based probability  
203 estimation. *arXiv preprint arXiv:1606.03439*, 2016.
- 204 [20] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*  
205 *arXiv:1412.6980*, 2014.
- 206 [21] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.  
207 Technical report, Citeseer, 2009.
- 208 [22] Rithesh Kumar, Anirudh Goyal, Aaron Courville, and Yoshua Bengio. Maximum entropy  
209 generators for energy-based models. *arXiv preprint arXiv:1901.08508*, 2019.
- 210 [23] Yingzhen Li, Richard E Turner, and Qiang Liu. Approximate inference with amortised mcmc.  
211 *arXiv preprint arXiv:1702.08343*, 2017.
- 212 [24] Takeru Miyato and Masanori Koyama. cgans with projection discriminator. *arXiv preprint*  
213 *arXiv:1802.05637*, 2018.
- 214 [25] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization  
215 for generative adversarial networks. *arXiv preprint arXiv:1802.05957*, 2018.
- 216 [26] Alfred Müller. Integral probability metrics and their generating classes of functions. *Advances*  
217 *in Applied Probability*, 29(2):429–443, 1997.
- 218 [27] Sebastian Nowozin, Botond Cseke, and Ryota Tomioka. f-gan: Training generative neural sam-  
219 plers using variational divergence minimization. In *Advances in neural information processing*  
220 *systems*, pages 271–279, 2016.
- 221 [28] Georg Ostrovski, Will Dabney, and Rémi Munos. Autoregressive quantile networks for genera-  
222 tive modeling. *arXiv preprint arXiv:1806.05575*, 2018.
- 223 [29] Suman Ravuri, Shakir Mohamed, Mihaela Rosca, and Oriol Vinyals. Learning implicit genera-  
224 tive models with the method of learned moments. *arXiv preprint arXiv:1806.11006*, 2018.
- 225 [30] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen.  
226 Improved techniques for training gans. In *Advances in neural information processing systems*,  
227 pages 2234–2242, 2016.
- 228 [31] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data  
229 distribution. In *Advances in Neural Information Processing Systems*, pages 11895–11907, 2019.
- 230 [32] Akinori Tanaka. Discriminator optimal transport. *arXiv preprint arXiv:1910.06832*, 2019.
- 231 [33] Yee Whye Teh, Max Welling, Simon Osindero, and Geoffrey E Hinton. Energy-based models  
232 for sparse overcomplete representations. *Journal of Machine Learning Research*, 4(Dec):  
233 1235–1260, 2003.
- 234 [34] Lucas Theis, Aäron van den Oord, and Matthias Bethge. A note on the evaluation of generative  
235 models. *arXiv preprint arXiv:1511.01844*, 2015.
- 236 [35] Ryan Turner, Jane Hung, Yunus Saatci, and Jason Yosinski. Metropolis-hastings generative  
237 adversarial networks. *arXiv preprint arXiv:1811.11357*, 2018.
- 238 [36] Masatoshi Uehara, Issei Sato, Masahiro Suzuki, Kotaro Nakayama, and Yutaka Matsuo.  
239 Generative adversarial nets from a density ratio estimation perspective. *arXiv preprint*  
240 *arXiv:1610.02920*, 2016.
- 241 [37] Aaron Van den Oord, Nal Kalchbrenner, Lasse Espeholt, Oriol Vinyals, Alex Graves, et al.  
242 Conditional image generation with pixelcnn decoders. In *Advances in neural information*  
243 *processing systems*, pages 4790–4798, 2016.
- 244 [38] Martin J Wainwright, Michael I Jordan, et al. Graphical models, exponential families, and  
245 variational inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008.

- 246 [39] Junbo Zhao, Michael Mathieu, and Yann LeCun. Energy-based generative adversarial network.  
247 *arXiv preprint arXiv:1609.03126*, 2016.
- 248 [40] Zhiming Zhou, Jiadong Liang, Yuxuan Song, Lantao Yu, Hongwei Wang, Weinan Zhang, Yong  
249 Yu, and Zhihua Zhang. Lipschitz generative adversarial nets. *arXiv preprint arXiv:1902.05687*,  
250 2019.

251 **A Related Works**

252 Both empirical [1] and theoretical [15] evidence has demonstrated that learning a discriminative model  
 253 with neural networks is relatively easy, and the neural generative model (sampler) is prone to reach its  
 254 bottleneck during the optimization. Hence, there is strong motivation to further improve the generated  
 255 distribution by exploring the remaining information. Two recent advancements are discriminator  
 256 rejection sampling (DRS) [2] and MH-GANs [35]. DRS conducts rejection sampling on the output  
 257 of the generator. The vital limitation that lies in the upper bound of  $D_\phi$  is needed to be estimated  
 258 for computing the rejection probability. MH-GAN sidesteps the above problem by introducing a  
 259 Metropolis-Hastings sampling procedure with generator acting as the independent proposal; the state  
 260 transition is estimated with a well-calibrated discriminator. However, the theoretical justification  
 261 of both the above two methods is based on the fact that the output of discriminator needs to be  
 262 viewed as an estimation of density ratio  $\frac{p_{\text{data}}}{p_{G_\theta}}$ . As pointed out by previous work [40], the output of a  
 263 discriminator in WGAN [1] suffers from the free offset and can not provide the density ratio, which  
 264 prevents the application of the above methods in WGAN.

265 Our work is inspired by recent theoretical studies on the property of discriminator in WGANs [13, 40].  
 266 [32] proposes discriminator optimal transport (DOT) to leverage the optimal transport plan implied  
 267 by WGANs’ discriminator, which is orthogonal to our method. Besides, turning the discriminator of  
 268 WGAN into an energy function is closely related to the amortized generation methods in energy-based  
 269 model (EBM) literature [19, 39, 22] where a separate network is proposed to learn to sample from  
 270 the partition function in [9]. Recent progress [31, 8] in the area of EBM has shown the feasibility  
 271 of generating high dimensional data with Langevin dynamics. From the perspective of EBM, our  
 272 proposed method can be seen as an intermediary between an amortized generative model and an  
 273 implicit generative model, *i.e.*, a semi-amortized generation method, which allows a trade-off between  
 274 speed and quality of generation. With a similar spirit, [11] also illustrates the potential connection  
 275 between neural classifier and energy-based model in supervised and semi-supervised scenarios.

276 **B Preliminaries**

277 **B.1 Generative Adversarial Networks**

278 Generative Adversarial Networks (GANs) [10] is an implicit generative model that aims to fit an  
 279 empirical data distribution  $p_{\text{data}}$  over sample space  $\mathcal{X}$ . The generative distribution  $p_{G_\theta}$  is implied by a  
 280 generated function  $G_\theta$ , which maps latent variable  $Z$  to sample  $X$ , *i.e.*,  $G_\theta : \mathcal{Z} \rightarrow \mathcal{X}$ . Typically, the  
 281 latent variable  $Z$  is distributed on a fixed prior distribution  $p(z)$ . With i.i.d samples available from  
 282  $p_{G_\theta}$  and  $p_{\text{data}}$ , the GAN typically learns the generative model through a min-max game between a  
 283 discriminator  $D_\phi$  and a generator  $G_\theta$ :

$$\min_{\theta} \max_{\phi} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [r(D_\phi(\mathbf{x}))] - \mathbb{E}_{\mathbf{x} \sim p_{G_\theta}} [m(D_\phi(\mathbf{x}))]. \quad (7)$$

284 With  $r$  and  $m$  as the function  $r(x) = m(x) = x$  and the  $D_\phi(x)$  is constrained as 1-Lipschitz  
 285 function, the Eq. 7 yields the WGANs objective which essentially minimizes the Wasserstein distance  
 286 between  $p_{\text{data}}$  and  $p_{G_\theta}$ . With  $r(x) = x$  and  $m(x)$  as the Fenchel conjugate[16] of a convex and lower-  
 287 semicontinuous function, the objective in Eq. 7 approximately minimize a variational estimation of  
 288  $f$ -divergence[27] between  $p_{\text{data}}$  and  $p_{G_\theta}$ .

289 **B.2 Energy Based Model and MCMC basics**

290 The energy-based model tends to learn an unnormalized probability model implied by an energy  
 291 function  $E_\theta(x)$  to prescribe the ground truth data distribution  $p_{\text{data}}$ . The corresponding normalized  
 292 density function is:

$$q_\theta(x) = \frac{e^{-E_\theta(x)}}{Z_\theta}, \quad Z_\theta = \int e^{-E_\theta(x)} dx, \quad (8)$$

293 where  $Z_\theta$  is so-called normalization constant. The objective of training an energy-based model with  
 294 maximum likelihood estimation is as:

$$\mathcal{L}_{\text{MLE}}(\theta; p) := -\mathbb{E}_{x \sim p_{\text{data}}(x)} [\log q_\theta(x)]. \quad (9)$$



295 The estimated gradient with respect to the maximum likelihood estimation objective is as follows:

$$\begin{aligned} \nabla_{\theta} \mathcal{L}_{\text{MLE}}(\theta; p) &= \nabla_{\theta} \mathbb{E}_{x \sim p_{\text{data}}(x)} [E_{\theta}(x)] - \frac{\int e^{-E_{\theta}(x)} \nabla_{\theta} E_{\theta}(x) dx}{Z_{\theta}} \\ &= \mathbb{E}_{x \sim p_{\text{data}}(x)} [\nabla_{\theta} E_{\theta}(x)] - \mathbb{E}_{x \sim q_{\theta}(x)} [\nabla_{\theta} E_{\theta}(x)]. \end{aligned} \quad (10)$$

296 The above method for gradient estimation in Equation 10 is called contrastive divergence (CD).

297 Markov chain Monte Carlo is a powerful framework for drawing samples from a given distribution.  
 298 An MCMC is specified by a transition kernel  $\mathcal{K}(x'|x)$  which corresponds to a unique stationary  
 299 distribution  $p$ . More specifically, MCMC can be viewed as drawing  $x_0$  from the initial distribution  $q_0$   
 300 and iteratively get sample  $x_t$  at the  $t$ -th iteration by applied the transition kernel on the previous step,  
 301 *i.e.*,  $x_t | x_{t-1} \sim \mathcal{K}(x_t | x_{t-1})$ . Following [23], we formalized the distribution  $q_t$  of  $x_t$  as obtained by a  
 302 fixed point update of form  $q_t(x) \leftarrow \mathcal{K} q_{t-1}(x)$ , and  $\mathcal{K} q_{t-1}(x)$ :

$$\mathcal{K} q_{t-1}(x) := \int q_{t-1}(x') \mathcal{K}(x|x') dx'.$$

303 As indicated by the standard theory of MCMC, the following monotonic property is satisfied:

$$D_{\text{KL}}(q_t || p) \leq D_{\text{KL}}(q_{t-1} || p). \quad (11)$$

304 And  $q_t$  converges to the stationary distribution  $p$  as  $t \rightarrow \infty$ .

### 305 B.3 Informativeness of WGAN Discriminator

306 So far, it is well known that the discriminator  $D_{\phi}$  in  $f$ -GAN is optimized to estimate a statistic related  
 307 to the density ratio between  $\frac{p_{\text{data}}}{p_{G_{\theta}}}$  [2]. In this section, we seek to investigate the following questions:

- 308 • What kind of information is contained in the discriminator of WGANs?
- 309 • Why and how can the information be utilized to further improved the quality of generated  
 310 distribution?

311 Different from  $f$ -GANs, the objective of WGANs is derived from the Integral Probability Metric,  
 312 and the discriminator can not naturally be derived as an estimated density ratio. Before leveraging  
 313 the remaining information in the discriminator, the property of the discriminator in WGANs needs to  
 314 be investigated first. We introduce the primal problem implied by WGANs objective as follows:

315 Let  $\pi$  denote the joint probability for transportation between  $P$  and  $Q$ , which satisfies the marginality  
 316 conditions,

$$\int d\mathbf{y} \pi(\mathbf{x}, \mathbf{y}) = p(\mathbf{x}), \quad \int d\mathbf{x} \pi(\mathbf{x}, \mathbf{y}) = q(\mathbf{y}) \quad (12)$$

317 The primal form first-order Wasserstein distance  $W_1$  is defined as:

$$W_1(P, Q) = \inf_{\pi \in \Pi(P, Q)} \mathbb{E}_{(x,y) \sim \pi} [\|x - y\|_2]$$

318 the objective function of the discriminator in Wasserstein GANs is the Kantorovich-Rubinstein duality  
 319 of Eq. 12, and the optimal discriminator has the following property[13]:

320 **Theorem 2.** *Let  $\pi^*$  as the optimal transport plan in Eq. 12 and  $x_t = tx + (1-t)y$  with  $0 \leq t \leq 1$ .  
 321 With the optimal discriminator  $D_{\phi}$  as a differentiable function and  $\pi^*(x, x) = 0$  for all  $x$ , then it  
 322 holds that:*

$$\mathbb{P}_{(x,y) \sim \pi^*} \left[ \nabla_{x_i} D_{\phi}^*(x_t) = \frac{y - x}{\|y - x\|} \right] = 1$$

323 Theorem. 2 states that for each sample  $x$  in the generated distribution  $p_{G_{\theta}}$ , the gradient on the  $x$   
 324 directly points to a sample  $y$  in the  $p_{\text{data}}$ , where the  $(x, y)$  pairs are consistent with the optimal  
 325 transport plan  $\pi^*$ . All the linear interpolations  $x_t$  between  $x$  and  $y$  satisfy that  $\nabla_{x_k} D_{\phi}^*(x_t) = \frac{y-x}{\|y-x\|}$ .  
 326 It should also be noted that similar results can also be drawn in some variants of WGANs, whose loss  
 327 functions may have a slight difference with standard WGAN [40]. For example, the SNGAN uses  
 328 the hinge loss during the optimization of the discriminator, *i.e.*,  $r(\cdot)$  and  $g(\cdot)$  in Eq. 7 is selected as

329  $\max(0, -1 - u)$  for stabilizing the training procedure. We provide a detailed discussion on several  
 330 surrogate objectives in Appendix. H.

331 The above property of discriminator in WGANs can be interpreted as that given a sample  $x$  from  
 332 generated distribution  $p_{G_\theta}$  we can obtain a corresponding  $y$  in data distribution  $p_{data}$  by directly  
 333 conducting gradient decent with the optimal discriminator  $D_\phi^*$ :

$$y = x + w_x * \nabla_x D_\phi^*, \quad w_x \geq 0 \quad (13)$$

334 It seems to be a simple and appealing solution to improve  $p_{G_\theta}$  with the guidance of discriminator  $D_\phi$ .  
 335 However, the following issues exist:

336 1) there is no theoretical indication on how to set  $w_x$  for each sample  $x$  in generated distribution.  
 337 We noticed that a concurrent work [32] introduce a search process called Discriminator Optimal  
 338 Transport(DOT) by finding the corresponding  $y^*$  through the following:

$$y_x = \arg \min_{\mathbf{y}} \{ \|\mathbf{y} - \mathbf{x}\|_2 - D_\phi^*(\mathbf{y}) \} \quad (14)$$

339 However, it should be noticed that Eq. 14 has a non-unique solution. We further extend the fact into  
 340 the following theorem:

341 **Theorem 3.** *With the  $\pi^*$  and  $D_\phi^*$  as the optimal solutions of the primal problem in Eq. 12 and*  
 342 *Kantorovich-Rubinstein duality of Eq. 12, the distribution  $p_{ot}$  implied by the generated distribution*  
 343  *$p_{G_\theta}$  and the discriminator  $D_\phi^*$  is defined as ( $y_x$  is defined in Eq. 14):*

$$p_{ot}(\mathbf{y}) = \int d\mathbf{x} \delta(\mathbf{y} - y_x) p_{G_\theta}(\mathbf{x})$$

344 *when  $p_{data} \neq p_{G_\theta}$ , there exists infinite numbers of  $p_{ot}$  with  $p_{data}$  as a special case.*

345 Theorem 3 provides a theoretical justification for the poor empirical performance of conducting DOT  
 346 in the sample space, as shown in their paper.

347 2) Another problem lies in that samples distributed outside the generated distribution ( $p_{G_\theta}$ ) are never  
 348 explored during training, which results in much adversarial noise during the gradient-based search  
 349 process, especially when the sample space is high dimensional such as real-world images.

350 To fix the issues mentioned above in leveraging the information of discriminator in Wasserstein  
 351 GANs, we propose viewing the discriminator as an energy function. With the discriminator as an  
 352 energy function, the stationary distribution is unique, and Langevin dynamics can approximately  
 353 conduct sampling from the stationary distribution. Due to the monotonic property of MCMC, there  
 354 will not be issues like setting  $w_x$  in Eq. 13. Besides, the second issue can also be easily solved by  
 355 fine-tuning the energy spaces with contrastive divergence. In addition to the benefits illustrated above,  
 356 if the discriminator is an energy function, the samples from the corresponding energy-based model  
 357 can be obtained through Langevin dynamics by using the gradients of the discriminator which takes  
 358 advantage of the property of discriminator as shown in Theorem 2. With all the facts as mentioned  
 359 above, there is strong motivation to explore further and bridge the gap between discriminator in  
 360 WGAN and the energy-based model.

## 361 B.4 Semi-Amortized Generation with Langevin Dynamics

## 362 B.5 Real-World Image Generation

363 To quantitatively and empirically study the proposed DCD approach, in this section, we conduct  
 364 experiments on unsupervised real-world image generation with DCD and its related counterparts. On  
 365 several commonly used image datasets, experiments demonstrate that our proposed DCD algorithm  
 366 can always achieve better performance on different benchmarks with a significant margin.

## 367 B.6 Experimental setup

368 **Baselines.** We evaluated the following models as our baselines: we take PixelCNN [37], Pix-  
 369 elIQN [28], and MoLM [29] as representatives of other types of generative models. For the energy-  
 370 based model, we compared the proposed method with EBM [8] and NCSN [31]. For GAN models,

---

**Algorithm 1** Discriminator Contrastive Divergence

---

- 1: **Input:** Pretrained generator  $G_\theta$ , discriminator  $D_\phi$ .
  - 2: Set the step size  $\epsilon$ , the length of MCMC steps  $K$  and the total iterations  $T$ .
  - 3: **for** iteration  $i = 1, \dots, T$  **do**
  - 4:   Sample a batch of data samples  $\{x_t\}_{t=1}^m$  for empirical data distribution  $p_{\text{data}}$  and  $\{z_t\}_{t=1}^m$  for the prior distribution  $p(z)$ .
  - 5:   **for** iteration  $l = 1, \dots, K$  **do**
  - 6:     **Pixel Space:**  $G_\theta(z_t)^l = G_\theta(z_t)^{l-1} - \frac{\epsilon}{2} \nabla_x D_\phi(G_\theta(z_t)^{l-1}) + \sqrt{\epsilon} \omega, \omega \sim \mathcal{N}(0, \mathcal{I})$  **or**
  - 7:     **Latent Space:**  $z_t^l = z_t^{l-1} - \frac{\epsilon}{2} \nabla_z D_\phi(G_\theta(z_t)^{l-1}) + \sqrt{\epsilon} \omega, \omega \sim \mathcal{N}(0, \mathcal{I})$
  - 8:   **end for**
  - 9:   Optimized the following objective w.r.t.  $\phi$ :
  - 10:   **Pixel Space:**  $L = \frac{1}{m} \sum_t (D_\phi(x_t) - D_\phi(G_\theta(z_t)^K))$  **or**
  - 11:   **Latent Space:**  $L = \frac{1}{m} \sum_t (D_\phi(x_t) - D_\phi(G_\theta(z_t^K)))$
  - 12: **end for**
- 

Model	Inception	FID
<b>CIFAR-10 Unconditional</b>		
PixelCNN [37]	4.60	65.93
PixelIQN [28]	5.29	49.46
EBM [8]	6.02	40.58
WGAN-GP [13]	7.86 ± .07	18.12
MoLM [29]	7.90 ± .10	18.9
SNGAN [25]	8.22 ± .05	21.7
ProgressiveGAN [18]	8.80 ± .05	-
NCSN [31]	8.87 ± .12	25.32
<hr/>		
DCGAN w/ DRS [2]	3.073	-
DCGAN w/ MH-GAN [35]	3.379	-
ResNet-SAGAN w/ DOT [32]	8.50 ± .12	19.71
<hr/>		
<b>SNGAN-DCD (Pixel)</b>	8.54 ± .11	21.67
<b>SNGAN-DCD (Latent)</b>	<b>9.11 ± .04</b>	<b>16.24</b>
<hr/>		
<b>CIFAR-10 Conditional</b>		
EBM [8]	8.30	37.9
SNGAN [25]	8.43 ± .09	15.43
<b>SNGAN-DCD (Pixel)</b>	8.73 ± .13	22.84
<b>SNGAN-DCD (Latent)</b>	8.81 ± .11	15.05
BigGAN [3]	<b>9.22</b>	<b>14.73</b>

Table 3: Inception and FID scores for CIFAR-10.

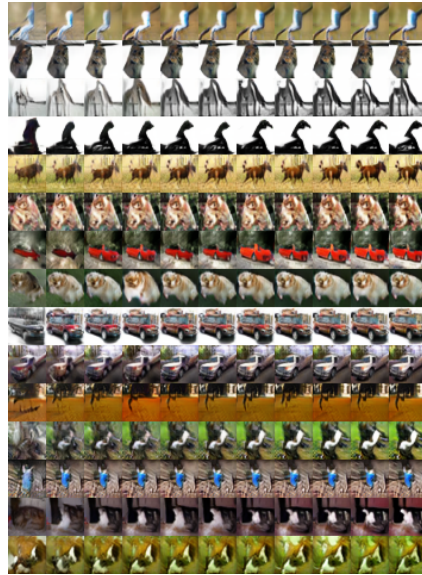


Figure 2: Unconditional CIFAR-10 Langevin dynamics visualization.

371 we take WGAN-GP [13], Spectral Normalization GAN (SNGAN) [25], and Progressiv eGAN [18]  
372 for comparison. We also take the aforementioned DRS [2], DOT [32] and MH-GAN [35] into  
373 consideration. The choices of EBM and GANs are due to their close relation to our proposed method,  
374 as analyzed in Section 2. We omit other previous GAN methods since as a representative of a  
375 state-of-the-art GAN model, SNGAN and Progressive GAN has been shown to rival or outperform  
376 several former methods such as the original GAN [10], the energy-based generative adversarial  
377 network [39], and the original WGAN with weight clipping [1].

378 **Evaluation Metrics.** For evaluation, we concentrate on comparing the quality of generated images  
379 since it is well known that GAN models cannot perform reliable likelihood estimations [34]. We  
380 choose to compare the Inception Scores [30] and Frchet Inception Distances (FID) [15] reached  
381 during training iterations, both computed from 50K samples. A high image quality corresponds to  
382 high Inception and low FID scores. Specifically, the intuition of IS is that high-quality images should  
383 lead to high confidence in classification, while FID aims to measure the computer-vision-specific  
384 similarity of generated images to real ones through Frchet distance.

385 **Data.** We use CIFAR-10 [21], STL-10 [4] and ImageNet [6], which are all standard datasets widely  
386 used in generative literature. STL-10 consists of unlabeled real-world color images, while CIFAR-10

387 and ImageNet is provided with class labels, which enables us to conduct conditional generation tasks.  
 388 For STL-10, we also shrink the images into  $32 \times 32$  as in previous works.

389 **Network Architecture.** For all experiment settings, we follow Spectral Normalization GAN  
 390 (SNGAN) [25] and adopt the same Residual Network (ResNet) [14] structures and hyperparameters,  
 391 which presently is the state-of-the-art implementation of WGAN. Details can be found in Appendix. G.  
 392 We take their open-source code and pre-trained model as the base model for the experiments on  
 393 CIFAR-10 and ImageNet. For STL-10, since there is no pre-trained model available to reproduce the  
 394 results, we train the SNGAN from scratch and take it as the base model.

### 395 B.6.1 Results

396 For quantitative evaluation, we report the inception  
 397 score [30] and FID [15] scores on CIFAR-  
 398 10 in Tab. 3. As shown in the Tab. 3, in pixel  
 399 space, by introducing the proposed DCD algo-  
 400 rithm, we achieve a significant improvement of  
 401 inception score over the SNGAN. The reported  
 402 inception score is even higher than most values  
 403 achieved by class-conditional generative models.  
 404 Our FID score of 21.67 on CIFAR-10 is competi-  
 405 tive with other top generative models. When the  
 406 DCD is conducted in the latent space, we further  
 407 achieve a 9.11 inception score and a 16.24 FID,  
 408 which is a new state-of-the-art performance of  
 409 IS. When combined with label information to  
 410 perform conditional generation, we further im-  
 411 prove the FID to 15.05, which is comparable  
 412 with current state-of-the-art large-scale trained  
 413 models [3]. Some visualization of generated  
 414 examples can be found in Fig 2, which demonstrates that the Markov chain is able to generate  
 415 more realistic samples, suggesting that the MCMC process is meaningful and effective. Tab. 4 and  
 416 Tab. 5 shows the performance on STL-10 and ImageNet respectively, which demonstrate that as a  
 417 generalized method, DCD is not over-fitted to the specific dataset. More experiment details and the  
 418 generated samples can be found in Appendix. I.

Model	Inception	FID
SNGAN [25]	$8.90 \pm .12$	18.73
<b>SNGAN-DCD (Pixel)</b>	$9.25 \pm .09$	22.25
<b>SNGAN-DCD (Latent)</b>	<b><math>9.33 \pm .04</math></b>	<b>17.68</b>

Table 4: Inception and FID scores for STL-10

Model	Inception
cGAN	36.23
cGAN w/ DOT [2]	37.29
SNGAN [25]	36.8
<b>SNGAN-DCD</b>	<b>38.9</b>

Table 5: Inception scores for ImageNet

## 419 C Broader Impact

420 It should be noted that the semi-amortized generation allows a trade-off between the generation  
 421 quality and sampling speed, which holds a slower sampling speed than a direct generation with a  
 422 generator. Hence the proposed method is suitable to the application scenario where the generation  
 423 quality is given vital importance. Another interesting observation during the experiments is the  
 424 discriminator contrastive divergence surprisingly reduces the occurrence of adversarial samples  
 425 during training, so it should be a promising future direction to investigate the relationship between  
 426 our method and bayesian adversarial learning.

427 However, negative consequences also exist since advances in generative models may lead to more  
 428 realistic fake images, which have the capacity to deceive, emotionally distress, and affect public  
 429 opinions and actions. To mitigate the risks associated with deep generative models, we encourage  
 430 researchers to understand and avoid the bad influence of using generative models in particular  
 431 real-world scenarios.

## 432 D Proof of Theorem 2

433 It should be noticed that Theorem. 2 can be generalized to that Lipschitz continuity with  $l_2$ -norm  
 434 (Euclidean Distance) can guarantee that the gradient is directly pointing towards some sample[40].  
 435 We introduce the following lemmas, and Theorem. 2 is a special case.

436 Let  $(x, y)$  be such that  $y \neq x$ , and we define  $x_t = x + t \cdot (y - x)$  with  $t \in [0, 1]$ .

437 **Lemma 1.** If  $f(x)$  is  $k$ -Lipschitz with respect to  $\|\cdot\|_p$  and  $f(y) - f(x) = k\|y - x\|_p$ , then  $f(x_t) =$   
438  $f(x) + t \cdot k\|y - x\|_p$ .

439 *Proof.* As we know  $f(x)$  is  $k$ -Lipschitz, with the property of norms, we have

$$\begin{aligned} f(y) - f(x) &= f(y) - f(x_t) + f(x_t) - f(x) \\ &\leq f(y) - f(x_t) + k\|x_t - x\|_p = f(y) - f(x_t) + t \cdot k\|y - x\|_p \\ &\leq k\|y - x_t\|_p + t \cdot k\|y - x\|_p = k \cdot (1 - t)\|y - x\|_p + t \cdot k\|y - x\|_p \\ &= k\|y - x\|_p. \end{aligned} \quad (15)$$

440  $f(y) - f(x) = k\|y - x\|_p$  implies all the inequalities is equalities. Therefore,  $f(x_t) = f(x) + t \cdot$   
441  $k\|y - x\|_p$ .  $\square$

442 **Lemma 2.** Let  $v$  be the unit vector  $\frac{y-x}{\|y-x\|_2}$ . If  $f(x_t) = f(x) + t \cdot k\|y - x\|_2$ , then  $\frac{\partial f(x_t)}{\partial v}$  equals to  
443  $k$ .

*Proof.*

$$\begin{aligned} \frac{\partial f(x_t)}{\partial v} &= \lim_{h \rightarrow 0} \frac{f(x_t + hv) - f(x_t)}{h} = \lim_{h \rightarrow 0} \frac{f(x_t + h \frac{y-x}{\|y-x\|_2}) - f(x_t)}{h} \\ &= \lim_{h \rightarrow 0} \frac{f(x_t + \frac{h}{\|y-x\|_2}) - f(x_t)}{h} = \lim_{h \rightarrow 0} \frac{\frac{h}{\|y-x\|_2} \cdot k\|y - x\|_2}{h} = k. \quad \square \end{aligned}$$

444 Then we derive the formal proof of Theorem 2.

445 *Proof.* Assume  $p = 2$ , if  $f(x)$  is  $k$ -Lipschitz with respect to  $\|\cdot\|_2$  and  $f(x)$  is differentiable at  $x_t$ ,  
446 then  $\|\nabla f(x_t)\|_2 \leq k$ . Let  $v$  be the unit vector  $\frac{y-x}{\|y-x\|_2}$ . We have

$$k^2 = k \frac{\partial f(x_t)}{\partial v} = k \langle v, \nabla f(x_t) \rangle = \langle kv, \nabla f(x_t) \rangle \leq \|kv\|_2 \|\nabla f(x_t)\|_2 = k^2. \quad (16)$$

447 Because the equality holds only when  $\nabla f(x_t) = kv = k \frac{y-x}{\|y-x\|_2}$ , we have that  $\nabla f(x_t) = k \frac{y-x}{\|y-x\|_2}$ .  
448  $\square$

### 449 E Proof of Theorem 3

450 Theorem. 3 states that following the following procedure as introduced in [32], there is non-unique  
451 stationary distribution. The complete procedure is to find the following  $y$  for  $x \sim P_{G_\theta}$ :

$$y^* = \arg \min_x \{\|x - y\|_2 - D(x)\}. \quad (17)$$

452 To find the corresponding  $y^*$ , the following gradient based update is conducted:

$$\{x \leftarrow x - \epsilon \nabla_x \{\|x - y\|_2 - D(x)\}\}. \quad (18)$$

453 For all the points  $x_t$  in the linear interpolation of  $x$  and target  $y^*$  as defined in the proof of Theorem 2,

$$\nabla_{x_t} \{\|x_t - y\|_2 - D(x_t)\} = \frac{y - x}{\|y - x\|_2} - \frac{y - x}{\|y - x\|_2} = 0, \quad (19)$$

454 which indicates all points in the linear interpolation satisfy the stationary condition.

### 455 F Proof of Proposition 1

456 Proposition. 1 is the direct result of the following Lemma. 3. Following [23], we provide the complete  
457 proof as following.

458 **Lemma 3.** [5] Let  $q$  and  $r$  be two distributions for  $z_0$ . Let  $q_t$  and  $r_t$  be the corresponded distributions  
459 of state  $z_t$  at time  $t$ , induced by the transition kernel  $\mathcal{K}$ . Then  $D_{KL}[q_t|r_t] \geq D_{KL}[q_{t+1}|r_{t+1}]$  for all  
460  $t \geq 0$ .

*Proof.*

$$\begin{aligned}
D_{\text{KL}}[q_t||r_t] &= \mathbb{E}_{q_t} \left[ \log \frac{q_t(\mathbf{z}_t)}{r_t(\mathbf{z}_t)} \right] \\
&= \mathbb{E}_{q_t(\mathbf{z}_t)\mathcal{K}(\mathbf{z}_{t+1}|\mathbf{z}_t)} \left[ \log \frac{q_t(\mathbf{z}_t)\mathcal{K}(\mathbf{z}_{t+1}|\mathbf{z}_t)}{r_t(\mathbf{z}_t)\mathcal{K}(\mathbf{z}_{t+1}|\mathbf{z}_t)} \right] \\
&= \mathbb{E}_{q_{t+1}(\mathbf{z}_{t+1})q_{t+1}(\mathbf{z}_t|\mathbf{z}_{t+1})} \left[ \log \frac{q_{t+1}(\mathbf{z}_{t+1})q(\mathbf{z}_t|\mathbf{z}_{t+1})}{r_{t+1}(\mathbf{z}_{t+1})r(\mathbf{z}_t|\mathbf{z}_{t+1})} \right] \\
&= D_{\text{KL}}[q_{t+1}||r_{t+1}] + \mathbb{E}_{q_{t+1}} D_{\text{KL}}[q_{t+1}(\mathbf{z}_t|\mathbf{z}_{t+1})||r_{t+1}(\mathbf{z}_t|\mathbf{z}_{t+1})].
\end{aligned}$$

461

□

## 462 G Network architectures

463 The ResNet architectures for CIFAR-10 and STL-10 datasets are shown in Tab. 6, which are similar  
464 to the ones in [13]. For the ImageNet datasets, we follow the ResNet architectures in [25]. The details  
465 are shown in Tab. 7.

Table 6: ResNet architectures for CIFAR-10 and STL-10 datasets.

$z \in \mathbb{R}^{128} \sim \mathcal{N}(0, I)$	RGB image $x \in \mathbb{R}^{32 \times 32 \times 3}$
dense, $4 \times 4 \times 256$	ResBlock down 128
ResBlock up 256	ResBlock down 128
ResBlock up 256	ResBlock 128
ResBlock up 256	ResBlock 128
BN, ReLU, $3 \times 3$ conv, 3 Tanh	ReLU
(a) Generator	Global sum pooling
	dense $\rightarrow 1$
	(b) Discriminator

Table 7: ResNet architectures of the Generator for ImageNet dataset. As for the model of the *projection discriminator*, we used the same architecture used in [24]. Please see the paper for the details.

$z \in \mathbb{R}^{128} \sim \mathcal{N}(0, I)$
dense, $4 \times 4 \times 1024$
ResBlock up 1024
ResBlock up 512
ResBlock up 256
ResBlock up 128
ResBlock up 64
BN, ReLU, $3 \times 3$ conv 3
Tanh
(a) Generator

## 466 H Discussions on Objective Functions

467 Optimization of the standard objective of WGAN, *i.e.* with  $r(x) = m(x) = x$  in Eq. 7, are found  
468 to be unstable due to the numerical issues and free offset [40, 25]. Instead, several surrogate losses  
469 are actually used in practice. For example, the logistic loss( $r(x) = m(x) = -\log(1 + e^{-x})$ ) and  
470 hinge loss( $r(x) = m(x) = \min(0, x)$ ) are two widely applied objectives. Such surrogate losses are  
471 valid due to that they are actually the lower bounds of the Wasserstein distance between the two  
472 distributions of interest. The statement can be easily derived by the fact that  $-\log(1 + e^{-x}) \leq x$  and  
473  $\min(0, x) \leq x$ . A more detailed discussion could also be found in [32].

474 Note that  $\min(0, -1 + x)$  and  $-\log(1 + e^{-x})$  are in the function family proposed in [40], and  
475 Theorem 4 in [40] guarantees the gradient property of discriminator.

## 476 I More Experiment Details

### 477 I.1 CIFAR-10

478 For the meta-parameters in DCD Algorithm 1, when the MCMC process is conducted in the pixel  
479 space, we choose 6–8 as the number of MCMC steps  $K$ , and set the step size  $\epsilon$  as 10 and the standard  
480 deviation of the Gaussian noise as 0.01, while for the latent space we set  $K$  as 50,  $\epsilon$  as 0.2 and the  
481 deviation as 0.1. Adam optimizer [20] is set with  $2 \times 10^{-4}$  learning rate with  $\beta_1 = 0, \beta_2 = 0.9$ . We  
482 use 5 critic updates per generator update, and a batch size of 64.

### 483 I.2 STL-10

484 We show generated samples of DCD during Langevin dynamics in Fig. 3. We run 150 steps of  
485 MCMC steps and plot generated samples for every 10 iterations. The step size is set as 0.05 and the  
486 noise is set as  $N(0, 0.1)$ .

### 487 I.3 ImageNet

488 We show generated samples of DCD during Langevin dynamics in Fig. 4. We run 1000 Langevin  
489 dynamics steps and plot generated samples for every 100 iterations. The initial step size and the  
490 Gaussian noise are set as 0.05 and  $N(0, 0.1)$  respectively. The step size and standard deviation of  
491 Gaussian noise are simultaneously decayed with a factor 0.3 for every 100 iterations.

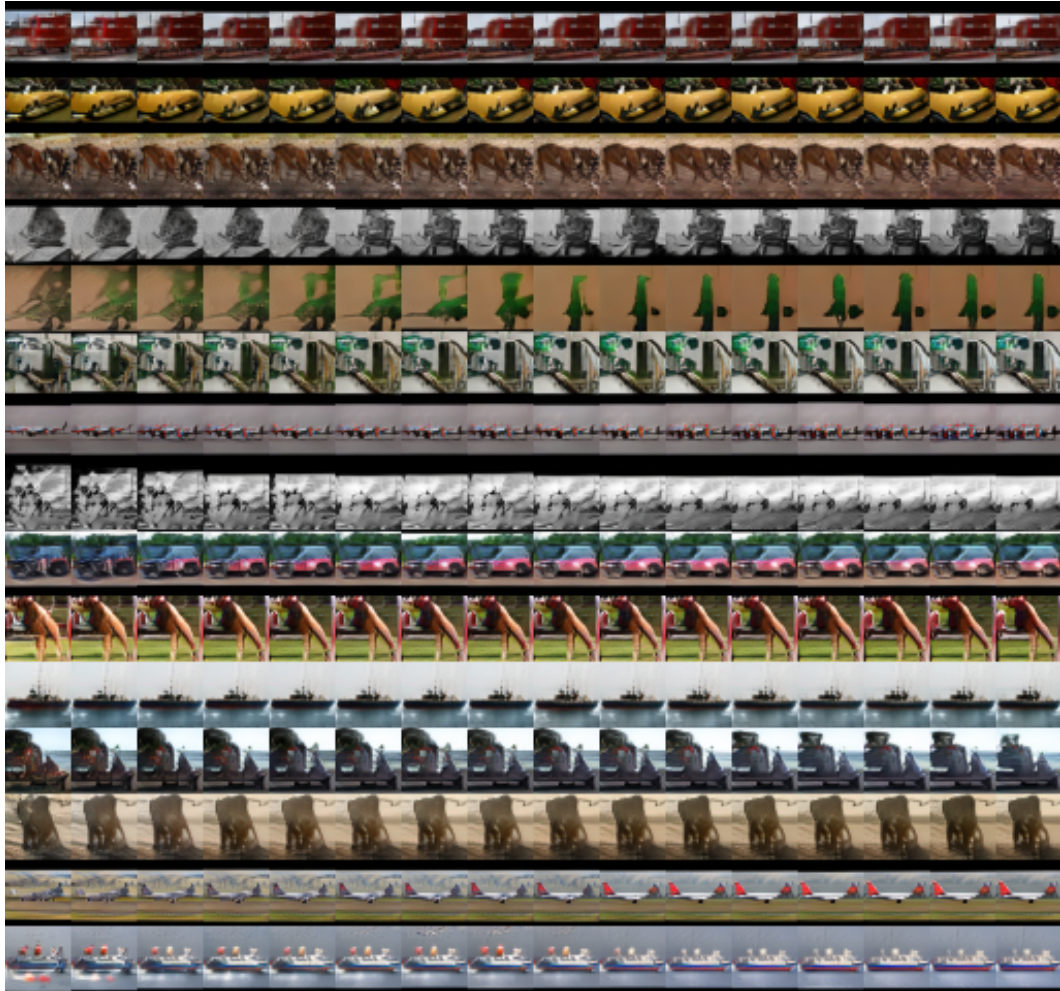


Figure 3: STL-10 Langevin dynamics visualization.





Figure 4: ImageNet Langevin dynamics visualization.