

# IMPROVED DIFFUSION-BASED GENERATIVE MODEL WITH BETTER ADVERSARIAL ROBUSTNESS

Zekun Wang<sup>\*1</sup>, Mingyang Yi<sup>\*†2</sup>, Shuchen Xue<sup>3,4</sup>, Zhenguo Li<sup>5</sup>, Ming Liu<sup>†1</sup>, Bing Qin<sup>1</sup>, Zhi-Ming Ma<sup>3,4</sup>

<sup>1</sup>Harbin Institute of Technology    <sup>2</sup>Renmin University of China

<sup>3</sup>Academy of Mathematics and Systems Science, Chinese Academy of Sciences

<sup>4</sup>University of Chinese Academy of Sciences    <sup>5</sup>Huawei Noah's Ark Lab

zkwang@ir.hit.edu.cn    yimingyang@ruc.edu.cn

xueshuchen17@mailsucas.ac.cn

## ABSTRACT

Diffusion Probabilistic Models (DPMs) have achieved significant success in generative tasks. However, their training and sampling processes suffer from the issue of distribution mismatch. During the denoising process, the input data distributions differ between the training and inference stages, potentially leading to inaccurate data generation. To obviate this, we analyze the training objective of DPMs and theoretically demonstrate that this mismatch can be alleviated through Distributionally Robust Optimization (DRO), which is equivalent to performing robustness-driven Adversarial Training (AT) on DPMs. Furthermore, for the recently proposed Consistency Model (CM), which distills the inference process of the DPM, we prove that its training objective also encounters the mismatch issue. Fortunately, this issue can be mitigated by AT as well. Based on these insights, we propose to conduct efficient AT on both DPM and CM. Finally, extensive empirical studies validate the effectiveness of AT in diffusion-based models. The code is available at [https://github.com/kugwzk/AT\\_Diff](https://github.com/kugwzk/AT_Diff).

## 1 INTRODUCTION

Diffusion Probabilistic Models (DPMs) (Ho et al., 2020; Song et al., 2020; Yi et al., 2024) have achieved remarkable success across a wide range of generative tasks such as image synthesis (Dhariwal & Nichol, 2021; Rombach et al., 2022; Ho et al., 2022a), video generation (Ho et al., 2022b; Blattmann et al., 2023), text-to-image generation (Nichol et al.; Ramesh et al., 2022; Saharia et al., 2022), *etc.* The core mechanism of DPMs involves a forward diffusion process that progressively injects noise into the data, followed by a reverse process that learns to generate data by denoising the noise. Unlike traditional generative models such as GANs (Goodfellow et al., 2014) or VAEs (Kingma & Welling, 2013), which directly map an easily sampled latent variable (e.g., Gaussian noise) to the target data through a single network function evaluation (NFE), DPMs adopt a gradual denoising approach that requires multiple NFEs (Song et al., 2022; Salimans & Ho, 2022; Lu et al., 2022b; Ma et al., 2024). However, this noising-then-denoising process introduces a distribution mismatch between the training and sampling stages, potentially leading to inaccuracies in the generated outputs.

Concretely, during the training stage, the model is learned to predict the noise in ground-truth noisy data derived from the training set. In contrast, during the inference stage, the input distribution is obtained from the output generated by the DPM in the previous step, which differs from the training phase, caused by the inaccurate estimation of the score function due to training (Song et al., 2021; Yi et al., 2023a) and the discretization error (Chen et al., 2022; Li et al., 2023; Xue et al., 2024b;a) brought by sampling. Such distribution mismatches are referred to as *Exposure Bias*, which has been discussed in auto-regressive language models (Bengio et al., 2015; Ranzato et al., 2016).

Recently, the aforementioned distribution mismatch problem in diffusion has been also recognized by (Ning et al., 2023; Li & van der Schaar, 2024; Ren et al., 2024; Ning et al., 2024; Li et al.,

<sup>\*</sup>Equal contribution

<sup>†</sup>Corresponding author

2024; Lou & Ermon, 2023). However, these studies either rely on strong mismatch distributional assumptions (e.g., Gaussian) (Ning et al., 2023; 2024; Ren et al., 2024) or incur significant additional computational costs (Li & van der Schaar, 2024). This indicates that a more practical solution to this problem has been overlooked until now. To bridge this gap, we begin with the discrete DPM introduced in (Ho et al., 2020). Intuitively, although there is a mismatch between training and inference, the distributions of intermediate noise generated during the inference stage are close to the ground-truth distributions observed during training. Therefore, improving the distributional robustness (Yi et al., 2021; Namkoong, 2019; Shapiro, 2017) (which measures the robustness of the model to distributional perturbations in training data) of DPM mitigates the distribution mismatch problem. To achieve this, we refer to Distribution Robust Optimization (DRO) (Shapiro, 2017; Namkoong, 2019), which aims to improve the distributional robustness of models. We then prove that applying DRO to DPM is mathematically equivalent to implementing *robustness-driven* Adversarial Training (AT) (Madry et al., 2018; Shafahi et al., 2019; Yi et al., 2021) on DPM.<sup>1</sup> Following the DRO framework, we also analyze the recently proposed diffusion-based Consistency Model (CM) (Song et al., 2023; Luo et al., 2023) which distills the trajectory of DPM into a model with one NFE generation. We first prove that the training objective of CM similarly suffers from the mismatch issue as in multi-step DPM. Moreover, the issue can also be mitigated by implementing AT. Therefore, for both DPM and CM, we propose to apply efficient AT (e.g., “Free-AT” (Shafahi et al., 2019)) during their training stages to mitigate the distribution mismatch problem.<sup>2</sup> Finally, we summarize our contributions as follows.

- We conduct an in-depth analysis of the diffusion-based models (DPM and CM) from a theoretical perspective and systematically characterize its distribution mismatch problem.
- For both DPM and CM, we theoretically show that their mismatch problem is mitigated by DRO, which is equivalent to implementing AT with proved error bounds during training.
- We propose to conduct efficient AT on both DPM and CM in various tasks, including image generation on CIFAR10  $32\times 32$  (Krizhevsky & Hinton, 2009) and ImageNet  $64\times 64$  (Deng et al., 2009), and zero-shot Text-to-Image (T2I) generation on MS-COCO  $512\times 512$  (Lin et al., 2014b). Extensive experimental results illustrate the effectiveness of the proposed AT training method in alleviating the distribution mismatch of DPM and CM.

## 2 RELATED WORK

**Distribution Mismatch in DPM.** The problem is analogous to the exposure bias in auto-regressive language models (Bengio et al., 2015; Ranzato et al., 2016; Shen et al., 2016; Rennie et al., 2017; Zhang et al., 2019c), whereas the next word prediction (Radford et al., 2019) relies on tokens predicted by the model in the inference stage, which may be mismatched with the ground-truth one taken in the training stage. The similarity to DPMs becomes evident due to their gradual denoising generation process. Ning et al. (2023) and Ning et al. (2024) propose adding extra Gaussian perturbation during the training stage or data-dependent perturbation during the inference stage, to mitigate this issue. Following this line of work, several methods are further proposed. For instance, to reduce the accumulated discrepancy between the intermediate noisy data in the training and inference stages, Li et al. (2024) search for a suboptimal mismatched input time step of the model to conduct inference. Similarly, Li & van der Schaar (2024) and Ren et al. (2024) directly minimize the difference between the generated intermediate noisy data and the ground-truth data. However, these methods either rely on strong assumptions (Ning et al., 2023; 2024; Li et al., 2024; Ren et al., 2024) or are computationally expensive (Li & van der Schaar, 2024). In contrast, we are the first to explore the distribution mismatch problem from the perspective of DRO. Meanwhile, our proposed AT with strong theoretical foundations is both simple and efficient, compared with the existing methods.

**Adversarial Training and DRO.** In this paper, we leverage the Distributionally Robust Optimization (DRO) (Shapiro, 2017; Namkoong, 2019; Yi et al., 2021; Sinha et al., 2018; Wang et al., 2022; Yi et al., 2023b) to improve the distributional robustness of DPM and CM, thereby mitigating

<sup>1</sup>Note that the “adversarial” here refers to perturbation to input training data, instead of the adversarial of generator-discriminator in GAN (Goodfellow et al., 2014).

<sup>2</sup>Notably, the standard AT (Madry et al., 2018) solves a minimax problem that slows the training process. The efficient AT has no extra computational cost compared to the standard training ones (Shafahi et al., 2019).

the distribution mismatch problem. As demonstrated in (Sinha et al., 2018; Yi et al., 2021; Lee & Raginsky, 2018), we link the DRO with AT (Madry et al., 2018; Goodfellow et al., 2015), which is designed to improve the input (instead of distributional) robustness of the model. In supervised learning, the adversarial examples generated by efficient AT methods (Shafahi et al., 2019; Zhang et al., 2019a;b; Zhu et al., 2020; Jiang et al., 2020) have been proven to be efficient augmented data to improve the robustness and generalization performance of models (Rebuffi et al., 2021; Wu et al., 2020; Yi et al., 2021). In this paper, we further verify that the AT generated adversarial augmented examples are also beneficial for generative models DPM and CM.

In addition, recent studies (Nie et al., 2022; Wang et al., 2023; Zhang et al., 2023) utilize DPM to generate examples in adversarial training to improve the robustness of the classification model. This is quite different from the method in this paper, as we focus on employing AT during training of diffusion-based model to improve its distributional robustness to alleviate the distribution mismatching.

### 3 PRELIMINARY

**Diffusion Probabilistic Models.** DPM (Sohl-Dickstein et al., 2015; Ho et al., 2020) constructs the Markov chain  $\mathbf{x}_t$  by transition kernel  $q(\mathbf{x}_{t+1} | \mathbf{x}_t) = \mathcal{N}(\sqrt{\alpha_{t+1}}\mathbf{x}_t, (1 - \alpha_{t+1})\mathbf{I})$ , where  $\alpha_1, \dots, \alpha_T$  are in  $[0, 1]$ . Let  $\bar{\alpha}_t := \prod_{s=1}^t \alpha_s$ , and  $\mathbf{x}_0 \sim q$  be ground-truth data. Then, for  $\mathbf{x}_t$ , it holds

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}_t \quad t = 1, \dots, T, \quad (1)$$

with  $\boldsymbol{\epsilon}_t \sim \mathcal{N}(0, \mathbf{I})$ . The reverse process  $p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1})$  is parameterized as

$$p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1}) = \mathcal{N}(\mu_\theta(\mathbf{x}_{t+1}, t + 1), \sigma_{t+1}^2 \mathbf{I}), \quad (2)$$

where  $\sigma_{t+1}^2 = 1 - \alpha_{t+1}$ . To learn  $p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1})$ , a standard method is to minimize the following evidence lower bound of negative log-likelihood (NLL) (Ho et al., 2020),

$$-\mathbb{E}_q[\log p_\theta(\mathbf{x}_0)] \leq \mathbb{E}_q \left[ -\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \right]. \quad (3)$$

Here, minimizing the ELBO in the r.h.s. of above inequality links to  $p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1})$  since it is equivalent to minimizing the following rewritten objective

$$\min_{\theta} \left\{ D_{KL}(q(\mathbf{x}_T) \| p_\theta(\mathbf{x}_T)) + \sum_{t=0}^{T-1} \underbrace{D_{KL}(q(\mathbf{x}_t | \mathbf{x}_{t+1}) \| p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1}))}_{L_t} \right\}, \quad (4)$$

as in (Ho et al., 2020; Bao et al., 2022; Yi et al., 2023a). Here, the conditional Kullback–Leibler (KL) divergence  $D_{KL}(q(\mathbf{x}_t | \mathbf{x}_{t+1}) \| p(\mathbf{x}_t | \mathbf{x}_{t+1})) = \int q(\mathbf{x}_t | \mathbf{x}_{t+1}) \log \frac{q(\mathbf{x}_t | \mathbf{x}_{t+1})}{p(\mathbf{x}_t | \mathbf{x}_{t+1})} d\mathbf{x}_t d\mathbf{x}_{t+1}$  (Duchi, 2016), and minimizing  $L_t$  is equivalent to solve the following noise prediction problem

$$\min_{\theta} \mathbb{E} \left[ \|\boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}_t, t) - \boldsymbol{\epsilon}_t\|^2 \right]. \quad (5)$$

We use  $\|\cdot\|_p$  to denote  $\ell_p$ -norm. Unless specified, the norm  $\|\cdot\|$  refers to the  $\ell_2$ -norm  $\|\cdot\|_2$ . Since  $\bar{\alpha}_t \rightarrow 0$  for  $t \rightarrow T$ ,  $\mathbf{x}_0$  is obtained by conducting the reverse diffusion process  $p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1})$  starting from  $\mathbf{x}_T \sim \mathcal{N}(0, \mathbf{I})$  and  $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I})$ , under the learned model  $\boldsymbol{\epsilon}_\theta$  with

$$\mathbf{x}_t = \frac{1}{\sqrt{\alpha_{t+1}}} \left( \mathbf{x}_{t+1} - \frac{1 - \alpha_{t+1}}{\sqrt{1 - \bar{\alpha}_{t+1}}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_{t+1}, t + 1) \right) + \sqrt{1 - \alpha_{t+1}}\boldsymbol{\epsilon}. \quad (6)$$

**Wasserstein Distance.** For integer  $p > 0$ ,  $\Gamma(\mu, \nu)$  as the set of union distributions with marginal  $\mu$  and  $\nu$ , the Wasserstein  $p$ -distance (Villani et al., 2009) between distributions  $\mu$  and  $\nu$  with finite  $p$ -moments is

$$W_p^p(\mu, \nu) = \inf_{\gamma \in \Gamma(\mu, \nu)} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \gamma} \|\mathbf{x} - \mathbf{y}\|_p^p. \quad (7)$$

### 4 ROBUSTNESS-DRIVEN ADVERSARIAL TRAINING OF DIFFUSION MODELS

In this section, we formally show that the success of DPM relies on specific conditions, i.e.,  $\mathbf{x}_t$  is close to  $\mathbf{x}_{t+1}$ . Next, to mitigate the drawbacks brought by the restriction, we propose to consider the distribution mismatch problem as discussed in Section 1, and connect the problem to a rewritten ELBO. Finally, we apply DRO for this ELBO to mitigate the distribution mismatch problem and finally link it to AT to be implemented in practice.

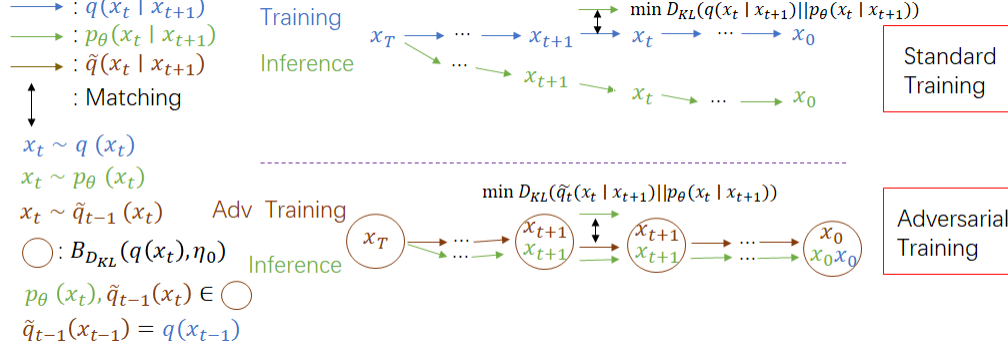


Figure 1: A comparison between standard training and the proposed distributional robust optimization in (12). When minimizing  $D_{KL}(\tilde{q}_t(x_t | x_{t+1}) || p_\theta(x_t | x_{t+1}))$ , the  $x_{t+1}$  is sampled from  $\tilde{q}_t(x_{t+1})$ , such that both  $\tilde{q}_t(x_{t+1})$  in training stage and  $p_\theta(x_{t+1})$  in inference stage are in  $B_{D_{KL}}(q(x_{t+1}), \eta_0)$ , so that  $p_\theta(x_t)$  tends to locates in  $B_{D_{KL}}(q(x_t), \eta_0)$  as well as  $\tilde{q}_t(x_t)$ . Then, the distributional robustness captured by (12) guarantees the generated  $p_\theta(x_t)$  always locates around  $q(x_t)$  for all  $t$ .

#### 4.1 HOW DOES DPM WORKS IN PRACTICE?

Notably, minimizing (4) potentially obtains a sharp NLL under target distribution  $q(x_0)$ . However, in the following proposition, we show that (4) also implicitly minimizes the NLL of each  $x_t$ .

**Proposition 1.** *The minimization problem (4) is equivalent to minimizing an upper bound of  $\mathbb{E}_q[-\log p_\theta(x_t)]$  for any  $0 \leq t \leq T$ .*

The proof is provided in Appendix A. It shows that though (4) is proposed to generate  $x_0 \sim q(x_0)$ , it also guides the model to generate  $x_t$  such that  $p_\theta(x_t)$  approximates the ground-truth distribution  $q(x_t)$ . The conclusion is nontrivial as minimizing the ELBO of NLL  $\mathbb{E}_q[-\log p_\theta(x_0)]$  does not necessarily impose any restrictions on  $x_t$  for  $t \geq 1$ .

Next, we will further explain why (4) leads to a small NLL of  $x_t$ . In  $L_t$  of (4),  $p_\theta(x_t | x_{t+1})$  approximates  $q(x_t | x_{t+1})$  with  $x_{t+1} \sim q(x_{t+1})$  representing ground-truth data. Consequently,  $p_\theta(x_t)$  approximates  $q(x_t)$  by recursively applying such a relationship as in the following proposition.

**Proposition 2.** *Suppose  $p_\theta(x_t | x_{t+1})$  matches  $q(x_t | x_{t+1})$  well such that*

$$L_t = D_{KL}(q(x_t | x_{t+1}) || p_\theta(x_t | x_{t+1})) \leq \frac{\gamma}{T}, \quad (8)$$

and the discrepancy satisfies  $D_{KL}(q(x_T) || p_\theta(x_T)) \leq \gamma_0$ , then for any  $0 \leq t \leq T$ , we have

$$D_{KL}(q(x_t) || p_\theta(x_t)) \leq D_{KL}(q(x_T) || p_\theta(x_T)) + L_t \leq \gamma_0 + \frac{(T-t)\gamma}{T}. \quad (9)$$

The results is similarly obtained in (Chen et al., 2023), while their result is applied for  $D_{KL}(q(x_0) || p_\theta(x_0))$ , which is narrowed compared with Proposition 2. The proof is provided in Appendix A, which formally explains why (4) results in  $p_\theta(x_t)$  approximating  $q(x_t)$ . However, this proposition is built upon small  $L_t$ , and notably, the error introduced by  $L_t$  will be accumulated on the r.h.s. of (9), as it increases w.r.t.  $t$ . This phenomenon is caused by the *distribution mismatch problem* discussed in Section 1. Concretely, in (4), minimizing  $L_t$  learns the transition probability  $p_\theta(x_t | x_{t+1})$  based on  $x_{t+1} \sim q(x_{t+1})$ , while in practice,  $x_t$  in (6) is generated from  $x_{t+1} \sim p_\theta(x_{t+1})$ . The error between  $p_\theta(x_{t+1})$  and  $q(x_{t+1})$  will propagates into the error between  $p_\theta(x_t)$  and  $q(x_t)$  as in (9).

Therefore, owing to the existence of distribution mismatch, only if  $L_t$  is minimized, the gap between  $p_\theta(x_t)$  and  $q(x_t)$  can be guaranteed. However, the following proposition proved in Appendix A indicates that  $L_t$  is theoretically minimized with restrictions.

**Proposition 3.**  *$L_t$  in (4) is well minimized, only if  $q(x_{t+1})$  is Gaussian or  $\|x_{t+1} - x_t\| \rightarrow 0$ .*

In practice, the  $q(x_{t+1})$  is usually non-Gaussian. Besides, the gap  $\|x_{t+1} - x_t\|$  is not necessarily small, especially for samplers with few sampling steps, e.g., DDIM (Song et al., 2022), DPM-Solver (Lu et al., 2022a). Therefore, in practice, the accumulated error in (9) caused by the distribution mismatch problem may become large, and degenerate the quality of  $x_0$ .

## 4.2 DISTRIBUTIONAL ROBUSTNESS IN DPM

Inspired by the discussion above, we propose a new training objective as the sum of NLLs under  $\mathbf{x}_t$ ,

$$\min_{\theta} \mathcal{L}(\theta) = \sum_{t=0}^T \mathbb{E}_q [-\log p_{\theta}(\mathbf{x}_t)]. \quad (10)$$

Then the following proposition constructs ELBOs for each of  $\mathbb{E}_q[-\log p_{\theta}(\mathbf{x}_t)]$ .

**Proposition 4.** *For any distribution  $\tilde{q}$  satisfies  $\tilde{q}(\mathbf{x}_t) = q(\mathbf{x}_t)$  for specific  $t$ , we have*

$$\mathbb{E}_q [-\log p_{\theta}(\mathbf{x}_t)] \leq \underbrace{D_{KL}(\tilde{q}(\mathbf{x}_t | \mathbf{x}_{t+1}) \| p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t+1}))}_{L_t^{\tilde{q}}} + C, \quad (11)$$

for a constant  $C$  independent of  $\theta$ .

The proof is in Appendix A.2. This proposition generalizes the results in Proposition 1 since  $\tilde{q}$  can be taken as  $q$  in Proposition 1. During minimizing  $L_t^{\tilde{q}}$ , the transition probability  $p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t+1})$  matches  $\tilde{q}(\mathbf{x}_t | \mathbf{x}_{t+1})$ , while  $\mathbf{x}_{t+1} \sim \tilde{q}(\mathbf{x}_{t+1})$  in the training stage has no restriction. Thus, one may take  $\tilde{q}(\mathbf{x}_{t+1}) \approx p_{\theta}(\mathbf{x}_{t+1})$ , then in  $L_t^{\tilde{q}}$ ,  $p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t+1})$  matches  $\tilde{q}(\mathbf{x}_t | \mathbf{x}_{t+1})$  leads  $p_{\theta}(\mathbf{x}_t) \approx \tilde{q}(\mathbf{x}_t) = q(\mathbf{x}_t)$ , which mitigates the distribution mismatch problem, when minimizing such  $L_t^{\tilde{q}}$ .

Unfortunately, for each  $t$ , obtaining such specific  $\tilde{q}_t(\mathbf{x}_{t+1}) = p_{\theta}(\mathbf{x}_{t+1})$  is computationally expensive (Li & van der Schaar, 2024), which prevents us using desired  $\tilde{q}_t(\mathbf{x}_{t+1})$ . However, we know  $p_{\theta}(\mathbf{x}_{t+1})$  is around  $q(\mathbf{x}_{t+1})$ . Therefore, by borrowing the idea from DRO (Shapiro, 2017), for each  $t$ , we propose to minimize the maximal value of  $L_t^{\tilde{q}_t}$  over all possible  $\tilde{q}_t(\mathbf{x}_{t+1})$  around  $q(\mathbf{x}_{t+1})$ . This leads to a small  $L_t^{p_{\theta}}$ , as  $p_{\theta}(\mathbf{x}_{t+1})$  locates around  $q(\mathbf{x}_{t+1})$ , so that is included in the ‘‘maximal range’’. Technically, the DRO-based EBLO of (11) is formulated as follows. Here  $p_{\theta}(\mathbf{x}_{t+1})$  is supposed in  $B_{D_{KL}}(q(\mathbf{x}_{t+1}), \eta_0)$ , and it captures the distributional robustness of  $p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t+1})$  w.r.t. input  $\mathbf{x}_{t+1}$ .

$$\begin{aligned} \min_{\theta} \sum_{t=0}^{T-1} L_t^{\text{DRO}}(\theta) &= \min_{\theta} \sum_{t=0}^{T-1} \sup_{\tilde{q}_t(\mathbf{x}_{t+1}) \in B_{D_{KL}}(q(\mathbf{x}_{t+1}), \eta_0)} D_{KL}(\tilde{q}_t(\mathbf{x}_t | \mathbf{x}_{t+1}) \| p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t+1})); \\ \text{s.t.} \quad \tilde{q}_t(\mathbf{x}_t) &= q(\mathbf{x}_t). \end{aligned} \quad (12)$$

Here  $\tilde{q}_t(\mathbf{x}_{t+1}) \in B_{D_{KL}}(q(\mathbf{x}_{t+1}), \eta_0)$  means  $D_{KL}(q(\mathbf{x}_{t+1}) \| \tilde{q}_t(\mathbf{x}_{t+1})) \leq \eta_0$ . By solving problem (12), if the desired  $\tilde{q}_t(\mathbf{x}_{t+1}) = p_{\theta}(\mathbf{x}_{t+1})$  is in  $B_{D_{KL}}(q(\mathbf{x}_{t+1}), \eta_0)$ , then the conditional probability in (12) transfers  $\mathbf{x}_{t+1} \sim p_{\theta}(\mathbf{x}_{t+1})$  to target  $\mathbf{x}_t \sim q(\mathbf{x}_t)$  is learned, which mitigates the distribution mismatch problem. The theoretical clarification is in the following Proposition proved in Appendix A.2, which indicates that small DRO loss (12) guarantees the quality of generated  $\mathbf{x}_0$ .

**Proposition 5.** *If  $L_t^{\text{DRO}}(\theta) \leq \eta_0$  in (12) for all  $t$ , and  $D_{KL}(q(\mathbf{x}_T) \| p_{\theta}(\mathbf{x}_T)) \leq \eta_0$ , then  $D_{KL}(q(\mathbf{x}_0) \| p_{\theta}(\mathbf{x}_0)) \leq \eta_0$ .*

Up to now, we do not know how to compute the DRO-based training objective (12) we derived. Fortunately, the following theorem corresponds (12) to a ‘‘perturbed’’ noise prediction problem similar to (5). The theorem is proved in Appendix A.2.

**Theorem 1.** *There exists  $\delta_t$  depends on  $\mathbf{x}_0$  and  $\epsilon_t$  makes (13) equivalent to problem (12).*

$$\min_{\theta} \sum_{t=0}^{T-1} \mathbb{E}_{q(\mathbf{x}_0), \epsilon_t} \left[ \left\| \epsilon_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_t + \delta_t, t) - \epsilon_t - \frac{\delta_t}{\sqrt{1 - \bar{\alpha}_t}} \right\|^2 \right]. \quad (13)$$

This theorem connects the proposed DRO problem (12) with noise prediction problem (13). Naturally, we can solve (13), if we know the exact  $\delta_t$ . Fortunately, we have the following proposition to characterize the range of  $\delta_t$ , and it is proved in Appendix A.2.

**Proposition 6.** *For  $\eta > 0$  and  $\delta_t$  in (13),  $\|\delta_t\|_1 \leq \eta$  holds with probability at least  $1 - \sqrt{2(1 - \bar{\alpha}_t)}/\eta$ .*

The proposition indicates that for any  $\delta_t$  depends on  $\mathbf{x}_0, \epsilon_t$  in (13), it is likely in a small range (measured under any  $\ell_p$ -norm, since they can bound each other in Euclidean space). Thus, to resolve (13) (so that (12)), we propose to directly consider the following adversarial training (Madry et al.,

2018) objective with the perturbation  $\delta$  is taken over its possible range as proved in Proposition 6, which captures the input (instead of distribution) robustness of model  $\epsilon_\theta$ .

$$\min_{\theta} \sum_{t=0}^{T-1} \mathbb{E}_{q(\mathbf{x}_0)} \left[ \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} \left[ \sup_{\delta: \|\delta\| \leq \eta} \left\| \epsilon_\theta(\sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1-\alpha_t} \epsilon_t + \delta) - \epsilon_t - \frac{\delta}{\sqrt{1-\alpha_t}} \right\|^2 \right] \right]. \quad (14)$$

We present a fine-grained connection between (14) and classical AT in Appendix C. Notably, our objective (14) is different from the ones in (Ning et al., 2023), whereas  $\delta$  in it is a Gaussian, and  $\epsilon_\theta$  predicts  $\epsilon_t$  instead of  $\epsilon_t + \delta/\sqrt{1-\alpha_t}$  as ours.

To make it clear, we summarize the rationale from DRO objective (12) to AT our objective (14). Since Theorem 1 shows solving (12) is equivalent to (13), which conducts noise prediction (5) with a perturbation  $\delta_t$  in a small range added (Proposition 6). Thus, we propose to minimize the maximal loss over the possible  $\delta_t$ , which is indeed our AT objective (14).

## 5 ADVERSARIAL TRAINING UNDER CONSISTENCY MODEL

Although the DPM generates high-quality target data  $\mathbf{x}_0$ , the multi-step denoising process (6) requires numerous model evaluations, which can be computationally expensive. To resolve this, the diffusion-based consistency model (CM) is proposed in (Song et al., 2023). Consistency model  $f_\theta(\mathbf{x}_t, t)$  transfers  $\mathbf{x}_t \sim q(\mathbf{x}_t)$  into a distribution that approximates the target  $q(\mathbf{x}_0)$ .  $f_\theta$  is optimized by the following consistency distillation (CD) loss<sup>3</sup>

$$\min_{\theta} \mathcal{L}_{CD}(\theta) = \sum_{t=0}^{T-1} \mathbb{E}_{\mathbf{x}_{t+1} \sim q(\mathbf{x}_{t+1})} [d(f_\theta(\Phi_t(\mathbf{x}_{t+1}), t), f_\theta(\mathbf{x}_{t+1}, t+1))], \quad (15)$$

where  $\Phi_t(\mathbf{x}_{t+1})$  is a solution of a specific ordinary differential equation (ODE) ((37) in Appendix B) which is a deterministic function transfers  $\mathbf{x}_{t+1}$  to  $\mathbf{x}_t$ , i.e.,  $\Phi_t(\mathbf{x}_{t+1}) \sim q(\mathbf{x}_t)$ , and  $d(\mathbf{x}, \mathbf{y})$  is a distance between  $\mathbf{x}$  and  $\mathbf{y}$  e.g.,  $\ell_1, \ell_2$  distance.

**Remark 1.** In (Song et al., 2023; Luo et al., 2023), the noisy data  $\mathbf{x}_t$  in (15) is described by an ODE (37) in Appendix B. However, we use the discrete  $\mathbf{x}_t$  (1) here to unify the notations with Section 4. The two frameworks are mathematically equivalent as all  $\mathbf{x}_t$  in (1) located in the trajectory of ODE in (Song et al., 2023). More details of this claim refer to Appendix B.

Next, we use the following theorem to illustrate that solving problem (15) indeed creates  $f_\theta(\mathbf{x}_t, t)$  with distribution close target  $q(\mathbf{x}_0)$ . The theorem is proved in Appendix B.

**Theorem 2.** For  $\mathcal{L}_{CD}(\theta)$  in (15) with  $d(\cdot, \cdot)$  is  $\ell_2$  distance, then  $W_1(f_\theta(\mathbf{x}_t, t), \mathbf{x}_0) \leq \sqrt{t \mathcal{L}_{CD}(\theta)}$ <sup>4</sup>.

Though solving problem (15) creates the desired CM  $f_\theta$ , computing the exact  $\Phi_t(\mathbf{x}_{t+1})$  involves solving an ODE as pointed out in Appendix B. Thus, in practice (Song et al., 2023; Luo et al., 2023), the  $\Phi_t(\mathbf{x}_{t+1})$  is approximated by a computable numerical estimation  $\hat{\Phi}_t(\mathbf{x}_{t+1}, \epsilon_\phi)$  of it, e.g., Euler ((42) in Appendix B.1) or DDIM (Song et al., 2023), where  $\epsilon_\phi$  is a pretrained noise prediction model as in (5). Therefore, the practical training objective of (15) becomes

$$\min_{\theta} \sum_{t=0}^{T-1} \hat{\mathcal{L}}_{CD}(\theta) = \mathbb{E}_{\mathbf{x}_{t+1} \sim q(\mathbf{z}_t)} \left[ d \left( f_\theta(\hat{\Phi}_t(\mathbf{x}_{t+1}, \epsilon_\phi), t), f_\theta(\mathbf{x}_{t+1}, t+1) \right) \right]. \quad (16)$$

In (16),  $\hat{\Phi}_t(\mathbf{x}_{t+1}, \epsilon_\phi)$  is an estimation to  $\Phi_t(\mathbf{x}_{t+1})$ , which causes an inaccurate training objective  $\hat{\mathcal{L}}_{CD}$  in (16), compared with target  $\mathcal{L}_{CD}$  (15). Thus, this results in the distribution mismatch problem in CM, as in DPM of Section 4. However, similar to Section 4.2, if we train  $f_\theta$  with robustness to the gap between  $\hat{\Phi}_t(\mathbf{x}_{t+1}, \epsilon_\phi)$  and  $\Phi_t(\mathbf{x}_{t+1})$ , the distribution mismatch problem in CM is mitigated.

Technically, suppose  $\Phi_t(\mathbf{x}_{t+1}) = \hat{\Phi}_t(\mathbf{x}_{t+1}, \epsilon_\phi) + \delta_t(\mathbf{x}_{t+1})$ , we can consider minimizing the following adversarial training objective of CM, if  $\|\delta_t(\mathbf{x}_{t+1})\| \leq \eta$  uniformly over  $t$ , for some constant  $\eta$ ,

<sup>3</sup>In practice, (15) is updated under target model  $f_{\theta^-}(\Phi_t(\mathbf{x}_{t+1}), t)$  with exponential moving average (EMA)  $\theta^-$  under a stop gradient operation. (Song et al., 2023) find that it greatly stabilizes the training process. In this section, we focus on the theory of consistency model and still use  $\theta$  in formulas.

<sup>4</sup>Here  $W_1(f_\theta(\mathbf{x}_t, t), \mathbf{x}_0)$  is the Wasserstein 1-distance between distributions of  $f_\theta(\mathbf{x}_t, t)$  and  $\mathbf{x}_0$ .

**Algorithm 1** Adversarial Training for Diffusion Model

---

```

1: Input: dataset  $\mathcal{D}$ , model parameter  $\theta$ , learning rate  $\kappa$ , loss weighting  $\lambda(\cdot)$ , adversarial steps  $K$ ,
   adversarial learning rate  $\alpha$ 
2: while do not converge do
3:   Sample  $\mathbf{x} \sim \mathcal{D}$  and  $t \sim \mathcal{U}[1, T]$ 
4:   Sample  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:    $\delta \leftarrow \mathbf{0}$ 
6:   for  $i = 1, 2, \dots, K$  do
7:      $\mathcal{L} \leftarrow \left\| \epsilon_{\theta}(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon + \delta) - \epsilon - \frac{\delta}{\sqrt{1 - \bar{\alpha}_t}} \right\|^2$  in (14)
8:      $\delta \leftarrow \delta + \alpha \cdot \frac{\nabla_{\delta} \mathcal{L}}{\|\nabla_{\delta} \mathcal{L}\|}$   $\triangleright$  maximize perturbation
9:      $\theta \leftarrow \theta - \kappa \cdot \nabla_{\theta} \mathcal{L}$   $\triangleright$  update model
10:  end for
11: end while

```

---

**Algorithm 2** Adversarial Training for Consistency Distillation

---

```

1: Input: dataset  $\mathcal{D}$ , initial model parameter  $\theta$ , learning rate  $\kappa$ , pretrained noise prediction model
    $\epsilon_{\phi}$ , ODE solver  $\hat{\Phi}(\cdot, \epsilon_{\phi}, \text{metric } d(\cdot, \cdot))$ , loss weighting  $\lambda(\cdot)$ , target model EMA  $\mu$ , adversarial
   steps  $K$ , adversarial learning rate  $\alpha$ 
2:  $\theta^- \leftarrow \theta$ 
3: while do not converge do
4:   Sample  $\mathbf{x} \sim \mathcal{D}$  and  $t \sim \mathcal{U}[0, T - 1]$ 
5:   Sample  $\mathbf{x}_{t+1}$  from (1)
6:    $\delta \leftarrow \mathbf{0}$ 
7:   for  $i = 1, 2, \dots, K$  do
8:      $\mathcal{L} \leftarrow \lambda(t)d(f_{\theta}(\mathbf{x}_{t+1}, t + 1), f_{\theta^-}(\hat{\Phi}_t(\mathbf{x}_{t+1}, \epsilon_{\phi}) + \delta, t))$  in (17)
9:      $\delta \leftarrow \delta + \alpha \cdot \frac{\nabla_{\delta} \mathcal{L}}{\|\nabla_{\delta} \mathcal{L}\|}$   $\triangleright$  maximize perturbation
10:     $\theta \leftarrow \theta - \kappa \cdot \nabla_{\theta} \mathcal{L}$   $\triangleright$  update model
11:     $\theta^- \leftarrow \text{stopgrad}(\mu\theta^- + (1 - \mu)\theta)$ 
12:  end for
13: end while

```

---

so that the target  $\Phi_t(\mathbf{x}_{t+1})$  is included in the maximal range as well.

$$\hat{\mathcal{L}}_{CD}^{Adv}(\theta) = \sum_{t=0}^{T-1} \mathbb{E}_{\mathbf{x}_{t+1}} \left[ \sup_{\|\delta\| \leq \eta} d\left(f_{\theta}(\hat{\Phi}_t(\mathbf{x}_{t+1}, \epsilon_{\phi}) + \delta, t), f_{\theta}(\mathbf{x}_{t+1}, t + 1)\right) \right]. \quad (17)$$

By doing so, the learned model  $f_{\theta}$  can be robust to the perturbation brought by  $\delta_t(\mathbf{x}_{t+1})$ , so that results in a small  $\mathcal{L}_{CD}(\theta)$ , as well as the small  $W_1(f_{\theta}(\mathbf{x}_T, T), \mathbf{x}_0)$  as proved in Theorem 2. Next, we use the following theorem to show that  $\|\delta_t(\mathbf{x}_{t+1})\|$  is indeed small, and minimizing  $\hat{\mathcal{L}}_{CD}^{Adv}(\theta)$  results in  $f_{\theta}(\mathbf{x}_T, T)$  with distribution approximates  $\mathbf{x}_0$ .

**Theorem 3.** *Under proper regularity conditions, for  $0 \leq t < T$ , we have  $\mathbb{E}_{\mathbf{x}_{t+1}}[\|\delta_t(\mathbf{x}_{t+1})\|] \leq o(1)$ . On the other hand, it holds*

$$W_1(f_{\theta}(\mathbf{x}_T, T), \mathbf{x}_0) \leq \sqrt{T \hat{\mathcal{L}}_{CD}^{Adv}(\theta) + o(1)}. \quad (18)$$

The theorem is proved in Appendix B.1, and it indicates that using the proposed adversarial training objective (17) of CM indeed guarantees the learned CM transfers  $\mathbf{x}_T$  into data from  $q(\mathbf{x}_0)$ .

## 6 EXPERIMENTS

### 6.1 ALGORITHMS

In the standard adversarial training method like Projected Gradient Descent (PGD) (Madry et al., 2018), the perturbation  $\delta$  is constructed by implementing numbers (3-8) of gradient ascents to  $\delta$

before updating the model, which slows down the training process. To resolve this, we adopt an efficient implementation (Shafahi et al., 2019) in Algorithms 1, 2 to solve AT (14) and (17) of DPM and CM, *which has similar computational cost compared to standard training*, and significantly accelerate standard AT. Notably, unlike PGD, in Algorithms 1 and 2, every maximization step of perturbation  $\delta$  follows an update step of the model  $\theta$ . Thus, the efficient AT do not require further back propagations to construct adversarial samples as in PGD. We provide a comparison between our efficient AT and standard AT (PGD) with the same update iterations of model  $\theta$  in Appendix G.1. Moreover, we observe that efficient AT can yield comparable and even better performance than PGD while accelerating the training ( $2.6\times$  speed-up), further verifying the benefits of our efficient AT.<sup>5</sup>

## 6.2 PERFORMANCE ON DPM

**Settings.** The experiments are conducted on the unconditional generation on CIFAR-10  $32\times 32$  (Krizhevsky & Hinton, 2009) and the class-conditional generation on ImageNet  $64\times 64$  (Deng et al., 2009). Our model and training pipelines in adopted from ADM (Dhariwal & Nichol, 2021) paper, where ADM is a UNet-type network (Ronneberger et al., 2015), with strong performance in image generation under diffusion model.

To save training costs, our methods and baselines are fine-tuned from pretrained models, rather than training from scratch. By doing so, we can efficiently assess the performance of methods, which is more practical for general scenarios. We also explore training from scratch in Appendix G.2, which also verifies the effectiveness of our method in this regime. During training, we fine-tune the pretrained models (details are in Appendix E.1) with batch size 128 for 150K iterations under learning rate  $1e-4$  on CIFAR-10, and batch size 1024 for 50K iterations under learning rate of  $3e-4$  on ImageNet. For the hyperparameters of AT, we select the adversarial learning rate  $\alpha$  from  $\{0.05, 0.1, 0.5\}$  and the adversarial step  $K$  from  $\{3, 5\}$ . More details are in Appendix E.1.

We use the Frechet Inception Distance (FID) (Heusel et al., 2017) to evaluate image quality. Unless otherwise specified, 50K images are sampled for evaluation. Other results of metric Classification Accuracy Score (CAS) (Ravuri & Vinyals, 2019), sFID, Inception Score, Precision, and Recall are in Appendix F.1 and F.4 for comprehensive evaluation.

**Baselines.** For experiments on diffusion models, we consider the following baselines. 1): the original pretrained model. Compared with it, we verify whether the models are overfitting during fine-tuning. 2): continue fine-tuning the pretrained model, which is fine-tuned with the standard diffusion objective (5). Compared to it, we validate whether performance improvements come only from more training costs. We also compare with the existing typical method to alleviate the DPM distribution mismatch, 3): ADM-IP (Ning et al., 2023), which adds a Gaussian perturbation to the input data to simulate mismatch errors during the training process. The last two fine-tuning baselines are based on **the same** pretrained model and hyperparameters as in the original literature.

**Results.** To verify the effectiveness of our AT method, we conduct experiments with four diffusion samplers: IDDPM (Dhariwal & Nichol, 2021), DDIM (Song et al., 2022), DPM-Solver (Lu et al., 2022b), and ES (Ning et al., 2024) under various NFEs. The sampler choices contain the three most popular samplers: IDDPM, DDIM, DPM-Solver, and ES, a sampler that scales down the norm of predicted noise to mitigate the distribution mismatch from the perspective of sampling. The experimental results of CIFAR-10 and ImageNet are shown in Table 1 and Table 2, respectively. Results of more than hundreds of NFEs are shown in Appendix F.3

As can be seen, the proposed AT for DPM significantly improves the performance of the original pretrained model and outperforms the other baselines (continue fine-tuning and ADM-IP) overall for all diffusion samplers and NFEs we take. Moreover, we have the following observations.

1): Fewer (practically used) sampling steps (5,10) will result in larger mismatching errors, while our AT method demonstrates significant improvements in this regime across various samplers, e.g., AT improves FID 27.72 to 17.36 under 5 NFEs DPM-Solver on ImageNet. This suggests that our method is indeed effective in alleviating the distribution mismatch of DPM. The results also indicate that our method consistently beats the baseline methods, regardless of stochastic (IDDPM)

<sup>5</sup>For the experts in AT, they would recognize that the AT in Algorithms 1, 2 actually constructs the adversarial augmented data to improve the performance of the model (Zhu et al., 2020; Jiang et al., 2020; Yi et al., 2021).



Table 1: Sample quality measured by FID  $\downarrow$  of different sampling methods of DPM under different NFEs on CIFAR10 32x32. All models are trained with same iterations (computational costs).

(a) IDDPM						(b) DDIM					
Methods \ NFEs	5	8	10	20	50	Methods \ NFEs	5	8	10	20	50
ADM (original)	37.99	26.75	22.62	10.52	4.55	ADM (original)	34.28	14.34	11.66	7.00	4.68
ADM (finetune)	<b>36.91</b>	26.06	21.94	10.58	4.34	ADM (finetune)	29.30	15.08	12.06	6.80	4.15
ADM-IP	47.57	26.91	20.09	7.81	3.42	ADM-IP	43.15	15.72	10.47	4.58	4.89
ADM-AT (Ours)	37.15	<b>23.59</b>	<b>15.88</b>	<b>6.60</b>	<b>3.34</b>	ADM-AT (Ours)	<b>26.38</b>	<b>12.98</b>	<b>9.30</b>	<b>4.40</b>	<b>3.07</b>

(c) ES						(d) DPM-Solver					
Methods \ NFEs	5	8	10	20	50	Methods \ NFEs	5	8	10	20	50
ADM (original)	82.18	29.28	17.73	5.11	2.70	ADM (original)	23.95	8.00	5.46	3.46	3.14
ADM (finetune)	63.46	24.80	17.03	5.19	2.52	ADM (finetune)	22.98	7.61	5.29	3.41	3.12
ADM-IP	91.10	31.44	18.72	5.19	2.89	ADM-IP	43.83	6.70	6.80	9.78	10.91
ADM-AT (Ours)	<b>41.07</b>	<b>21.62</b>	<b>14.68</b>	<b>4.36</b>	<b>2.48</b>	ADM-AT (Ours)	<b>18.40</b>	<b>5.84</b>	<b>4.81</b>	<b>3.28</b>	<b>3.01</b>

Table 2: Sample quality measured by FID  $\downarrow$  of different sampling methods of DPM under different NFEs on ImageNet 64x64. All models are trained with the same iterations (computational costs).

(a) IDDPM						(b) DDIM					
Methods \ NFEs	5	8	10	20	50	Methods \ NFEs	5	8	10	20	50
ADM (original)	76.92	33.74	27.63	12.85	5.30	ADM (original)	60.07	20.10	14.97	8.41	5.65
ADM (finetune)	78.87	33.99	27.82	12.80	5.26	ADM (finetune)	60.32	20.26	15.04	8.32	5.48
ADM-IP	67.12	29.96	22.60	8.66	<b>3.83</b>	ADM-IP	76.51	26.25	18.05	8.40	6.94
ADM-AT (Ours)	<b>45.65</b>	<b>23.79</b>	<b>19.18</b>	<b>8.28</b>	4.01	ADM-AT (Ours)	<b>43.04</b>	<b>16.08</b>	<b>12.15</b>	<b>6.20</b>	<b>4.67</b>

(c) ES						(d) DPM-Solver					
Methods \ NFEs	5	8	10	20	50	Methods \ NFEs	5	8	10	20	50
ADM (original)	71.31	28.97	21.10	8.23	3.76	ADM (original)	27.72	10.06	7.21	4.69	4.24
ADM (finetune)	72.30	29.24	21.58	8.25	3.64	ADM (finetune)	27.82	9.97	7.22	4.64	<b>4.15</b>
ADM-IP	88.37	33.91	23.32	7.80	3.54	ADM-IP	32.43	9.94	8.87	9.16	9.68
ADM-AT (Ours)	<b>43.95</b>	<b>19.57</b>	<b>14.12</b>	<b>6.16</b>	<b>3.45</b>	ADM-AT (Ours)	<b>17.36</b>	<b>6.55</b>	<b>5.78</b>	<b>4.56</b>	4.34

or deterministic samplers (DDIM, DPM-Solver). 2): The ES sampler results show that our AT is orthogonal to the sampling-based method to mitigate the distribution mismatch problem and can be combined to further alleviate the issue. Notably, we further verify in Appendix G.2 that our methods will not slow the convergence unlike AT in classification (Madry et al., 2018). We also perform ablation analysis of hyperparameters in our AT framework in Appendix G.3.

### 6.3 PERFORMANCE ON LATENT CONSISTENCY MODELS

**Settings.** We further evaluate the proposed AT for consistency models on text-to-image generation tasks with Latent Consistency Models (Luo et al., 2023) Stable Diffusion (SD) v1.5 (Rombach et al., 2022) backbone, which generates  $512 \times 512$  images. Both our AT and the original LCM training (baseline) are trained from scratch with the same hyperparameters (the IP method (Ning et al., 2023) is not applied straightforwardly). The training set is LAION-Aesthetics-6.5+ (Schuhmann et al., 2022) with hyperparameters following Song et al. (2023); Luo et al. (2023). We select the adversarial learning rate  $\alpha$  from  $\{0.02, 0.05\}$  and adversarial step  $K$  from  $\{2, 3\}$ . The models are trained with a batch size of 64 for 100K iterations. More details are shown in Appendix E.2.

Following Luo et al. (2023) and Chen et al. (2024), we evaluate models on MS-COCO 2014 (Lin et al., 2014a) at a resolution of  $512 \times 512$  by randomly drawing 30K prompts from its validation set. Then, we report the FID between the generated samples under these prompts and the reference samples from the full validation set following Saharia et al. (2022). We also report CLIP scores (Hessel et al., 2021) to evaluate the text-image alignment by CLIP-ViT-B/16.

Table 3: Results of LCM on MS-COCO 2014 validation set at 512×512 resolution in terms of FID ↓ and CLIP score ↑. All models are trained with the same setting (computational costs).

Methods	FID ↓				CLIP Score ↑			
	1 step	2 step	4 step	8 step	1 step	2 step	4 step	8 step
LCM	25.43	12.61	11.61	12.62	29.25	30.24	30.40	30.47
LCM-AT (Ours)	<b>23.34</b>	<b>11.28</b>	<b>10.31</b>	<b>10.68</b>	<b>29.63</b>	<b>30.43</b>	<b>30.49</b>	<b>30.53</b>

**Results.** The methods are evaluated under various sampling steps in Table 3, which shows that the LCM with AT consistently improves FID under various sampling steps. Besides, though the AT is not specified to improve text-image alignment, we observe that it has comparable or even better CLIP scores across various sampling steps, which shows that AT will not degenerate text-image alignment.

## 7 CONCLUSION

In this paper, we novelly introduce efficient Adversarial Training (AT) in the training of DPM and CM to mitigate the issue of distribution mismatch between training and sampling. We conduct an in-depth analysis of the DPM training objective and systematically characterize the distribution mismatch problem. Furthermore, we prove that the training objective of CM similarly faces the distribution mismatch issue. We theoretically prove that DRO can mitigate the mismatch for both DPM and CM, which is equivalent to conducting AT. Experiments on image generation and text-to-image generation benchmarks verify the effectiveness of the proposed AT method in alleviating the distribution mismatch of DPM and CM.

### ACKNOWLEDGMENTS

We thank anonymous reviewers for insightful feedback that helped improve the paper. Zekun Wang, Ming Liu, Bing Qin are supported by the National Science Foundation of China (U22B2059, 62276083), the Human-Machine Integrated Consultation System for Cardiovascular Diseases (2023A003). They also appreciate the support from China Mobile Group Heilongjiang Co., Ltd.

### REFERENCES

- Fan Bao, Chongxuan Li, Jun Zhu, and Bo Zhang. Analytic-dpm: an analytic estimate of the optimal reverse variance in diffusion probabilistic models. In *International Conference on Learning Representations*, 2022.
- Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. Scheduled sampling for sequence prediction with recurrent neural networks. In *Advances in Neural Information Processing Systems*, 2015.
- Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Conference on Computer Vision and Pattern Recognition*, 2023.
- Hongrui Chen, Holden Lee, and Jianfeng Lu. Improved analysis of score-based generative modeling: User-friendly bounds under minimal smoothness assumptions. In *International Conference on Machine Learning*, 2023.
- Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Zhongdao Wang, James T. Kwok, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart- $\alpha$ : Fast training of diffusion transformer for photorealistic text-to-image synthesis. In *International Conference on Learning Representations*, 2024.
- 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 *International Conference on Learning Representations*, 2022.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Conference on Computer Vision and Pattern Recognition*, 2009.

- Prafulla Dhariwal and Alexander Quinn Nichol. Diffusion models beat GANs on image synthesis. In *Advances in Neural Information Processing Systems*, volume 34, pp. 8780–8794, 2021.
- John Duchi. Lecture notes for statistics 311/electrical engineering 377. URL: [https://stanford.edu/class/stats311/Lectures/full\\_notes.pdf](https://stanford.edu/class/stats311/Lectures/full_notes.pdf). Last visited on, 2:23, 2016.
- Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, 2014.
- Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*, 2015.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Conference on Computer Vision and Pattern Recognition*, 2016.
- Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2021.
- 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. 2017.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, 2020.
- Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *Journal of Machine Learning Research*, 23(47):1–33, 2022a.
- Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. In *Advances in Neural Information Processing Systems*, 2022b.
- Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Tuo Zhao. SMART: robust and efficient fine-tuning for pre-trained natural language models through principled regularized optimization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020.
- Diederik P Kingma and Max Welling. Auto-encoding variational {Bayes}. In *International Conference on Learning Representations*, 2013.
- Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. 2009.
- Jaeho Lee and Maxim Raginsky. Minimax statistical learning with wasserstein distances. In *Advances in Neural Information Processing Systems*, 2018.
- Gen Li, Yuting Wei, Yuxin Chen, and Yuejie Chi. Towards faster non-asymptotic convergence for diffusion-based generative models. Preprint arXiv:2306.09251, 2023.
- Mingxiao Li, Tingyu Qu, Wei Sun, and Marie-Francine Moens. Alleviating exposure bias in diffusion models through sampling with shifted time steps. In *International Conference on Learning Representations*, 2024.
- Yangming Li and Mihaela van der Schaar. On error propagation of diffusion models. In *The Twelfth International Conference on Learning Representations*, 2024.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, 2014a.
- Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. In *ECCV*, 2014b.

- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019.
- Aaron Lou and Stefano Ermon. Reflected diffusion models. In *International Conference on Machine Learning*, 2023.
- 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 *Advances in Neural Information Processing Systems*, 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 *Advances in Neural Information Processing Systems*, 2022b.
- Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models: Synthesizing high-resolution images with few-step inference, 2023.
- Jiajun Ma, Shuchen Xue, Tianyang Hu, Wenjia Wang, Zhaoqiang Liu, Zhenguo Li, Zhi-Ming Ma, and Kenji Kawaguchi. The surprising effectiveness of skip-tuning in diffusion sampling. *arXiv preprint arXiv:2402.15170*, 2024.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018.
- Hongseok Namkoong. *Reliable machine learning via distributional robustness*. PhD thesis, Stanford University, 2019.
- Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. GLIDE: towards photorealistic image generation and editing with text-guided diffusion models. In *International Conference on Machine Learning*.
- Weili Nie, Brandon Guo, Yujia Huang, Chaowei Xiao, Arash Vahdat, and Animashree Anandkumar. Diffusion models for adversarial purification. In *International Conference on Machine Learning*, 2022.
- Mang Ning, Enver Sangineto, Angelo Porrello, Simone Calderara, and Rita Cucchiara. Input perturbation reduces exposure bias in diffusion models. In *International Conference on Machine Learning*, 2023.
- Mang Ning, Mingxiao Li, Jianlin Su, Albert Ali Salah, and Itir Önal Ertugrul. Elucidating the exposure bias in diffusion models. In *International Conference on Learning Representations*, 2024.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents, 2022.
- Marc’Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. Sequence level training with recurrent neural networks. In *International Conference on Learning Representations*, 2016.
- Suman V. Ravuri and Oriol Vinyals. Classification accuracy score for conditional generative models. In *Advances in Neural Information Processing Systems*, 2019.
- Sylvestre-Alvise Rebuffi, Sven Gowal, Dan A. Calian, Florian Stimberg, Olivia Wiles, and Timothy A. Mann. Fixing data augmentation to improve adversarial robustness. Preprint arXiv:2103.01946, 2021.
- Zhiyao Ren, Yibing Zhan, Liang Ding, Gaoang Wang, Chaoyue Wang, Zhongyi Fan, and Dacheng Tao. Multi-step denoising scheduled sampling: Towards alleviating exposure bias for diffusion models. In *AAAI Conference on Artificial Intelligence*, 2024.

- Steven J. Rennie, Etienne Marcheret, Youssef Mroueh, Jerret Ross, and Vaibhava Goel. Self-critical sequence training for image captioning. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Conference on Computer Vision and Pattern Recognition*, 2022.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18*, pp. 234–241. Springer, 2015.
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar 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 *Advances in Neural Information Processing Systems*, 2022.
- Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. In *International Conference on Learning Representations*, 2022.
- Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. LAION-5B: an open large-scale dataset for training next generation image-text models. In *Advances in Neural Information Processing Systems*, 2022.
- Ali Shafahi, Mahyar Najibi, Amin Ghiasi, Zheng Xu, John P. Dickerson, Christoph Studer, Larry S. Davis, Gavin Taylor, and Tom Goldstein. Adversarial training for free! In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, 2019.
- Alexander Shapiro. Distributionally robust stochastic programming. *SIAM Journal on Optimization*, 27(4):2258–2275, 2017.
- Shiqi Shen, Yong Cheng, Zhongjun He, Wei He, Hua Wu, Maosong Sun, and Yang Liu. Minimum risk training for neural machine translation. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, 2016.
- Albert N Shiryaev. *Probability-1*, volume 95. Springer, 2016.
- Aman Sinha, Hongseok Namkoong, and John Duchi. Certifying some distributional robustness with principled adversarial training. In *International Conference on Learning Representations*, 2018.
- Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, 2015.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2022.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2020.
- Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of score-based diffusion models. 2021.
- Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. In *International Conference on Machine Learning*, 2023.
- Cédric Villani et al. *Optimal transport: old and new*, volume 338. Springer, 2009.

- Martin J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge university press, 2019.
- Ruoyu Wang, Mingyang Yi, Zhitang Chen, and Shengyu Zhu. Out-of-distribution generalization with causal invariant transformations. In *Conference on Computer Vision and Pattern Recognition*, 2022.
- Zekai Wang, Tianyu Pang, Chao Du, Min Lin, Weiwei Liu, and Shuicheng Yan. Better diffusion models further improve adversarial training. In *International Conference on Machine Learning*, 2023.
- Dongxian Wu, Shu-Tao Xia, and Yisen Wang. Adversarial weight perturbation helps robust generalization. In *Advances in Neural Information Processing Systems*, 2020.
- Shuchen Xue, Zhaoqiang Liu, Fei Chen, Shifeng Zhang, Tianyang Hu, Enze Xie, and Zhenguo Li. Accelerating diffusion sampling with optimized time steps. *arXiv preprint arXiv:2402.17376*, 2024a.
- Shuchen Xue, Mingyang Yi, Weijian Luo, Shifeng Zhang, Jiacheng Sun, Zhenguo Li, and Zhi-Ming Ma. Sa-solver: Stochastic adams solver for fast sampling of diffusion models. *Advances in Neural Information Processing Systems*, 36, 2024b.
- Mingyang Yi, Lu Hou, Jiacheng Sun, Lifeng Shang, Xin Jiang, Qun Liu, and Zhiming Ma. Improved ood generalization via adversarial training and pretraing. In *International Conference on Machine Learning*, 2021.
- Mingyang Yi, Jiacheng Sun, and Zhenguo Li. On the generalization of diffusion model. Preprint arXiv:2305.14712, 2023a.
- Mingyang Yi, Ruoyu Wang, Jiacheng Sun, Zhenguo Li, and Zhi-Ming Ma. Breaking correlation shift via conditional invariant regularizer. In *The International Conference on Learning Representations*, 2023b.
- Mingyang Yi, Aoxue Li, Yi Xin, and Zhenguo Li. Towards understanding the working mechanism of text-to-image diffusion model. Preprint arXiv:2405.15330, 2024.
- Tianwei Yin, Michaël Gharbi, Richard Zhang, Eli Shechtman, Frédo Durand, William T Freeman, and Taesung Park. One-step diffusion with distribution matching distillation. In *CVPR*, 2024.
- Boya Zhang, Weijian Luo, and Zhihua Zhang. Enhancing adversarial robustness via score-based optimization. In *Advances in Neural Information Processing Systems*, 2023.
- Dinghui Zhang, Tianyuan Zhang, Yiping Lu, Zhanxing Zhu, and Bin Dong. You only propagate once: Accelerating adversarial training via maximal principle. In *Advances in Neural Information Processing Systems*, 2019a.
- Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric P. Xing, Laurent El Ghaoui, and Michael I. Jordan. Theoretically principled trade-off between robustness and accuracy. In *International Conference on Machine Learning*, 2019b.
- Wen Zhang, Yang Feng, Fandong Meng, Di You, and Qun Liu. Bridging the gap between training and inference for neural machine translation. In *Proceedings of the 57th Conference of the Association for Computational Linguistics*, 2019c.
- Chen Zhu, Yu Cheng, Zhe Gan, Siqi Sun, Tom Goldstein, and Jingjing Liu. Freelib: Enhanced adversarial training for natural language understanding. In *International Conference on Learning Representations*, 2020.

## A PROOFS IN SECTION 4

In this section, we present the proofs of the results in Section 4.

### A.1 PROOFS IN SECTION 4.2

**Proposition 1.** *The minimization problem (4) is equivalent to minimizing an upper bound of  $\mathbb{E}_q[-\log p_\theta(\mathbf{x}_t)]$  for any  $0 \leq t \leq T$ .*

*Proof.* We prove the first equivalence, by Jensen’s inequality. For any  $0 \leq t < T$ , we have

$$\begin{aligned}
& -\mathbb{E}_q[\log p_\theta(\mathbf{x}_t)] \\
& \leq \mathbb{E}_q \left[ -\log \frac{p_\theta(\mathbf{x}_{t:T})}{q(\mathbf{x}_{t+1:T} | \mathbf{x}_t)} \right] \\
& = \mathbb{E}_q \left[ -\log p_\theta(\mathbf{x}_T) - \sum_{t \leq s < T} \log \frac{p_\theta(\mathbf{x}_s | \mathbf{x}_{s+1})}{q(\mathbf{x}_{s+1} | \mathbf{x}_s)} \right] \\
& = \mathbb{E}_q \left[ -\log p_\theta(\mathbf{x}_T) - \sum_{t \leq s < T} \log \frac{p_\theta(\mathbf{x}_s | \mathbf{x}_{s+1})}{q(\mathbf{x}_s | \mathbf{x}_{s+1})} \cdot \frac{q(\mathbf{x}_s)}{q(\mathbf{x}_{s+1})} \right] \tag{19} \\
& = \mathbb{E}_q \left[ -\log \frac{p_\theta(\mathbf{x}_T)}{q(\mathbf{x}_T)} - \sum_{t \leq s < T} \log \frac{p_\theta(\mathbf{x}_s | \mathbf{x}_{s+1})}{q(\mathbf{x}_s | \mathbf{x}_{s+1})} - \log q(\mathbf{x}_t) \right] \\
& = D_{KL}(q(\mathbf{x}_T) \| p_\theta(\mathbf{x}_T)) + \mathbb{E}_q \left[ \sum_{s=t}^{T-1} \underbrace{D_{KL}(q(\mathbf{x}_s | \mathbf{x}_{s+1}) \| p_\theta(\mathbf{x}_s | \mathbf{x}_{s+1}))}_{L_t} \right] + H(\mathbf{x}_t)
\end{aligned}$$

Taking  $t = 0$ , we prove the first equivalence. Besides that, the entropy  $H(\mathbf{x}_t)$  of  $\mathbf{x}_t$  is a constant for  $\theta$  given data distribution  $\mathbf{x}_0$  for any  $0 \leq t < T$ . The second conclusion holds due to the non-negative property of KL-divergence.  $\square$

**Proposition 2.** *Suppose  $p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1})$  matches  $q(\mathbf{x}_t | \mathbf{x}_{t+1})$  well such that*

$$L_t = D_{KL}(q(\mathbf{x}_t | \mathbf{x}_{t+1}) \| p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1})) \leq \frac{\gamma}{T}, \tag{8}$$

*and the discrepancy satisfies  $D_{KL}(q(\mathbf{x}_T) \| p_\theta(\mathbf{x}_T)) \leq \gamma_0$ , then for any  $0 \leq t \leq T$ , we have*

$$D_{KL}(q(\mathbf{x}_t) \| p_\theta(\mathbf{x}_t)) \leq D_{KL}(q(\mathbf{x}_T) \| p_\theta(\mathbf{x}_T)) + L_t \leq \gamma_0 + \frac{(T-t)\gamma}{T}. \tag{9}$$

*Proof.* We have the following decomposition due to the chain rule of KL-divergence

$$\begin{aligned}
D_{KL}(q(\mathbf{x}_t, \mathbf{x}_{t+1}) \| p_\theta(\mathbf{x}_t, \mathbf{x}_{t+1})) &= D_{KL}(q(\mathbf{x}_t | \mathbf{x}_{t+1}) \| p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1})) + D_{KL}(q(\mathbf{x}_{t+1}) \| p_\theta(\mathbf{x}_{t+1})) \\
&= D_{KL}(q(\mathbf{x}_{t+1} | \mathbf{x}_t) \| p_\theta(\mathbf{x}_{t+1} | \mathbf{x}_t)) + D_{KL}(q(\mathbf{x}_t) \| p_\theta(\mathbf{x}_t)), \tag{20}
\end{aligned}$$

The transition probability  $p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1})$  matches  $q(\mathbf{x}_t | \mathbf{x}_{t+1})$ , so that the above equality implies

$$\begin{aligned}
& D_{KL}(q(\mathbf{x}_t) \| p_\theta(\mathbf{x}_t)) \\
& = D_{KL}(q(\mathbf{x}_{t+1}) \| p_\theta(\mathbf{x}_{t+1})) + D_{KL}(q(\mathbf{x}_t | \mathbf{x}_{t+1}) \| p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1})) - D_{KL}(q(\mathbf{x}_{t+1} | \mathbf{x}_t) \| p_\theta(\mathbf{x}_{t+1} | \mathbf{x}_t)) \\
& \leq D_{KL}(q(\mathbf{x}_{t+1}) \| p_\theta(\mathbf{x}_{t+1})) + \frac{\gamma}{T}. \tag{21}
\end{aligned}$$

The proposition holds due to initial condition  $D_{KL}(q(\mathbf{x}_T) \| p_\theta(\mathbf{x}_T)) \leq \gamma_0$  and simple induction.  $\square$

**Proposition 3.**  *$L_t$  in (4) is well minimized, only if  $q(\mathbf{x}_{t+1})$  is Gaussian or  $\|\mathbf{x}_{t+1} - \mathbf{x}_t\| \rightarrow 0$ .*

*Proof.* Due to Bayes' rule, we have

$$\begin{aligned}
q(\mathbf{x}_t | \mathbf{x}_{t+1}) &= \frac{q(\mathbf{x}_{t+1} | \mathbf{x}_t)q(\mathbf{x}_t)}{q(\mathbf{x}_{t+1})} \\
&\propto \exp\left(-\frac{\|\mathbf{x}_{t+1} - \sqrt{\alpha_{t+1}}\mathbf{x}_t\|^2}{2(1-\alpha_{t+1})} + \log q(\mathbf{x}_t) - \log q(\mathbf{x}_{t+1})\right) \\
&\propto \exp\left(-\frac{\|\mathbf{x}_{t+1} - \sqrt{\alpha_{t+1}}\mathbf{x}_t\|^2}{2(1-\alpha_{t+1})} + \langle \nabla_{\mathbf{x}} \log q(\mathbf{x}_{t+1}), \mathbf{x}_t - \mathbf{x}_{t+1} \rangle\right) \\
&\quad \exp\left(\frac{1}{2}(\mathbf{x}_t - \mathbf{x}_{t+1})^\top \nabla_{\mathbf{x}}^2 \log q(\mathbf{x}_{t+1})(\mathbf{x}_t - \mathbf{x}_{t+1}) + O(\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^3)\right).
\end{aligned} \tag{22}$$

As can be seen, the conditional probability can be approximated by Gaussian only if  $\nabla_{\mathbf{x}}^3 \log q(\mathbf{x}_{t+1})$  is zero or  $\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^3$  is extremely small with high probability. The two conditions can be respectively satisfied when  $q(\mathbf{x}_t)$  is a Gaussian or  $\mathbf{x}_t$  close to  $\mathbf{x}_{t+1}$ .  $\square$

## A.2 PROOFS IN SECTION 4.2

**Proposition 4.** For any distribution  $\tilde{q}$  satisfies  $\tilde{q}(\mathbf{x}_t) = q(\mathbf{x}_t)$  for specific  $t$ , we have

$$\mathbb{E}_q[-\log p_{\theta}(\mathbf{x}_t)] \leq \underbrace{D_{KL}(\tilde{q}(\mathbf{x}_t | \mathbf{x}_{t+1}) \| p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t+1}))}_{L_t^{\tilde{q}}} + C, \tag{11}$$

for a constant  $C$  independent of  $\theta$ .

*Proof.* W.o.l.g., suppose  $p_{\theta}(\mathbf{x}_t, \mathbf{x}_{t+1}) = p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t+1})q(\mathbf{x}_{t+1})$  and  $\tilde{q}(\mathbf{x}_t, \mathbf{x}_{t+1}) = \tilde{q}(\mathbf{x}_{t+1} | \mathbf{x}_t)q(\mathbf{x}_t)$ . By Jensen's inequality, we have

$$\begin{aligned}
&\mathbb{E}_q[-\log p_{\theta}(\mathbf{x}_t)] \\
&= -\int q(\mathbf{x}_t) \left( \log \int p_{\theta}(\mathbf{x}_t, \mathbf{x}_{t+1}) d\mathbf{x}_{t+1} \right) d\mathbf{x}_t \\
&= -\int q(\mathbf{x}_t) \left( \log \int \frac{p_{\theta}(\mathbf{x}_t, \mathbf{x}_{t+1})}{\tilde{q}(\mathbf{x}_{t+1} | \mathbf{x}_t)} \tilde{q}(\mathbf{x}_{t+1} | \mathbf{x}_t) d\mathbf{x}_{t+1} \right) d\mathbf{x}_t \\
&\leq -\int q(\mathbf{x}_t) \left( \int \log \frac{p_{\theta}(\mathbf{x}_t, \mathbf{x}_{t+1})}{\tilde{q}(\mathbf{x}_{t+1} | \mathbf{x}_t)} \tilde{q}(\mathbf{x}_{t+1} | \mathbf{x}_t) d\mathbf{x}_{t+1} \right) d\mathbf{x}_t \\
&= -\int q(\mathbf{x}_t) \left( \int \tilde{q}(\mathbf{x}_{t+1} | \mathbf{x}_t) \log \frac{p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t+1})}{\tilde{q}(\mathbf{x}_{t+1} | \mathbf{x}_t)} d\mathbf{x}_{t+1} \right) d\mathbf{x}_t \\
&\quad - \int q(\mathbf{x}_t) \left( \int \tilde{q}(\mathbf{x}_{t+1} | \mathbf{x}_t) \log \frac{q(\mathbf{x}_{t+1})}{\tilde{q}(\mathbf{x}_{t+1} | \mathbf{x}_t)} d\mathbf{x}_{t+1} \right) d\mathbf{x}_t \\
&= -\int \tilde{q}(\mathbf{x}_t, \mathbf{x}_{t+1}) \log \frac{p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t+1})}{\tilde{q}(\mathbf{x}_{t+1} | \mathbf{x}_t)} d\mathbf{x}_t d\mathbf{x}_{t+1} + C_1 \\
&= -\int \tilde{q}(\mathbf{x}_t, \mathbf{x}_{t+1}) \log \frac{p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t+1})}{\tilde{q}(\mathbf{x}_t | \mathbf{x}_{t+1})} \cdot \frac{q(\mathbf{x}_t)}{\tilde{q}(\mathbf{x}_{t+1})} d\mathbf{x}_t d\mathbf{x}_{t+1} + C_1 \\
&= -\int \tilde{q}(\mathbf{x}_t, \mathbf{x}_{t+1}) \log \frac{p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t+1})}{\tilde{q}(\mathbf{x}_t | \mathbf{x}_{t+1})} d\mathbf{x}_t d\mathbf{x}_{t+1} + C_1 + C_2 \\
&= D_{KL}(\tilde{q}(\mathbf{x}_t | \mathbf{x}_{t+1}) \| p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t+1})) + C \\
&= L_{\text{vib}}^{\tilde{q}}(\theta, t) + C,
\end{aligned} \tag{23}$$

where  $C, C_1, C_2$  are all constants independent of  $\theta$ .  $\square$

### A.2.1 PROOF OF THEOREM 1

In this section, we prove the Theorem 1. To simplify the notation, let  $p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t+1}) \sim \mathcal{N}(\boldsymbol{\mu}_{\theta}(\mathbf{x}_{t+1}, t+1), \sigma_{t+1})$ <sup>6</sup> in (6), then the optimal solution (Lemma 9 in (Bao et al., 2022)) of minimizing  $L_{t+1}^{\tilde{q}_t}$  is

$$\boldsymbol{\mu}_{\theta}(\mathbf{x}_{t+1}, t+1) = \mathbb{E}_{\tilde{q}_t}[\mathbf{x}_t | \mathbf{x}_{t+1}]. \tag{24}$$

<sup>6</sup>Here  $\sigma_{t+1}$  can be also optimized as in (Bao et al., 2022), but we find optimizing it in practice does not improve the empirical results.



For every specific  $t$ , we consider the following  $\tilde{q}_t$  in (12)<sup>7</sup>, such that

$$\begin{aligned}\tilde{q}_t(\mathbf{x}_{t+1} | \mathbf{x}_t) &\neq q(\mathbf{x}_{t+1} | \mathbf{x}_t); \\ \tilde{q}_t(\mathbf{x}_{t+1}) &\neq q(\mathbf{x}_{t+1}); \\ \tilde{q}_t(\mathbf{x}_{0:t}) &= q(\mathbf{x}_{0:t}). \\ \tilde{q}_t(\mathbf{x}_t | \mathbf{x}_0, \mathbf{x}_{t+1}) &= q(\mathbf{x}_t | \mathbf{x}_0, \mathbf{x}_{t+1}) = \mathcal{N}(\mu_{t+1}(\mathbf{x}_0, \mathbf{x}_{t+1}), \sigma_t).\end{aligned}\tag{25}$$

where  $\mu_{t+1}(\mathbf{x}_0, \mathbf{x}_{t+1}) = \frac{\sqrt{\bar{\alpha}_t(1-\alpha_{t+1})}}{1-\bar{\alpha}_{t+1}}\mathbf{x}_0 + \frac{\sqrt{\alpha_{t+1}(1-\bar{\alpha}_t)}}{1-\bar{\alpha}_{t+1}}\mathbf{x}_{t+1}$ . The  $\tilde{q}_t$  can be taken due to the Bayesian rule. Next, we analyze the optimal formulation in (24). Due to the property of conditional expectation, we have

$$\boldsymbol{\mu}_\theta(\mathbf{x}_{t+1}, t+1) = \mathbb{E}_{\tilde{q}_t} [\mathbb{E}_{\tilde{q}_t} [\mathbf{x}_t | \mathbf{x}_0, \mathbf{x}_{t+1}] | \mathbf{x}_{t+1}] = \mu_{t+1}(\mathbb{E}_{\tilde{q}_t} [\mathbf{x}_0 | \mathbf{x}_{t+1}], \mathbf{x}_{t+1}).\tag{26}$$

As can be seen, the optimal transition rule is decided by the conditional expectation  $\mathbb{E}_{\tilde{q}_t} [\mathbf{x}_0 | \mathbf{x}_{t+1}]$  for some  $\tilde{q}_t(\mathbf{x}_{t+1}) \in B_{D_{KL}}(\tilde{q}(\mathbf{x}_{t+1}), \eta_0)$  in (12). Then, we have the following lemma to get the desired conditional expectation.

**Lemma 1.** *There exists some  $\eta \geq \eta_0$  in (27) which makes (27) equivalent to problem (12).*

$$\min_{\theta} \sum_{t=0}^{T-1} \mathbb{E}_{\tilde{q}_t(\mathbf{x}_0)} \sup_{\tilde{q}_t(\mathbf{x}_{t+1}|\mathbf{x}_0) \in B_{D_{KL}}(\tilde{q}_t(\mathbf{x}_{t+1}|\mathbf{x}_0), \eta)} \mathbb{E}_{\tilde{q}_t(\mathbf{x}_{t+1}|\mathbf{x}_0)} [\|\boldsymbol{\mu}_\theta(\mathbf{x}_{t+1}, t+1) - \mathbf{x}_0\|^2],\tag{27}$$

where  $\mathbb{E}_{p_\theta} [\mathbf{x}_0 | \mathbf{x}_{t+1}] = \boldsymbol{\mu}_\theta(\mathbf{x}_{t+1}, t+1)$ .

*Proof.* Let us check the training objective  $\min_{\theta} \sup_{\tilde{q}_t \in B_{D_{KL}}(\tilde{q}_t(\mathbf{x}_t | \mathbf{x}_{t+1}), \eta)} D_{KL}(\tilde{q}_t(\mathbf{x}_t | \mathbf{x}_{t+1}) \| p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1}))$ . During this proof, we abbreviate  $B_{D_{KL}}(\tilde{q}_t(\mathbf{x}_{t+1}), \eta)$  as  $B$ . Since  $p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1}) \sim \mathcal{N}(\boldsymbol{\mu}_\theta(\mathbf{x}_{t+1}, t+1), \sigma_{t+1})$ , then

$$\begin{aligned}&\sup_{\tilde{q}_t(\mathbf{x}_{t+1}) \in B} D_{KL}(\tilde{q}_t(\mathbf{x}_t | \mathbf{x}_{t+1}) \| p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1})) \\ &\propto -\frac{d}{2} \log 2\pi\sigma_{t+1}^2 - \frac{1}{2\sigma_{t+1}^2} \sup_{\tilde{q}_t(\mathbf{x}_{t+1}) \in B} \mathbb{E}_{\tilde{q}_t(\mathbf{x}_t, \mathbf{x}_{t+1})} [\|\mathbf{x}_t - \boldsymbol{\mu}_\theta(\mathbf{x}_{t+1}, t+1)\|^2].\end{aligned}\tag{28}$$

As we consider  $\sigma_{t+1}$  as constant, an analysis of the expectation term is enough. Due to

$$\begin{aligned}\mathbb{E}_{\tilde{q}_t(\mathbf{x}_t, \mathbf{x}_{t+1})} [\|\mathbf{x}_t - \boldsymbol{\mu}_\theta(\mathbf{x}_{t+1}, t+1)\|^2] &\geq \inf_f \mathbb{E}_{\tilde{q}_t(\mathbf{x}_0, \mathbf{x}_t, \mathbf{x}_{t+1})} [\|\mathbf{x}_t - f(\mathbf{x}_0, \mathbf{x}_{t+1})\|^2] \\ &= \mathbb{E}_{\tilde{q}_t(\mathbf{x}_0, \mathbf{x}_t, \mathbf{x}_{t+1})} [\|\mathbf{x}_t - \mathbb{E}_{\tilde{q}_t}[\mathbf{x}_t | \mathbf{x}_0, \mathbf{x}_{t+1}]\|^2],\end{aligned}\tag{29}$$

where the last term is invariant over  $\tilde{q}_t \in B$  so that it is a uniform lower bound over all possible  $\tilde{q}_t$  and  $p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1})$ . The above inequality indicates that the optimal  $\boldsymbol{\mu}_\theta(\mathbf{x}_{t+1}, t+1)$  is achieved when the left in (29) becomes the right in (29).

On the other hand, for any  $\tilde{q}_t \in B$ , let us compute the gap such that

$$\begin{aligned}&\mathbb{E}_{\tilde{q}_t(\mathbf{x}_t, \mathbf{x}_{t+1})} [\|\mathbf{x}_t - \boldsymbol{\mu}_\theta(\mathbf{x}_{t+1}, t+1)\|^2] \\ &= \mathbb{E}_{\tilde{q}_t} [\|\mathbf{x}_t - \mathbb{E}_{\tilde{q}_t}[\mathbf{x}_t | \mathbf{x}_0, \mathbf{x}_{t+1}] + \mathbb{E}_{\tilde{q}_t}[\mathbf{x}_t | \mathbf{x}_0, \mathbf{x}_{t+1}] - \boldsymbol{\mu}_\theta(\mathbf{x}_{t+1}, t+1)\|^2] \\ &= \mathbb{E}_{\tilde{q}_t} [\|\mathbf{x}_t - \mathbb{E}_{\tilde{q}_t}[\mathbf{x}_t | \mathbf{x}_0, \mathbf{x}_{t+1}]\|^2] \\ &+ \mathbb{E}_{\tilde{q}_t} [\|\boldsymbol{\mu}_\theta(\mathbf{x}_{t+1}, t+1) - \mathbb{E}_{\tilde{q}_t}[\mathbf{x}_t | \mathbf{x}_0, \mathbf{x}_{t+1}]\|^2] \\ &- 2\mathbb{E}_{\tilde{q}_t} [\langle \mathbf{x}_t - \mathbb{E}_{\tilde{q}_t}[\mathbf{x}_t | \mathbf{x}_0, \mathbf{x}_{t+1}], \boldsymbol{\mu}_\theta(\mathbf{x}_{t+1}, t+1) - \mathbb{E}_{\tilde{q}_t}[\mathbf{x}_t | \mathbf{x}_0, \mathbf{x}_{t+1}] \rangle] \\ &= \mathbb{E}_{\tilde{q}_t} [\|\mathbf{x}_t - \mathbb{E}_{\tilde{q}_t}[\mathbf{x}_t | \mathbf{x}_0, \mathbf{x}_{t+1}]\|^2] \\ &+ \left( \sqrt{\bar{\alpha}_t} - \sqrt{1 - \bar{\alpha}_t - \sigma_{t+1}^2} \sqrt{\frac{\bar{\alpha}_{t+1}}{1 - \bar{\alpha}_{t+1}}} \right) \mathbb{E}_{\tilde{q}_t(\mathbf{x}_0, \mathbf{x}_{t+1})} [\|\mathbf{x}_0 - \boldsymbol{\mu}_\theta(\mathbf{x}_{t+1}, t+1)\|^2],\end{aligned}\tag{30}$$

where the equality is due to the property of conditional expectation leads to  $\mathbb{E}_{\tilde{q}_t} [\langle \mathbf{x}_t - \mathbb{E}_{\tilde{q}_t}[\mathbf{x}_t | \mathbf{x}_0, \mathbf{x}_{t+1}], \boldsymbol{\mu}_\theta(\mathbf{x}_{t+1}, t+1) - \mathbb{E}_{\tilde{q}_t}[\mathbf{x}_t | \mathbf{x}_0, \mathbf{x}_{t+1}] \rangle] = 0$ , and rewriting  $\mathbb{E}_{\tilde{q}_t} [\|\boldsymbol{\mu}_\theta(\mathbf{x}_{t+1}, t+1) - \mathbb{E}_{\tilde{q}_t}[\mathbf{x}_t | \mathbf{x}_0, \mathbf{x}_{t+1}]\|^2]$  as in equations (5)-(10) in (Ho et al., 2020). Due to this, we know that minimizing the

<sup>7</sup>We can do this since (12) only relates to  $\tilde{q}_t(\mathbf{x}_{t+1})$

square error is equivalent to minimizing the  $\mathbb{E}_{\tilde{q}_t(\mathbf{x}_t, \mathbf{x}_{t+1})} [\|\mathbf{x}_0 - \mathbf{x}_\theta(\mathbf{x}_{t+1}, t+1)\|^2]$ . On the other hand, since  $\tilde{q}_t^* \in B$ , then we have

$$\begin{aligned} & D_{KL}(q(\mathbf{x}_{t+1} | \mathbf{x}_0) \| \tilde{q}_t^*(\mathbf{x}_{t+1} | \mathbf{x}_0)) \\ &= D_{KL}(q(\mathbf{x}_0 | \mathbf{x}_{t+1}) \| \tilde{q}_t^*(\mathbf{x}_0 | \mathbf{x}_{t+1})) + D_{KL}(q(\mathbf{x}_{t+1}) \| \tilde{q}_t^*(\mathbf{x}_{t+1})) \\ &\geq \eta_0. \end{aligned} \quad (31)$$

Thus, we prove our conclusion.  $\square$

**Theorem 1.** *There exists  $\delta_t$  depends on  $\mathbf{x}_0$  and  $\epsilon_t$  makes (13) equivalent to problem (12).*

$$\min_{\theta} \sum_{t=0}^{T-1} \mathbb{E}_{q(\mathbf{x}_0), \epsilon_t} \left[ \left\| \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_t + \delta_t, t) - \epsilon_t - \frac{\delta_t}{\sqrt{1 - \bar{\alpha}_t}} \right\|^2 \right]. \quad (13)$$

*Proof.* By combining Lemma 1, suppose the supreme is attained under  $\tilde{q}_{t-1}$  such that  $\mathbf{x}_t \sim \tilde{q}_{t-1}(\mathbf{x}_t)$  with

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_t + \delta_t, \quad (32)$$

with  $\delta_t$  depends on  $\mathbf{x}_0$  and  $\mathbf{x}_t$ . Then we prove the conclusion.  $\square$

### A.2.2 PROOF OF PROPOSITION 5

**Proposition 5.** *If  $L_t^{\text{DRO}}(\theta) \leq \eta_0$  in (12) for all  $t$ , and  $D_{KL}(q(\mathbf{x}_T) \| p_\theta(\mathbf{x}_T)) \leq \eta_0$ , then  $D_{KL}(q(\mathbf{x}_0) \| p_\theta(\mathbf{x}_0)) \leq \eta_0$ .*

*Proof.* This theorem can proved by induction. Since  $D_{KL}(q(\mathbf{x}_T) \| p_\theta(\mathbf{x}_T)) \leq \eta_0$ , then, let  $\tilde{q}_{T-1}^*(\mathbf{x}_T) = p_\theta(\mathbf{x}_T)$  and satisfies  $\tilde{q}_{T-1}^*(\mathbf{x}_T) = q(\mathbf{x}_{T-1})$ . The existence of such distribution is due to Kolmogorov existence theorem (Shiryaev, 2016). Then, we have

$$\begin{aligned} D_{KL}(\tilde{q}_{T-1}^*(\mathbf{x}_{T-1}) \| p_\theta(\mathbf{x}_{T-1})) &\leq D_{KL}(\tilde{q}_{T-1}^*(\mathbf{x}_T) \| p_\theta(\mathbf{x}_T)) \\ &\quad + D_{KL}(\tilde{q}_{T-1}^*(\mathbf{x}_{T-1} | \mathbf{x}_T) \| p_\theta(\mathbf{x}_{T-1} | \mathbf{x}_T)) \\ &\leq L_t^{\text{DRO}}(\theta) \\ &\leq \eta_0, \end{aligned} \quad (33)$$

where the first inequality is due to the definition of  $L_t^{\text{DRO}}(\theta)$  and  $\tilde{q}_{T-1}^*(\mathbf{x}_T) = p_\theta(\mathbf{x}_T)$ . Then, we prove our conclusion by induction over  $t$ .  $\square$

### A.2.3 PROOF OF PROPOSITION 6

**Proposition 6.** *For  $\eta > 0$  and  $\delta_t$  in (13),  $\|\delta_t\|_1 \leq \eta$  holds with probability at least  $1 - \sqrt{2(1 - \bar{\alpha}_t)/\eta}$ .*

*Proof.* Due to the definition of the first order Wasserstein distance  $W_1(\cdot, \cdot)$  (Villani et al., 2009) for any specific  $\mathbf{x}_0$ , suppose

$$\pi^* \in \arg \min_{\pi(\mathbf{x}_t, \tilde{\mathbf{x}}_t) \in q_t(\mathbf{x}_t | \mathbf{x}_0) \times \tilde{q}_t(\tilde{\mathbf{x}}_t | \mathbf{x}_0)} \mathbb{E} [\|\tilde{\mathbf{x}}_t - \mathbf{x}_t\|_1], \quad (34)$$

so that

$$\mathbb{E}_{\pi^*} [\|\tilde{\mathbf{x}}_t - \mathbf{x}_t\|_1] = W_1(q_t(\mathbf{x}_t | \mathbf{x}_0), \tilde{q}_t(\mathbf{x}_t | \mathbf{x}_0)). \quad (35)$$

Let  $\delta_t$  be the one of (13) under  $\pi^*$  derived by Lemma 1, then

$$\begin{aligned} \mathbb{P}(\|\delta_t\|_1 \geq \eta | \mathbf{x}_0) &\leq \frac{\mathbb{E}_{\pi^*} [\|\delta_t\|_1]}{\eta} \\ &= \frac{W_1(q_t(\mathbf{x}_t | \mathbf{x}_0), \tilde{q}_t(\mathbf{x}_t | \mathbf{x}_0))}{\eta} \\ &\leq \frac{\sqrt{2(1 - \bar{\alpha}_t)} D_{KL}(q_t(\mathbf{x}_t | \mathbf{x}_0) \| \tilde{q}_t(\mathbf{x}_t | \mathbf{x}_0))}{\eta} \\ &\leq \sqrt{\frac{2(1 - \bar{\alpha}_t)}{\eta}}, \end{aligned} \quad (36)$$

where inequality  $a$  is due to the Talagrand's inequality (Wainwright, 2019). Then we prove our conclusion.  $\square$

## B PROOFS IN SECTION 5

Next, we give the proof of results in Section 5. Firstly, let us check the definition of the  $\Phi_t(\mathbf{x}_{t+1})$ . For the variance-preserving stochastic differential equation in Song et al. (2022)

$$d\mathbf{z}_s = -\frac{\beta_s}{2}\mathbf{z}_s dt + \sqrt{\beta_s}dW_s. \quad (37)$$

Due to the solution of  $\mathbf{z}_s$  in Song et al. (2023), we know  $\mathbf{z}_{s_t}$  has the same distribution with  $\mathbf{x}_t$  in (1) for  $\{s_t\}_{t=1}^T$  satisfies

$$\exp\left(-\int_0^{s_t} \beta(u)du\right) = \bar{\alpha}_t \quad (s_0 = 0). \quad (38)$$

In the rest of this section, we use  $d(\mathbf{x}, \mathbf{y})$  in (15) as  $\ell_2$  distance  $\|\mathbf{x} - \mathbf{y}\|^2$ , whereas the conclusions under other distance can be similarly derived. Owing to the discussion in above, similar to (Song et al., 2023), when  $\mathbf{x}_{t+1} = \mathbf{z}_{s_{t+1}}$ , let  $\Phi_t(\mathbf{x}_{t+1}) = \Psi_{s_t}(\mathbf{z}_{s_{t+1}})$ , we can rewrite the objective (15) as follows.

$$\min_{\theta} \mathcal{L}_{CD}(\theta) = \min_{\theta} \sum_{t=0}^{T-1} \mathbb{E}_{\mathbf{z}_{s_t}} \left[ \|\mathbf{f}_{\theta}(\Psi_{s_t}(\mathbf{z}_{s_{t+1}}), t) - \mathbf{f}_{\theta}(\mathbf{z}_{s_{t+1}}, t+1)\|^2 \right]. \quad (39)$$

Here  $\mathbf{z}_s$  follows the following reverse time ODE of (37) with  $\mathbf{z}_0 \sim q(\mathbf{x}_0)$ ,

$$d\mathbf{z}_s = -\underbrace{\frac{\beta_s}{2}\left(\mathbf{z}_s + \frac{1}{2}\nabla_{\mathbf{z}} \log q_s(\mathbf{z}_s)\right)}_{\phi_s} ds, \quad (40)$$

and such  $\mathbf{z}_s$  has the same distribution with the ones in (37) (Song et al., 2022), where  $q_s$  is the density of  $\mathbf{z}_s$ .  $\Psi_{s_t}(\mathbf{z}_{s_{t+1}}) = \mathbf{z}_{s_{t+1}} - \int_{s_t}^{s_{t+1}} \phi_s(\mathbf{z}_s) ds$ , which is a deterministic function of  $\mathbf{z}_{s_{t+1}}$ , and  $\mathbf{f}_{\theta}(\mathbf{z}_{s_0}, 0) = \mathbf{z}_{s_0} = \mathbf{z}_0$ .

Now, we are ready to prove the Theorem 2 as follows.

**Theorem 2.** For  $\mathcal{L}_{CD}(\theta)$  in (15) with  $d(\cdot, \cdot)$  is  $\ell_2$  distance, then  $W_1(\mathbf{f}_{\theta}(\mathbf{x}_t, t), \mathbf{x}_0) \leq \sqrt{t\mathcal{L}_{CD}(\theta)}$ <sup>8</sup>.

*Proof.* Owing to the definition of  $W_1$ -distance, and the discussion in above, we have

$$\begin{aligned} W_1(\mathbf{f}_{\theta}(\mathbf{x}_T, T), \mathbf{x}_0) &= W_1(\mathbf{f}_{\theta}(\mathbf{z}_{s_T}, T), \mathbf{z}_{s_0}) \\ &= W_1(\mathbf{f}_{\theta}(\mathbf{z}_{s_T}, T), \Psi_{s_0}(\Psi_{s_1}(\cdots \Psi_{s_{T-1}}(\mathbf{z}_{s_T})))) \\ &\leq \mathbb{E} \left[ \|\mathbf{f}_{\theta}(\mathbf{z}_{s_T}, T) - \Psi_{s_0}(\Psi_{s_1}(\cdots \Psi_{s_{T-1}}(\mathbf{z}_{s_T})))\| \right] \\ &\leq \sum_{t=0}^{T-1} \mathbb{E} \left[ \|\mathbf{f}_{\theta}(\mathbf{z}_{s_{t+1}}, t+1) - \mathbf{f}_{\theta}(\Psi_{s_t}(\mathbf{z}_{s_{t+1}}), t)\| \right] \\ &\leq \sqrt{T\mathcal{L}_{CD}(\theta)}, \end{aligned} \quad (41)$$

where the first inequality is due to the definition of Wasserstein distance, the second and last inequalities respectively use the triangle inequality and Schwarz's inequality.  $\square$

### B.1 PROOF OF THEOREM 3

As pointed out in the above, the used  $\hat{\Phi}_t(\mathbf{x}_{t+1}, \epsilon_{\phi})$  is a numerical estimator of  $\Phi_t(\mathbf{x}_{t+1})$ . In the sequel, let us consider  $\hat{\Phi}$  is an Euler estimator as follows, whereas our analysis can be similarly generalized to the other estimators.

$$\hat{\Phi}_t(\mathbf{x}_{t+1}, \epsilon_{\phi}) = \hat{\Psi}_{s_t}(\mathbf{z}_{s_{t+1}}, \epsilon_{\phi}) = \mathbf{z}_{s_{t+1}} + \underbrace{(s_{t+1} - s_t) \frac{\beta_{s_{t+1}}}{2} \left( \mathbf{z}_{s_{t+1}} + \epsilon_{\phi}(\mathbf{z}_{s_{t+1}}, t+1) / \sqrt{1 - \bar{\alpha}_{t+1}} \right)}_{\hat{\phi}_{s_{t+1}}}, \quad (42)$$

<sup>8</sup>Here  $W_1(\mathbf{f}_{\theta}(\mathbf{x}_t, t), \mathbf{x}_0)$  is the Wasserstein 1-distance between distributions of  $\mathbf{f}_{\theta}(\mathbf{x}_t, t)$  and  $\mathbf{x}_0$ .

where  $\sqrt{1 - \bar{\alpha}_{t+1}} \epsilon_\phi(\mathbf{z}_{s_{t+1}}, t+1)$  estimates  $\nabla_{\mathbf{z}} \log q_{s_{t+1}}(\mathbf{z}_{s_{t+1}})$  as pointed out in (Song et al., 2020), and the condition  $\mathbf{x}_{t+1} = \mathbf{z}_{s_{t+1}}$  is hold.

Next, we illustrate the used regularity conditions to derive Theorem 3.

**Assumption 1.** The discretion error of  $\hat{\Psi}_{s_t}(\mathbf{z}_{s_{t+1}}, \epsilon_\phi)$  is smaller than  $C(s_{t+1} - s_t)^2$  for constant  $C$ , that says

$$\left\| \hat{\Psi}_{s_t}(\mathbf{z}_{s_{t+1}}, \epsilon_\phi) - \mathbf{z}_{s_{t+1}} - \int_{s_t}^{s_{t+1}} \hat{\phi}_s(\mathbf{z}_s) ds \right\| \leq C(s_{t+1} - s_t)^2 \quad (43)$$

**Assumption 2.** The estimated score  $\nabla_{\mathbf{z}} \log \hat{q}_s(\mathbf{z})$  has bounded expected error, i.e.,

$$\mathbb{E}_{\mathbf{z} \sim q_{s_t}(\mathbf{z})} \left[ \left\| \hat{\phi}_{s_t}(\mathbf{z}) - \phi_{s_t}(\mathbf{z}) \right\|^2 \right] \leq \epsilon. \quad (44)$$

for all  $0 \leq t < T$ .

**Assumption 3.** For the learned model  $f_\theta$ , it holds  $\|f_\theta\| \leq D$ .

The Assumption 1 describes the discretion error of the Euler method under ODE with drift term  $\hat{\phi}_s$ , which can be satisfied under proper continuity conditions of model  $\epsilon_\phi$ . On the other hand, Assumption 2 describes the estimation error of  $\hat{\phi}_{s_t}(\mathbf{z})$ , which terms out to be the training objective of obtaining it, see (Song et al., 2020) for more details. The Assumption 3 is natural, since  $f_\theta$  predicts  $\mathbf{x}_0$ , which is usually an image data with bounded norm. Now, we are ready to prove the Theorem 3, which is presented by proving the following formal version.

**Theorem 4.** Under Assumptions 1, 2, and 3, for all  $\delta_{s_t}$ , we have  $\mathbb{E}_{\mathbf{z}_{s_t}} [\|\delta_{s_t}(\mathbf{z}_{s_t})\|] \leq O(\Delta_{s_t}^2 + \epsilon \sqrt{\Delta_{s_t}})$  for  $\Delta_{s_t} = s_{t+1} - s_t$ . Besides that, we have

$$W_1(f_\theta(\mathbf{z}_T, T), \mathbf{z}_0) \leq \sqrt{T \hat{\mathcal{L}}_{CD}^{Adv}(\theta) + \frac{4D^2}{\eta} [C\Delta_{s_t}^2 + \epsilon O(\sqrt{\Delta_{s_t}})]}. \quad (45)$$

*Proof.* Noting that  $\Phi_t(\mathbf{x}_{t+1}) = \Psi_{s_t}(\mathbf{z}_{s_{t+1}})$  and  $\hat{\Phi}_t(\mathbf{x}_{t+1}, \epsilon_\phi) = \hat{\Psi}_{s_t}(\mathbf{z}_{s_{t+1}}, \epsilon_\phi)$ , the key problem is to upper bound the difference between  $\hat{\Psi}_{s_t}(\mathbf{z}, \epsilon_\phi)$  and  $\Psi_{s_t}(\mathbf{z})$  for all  $t$  and  $\mathbf{z}$ . To do so, we note that

$$\left\| \hat{\Psi}_{s_t}(\mathbf{z}, \epsilon_\phi) - \Psi_{s_t}(\mathbf{z}) \right\| \leq \left\| \hat{\Psi}_{s_t}(\mathbf{z}, \epsilon_\phi) - \mathbf{z} - \int_{s_t}^{s_{t+1}} \hat{\phi}_s(\mathbf{z}_s) ds \right\| + \left\| \mathbf{z} - \int_{s_t}^{s_{t+1}} \hat{\phi}_s(\mathbf{z}_s) ds - \Psi_{s_t}(\mathbf{z}) \right\|, \quad (46)$$

where the first one in r.h.s can be upper bounded by  $C(s_{t+1} - s_t)^2$  according to Assumption 1. On the other hand, define  $\frac{d\hat{\mathbf{z}}_s}{ds} = \hat{\phi}_s(\hat{\mathbf{z}}_s)$ , then when  $\hat{\mathbf{z}}_{s_{t+1}} = \mathbf{z}_{s_{t+1}} = \mathbf{z}$  and  $s \in [s_t, s_{t+1}]$ .

$$\begin{aligned} \frac{d}{ds} \|\hat{\mathbf{z}}_s - \mathbf{z}_s\|^2 &= \left\langle \hat{\mathbf{z}}_s - \mathbf{z}_s, \hat{\phi}_s(\hat{\mathbf{z}}_s) - \phi_s(\mathbf{z}_s) \right\rangle \\ &= \left\langle \hat{\mathbf{z}}_s - \mathbf{z}_s, \hat{\phi}_s(\hat{\mathbf{z}}_s) - \hat{\phi}_s(\mathbf{z}_s) + \hat{\phi}_s(\mathbf{z}_s) - \phi_s(\mathbf{z}_s) \right\rangle \\ &\leq L \|\hat{\mathbf{z}}_s - \mathbf{z}_s\|^2 + \left\langle \hat{\mathbf{z}}_s - \mathbf{z}_s, \hat{\phi}_s(\mathbf{z}_s) - \phi_s(\mathbf{z}_s) \right\rangle \\ &\leq \left( \frac{1}{2} + L \right) \|\hat{\mathbf{z}}_s - \mathbf{z}_s\|^2 + \frac{1}{2} \left\| \hat{\phi}_s(\mathbf{z}_s) - \phi_s(\mathbf{z}_s) \right\|^2. \end{aligned} \quad (47)$$

Taking expectation over  $\mathbf{z}$ , by Gronwall's inequality, Assumption 2 and  $\hat{\mathbf{z}}_{s_{t+1}} = \mathbf{z}_{s_{t+1}}$ , we have

$$\mathbb{E} [\|\hat{\mathbf{z}}_{s_t} - \mathbf{z}_{s_t}\|^2] \leq \int_{s_t}^{s_{t+1}} \frac{e^{(1/2+L)(s-s_t)}}{2} \mathbb{E} [\|\hat{\phi}_s(\mathbf{z}_s) - \phi_s(\mathbf{z}_s)\|^2] ds \leq \frac{\epsilon}{4} \int_{s_t}^{s_{t+1}} \beta_s e^{(1/2+L)(s-s_t)} ds. \quad (48)$$

Plugging this into (46), we know

$$\mathbb{E} \left[ \left\| \hat{\Psi}_{s_t}(\mathbf{z}_{s_t}, \epsilon_\phi) - \Psi_{s_t}(\mathbf{z}_{s_t}) \right\| \right] \leq C(s_{t+1} - s_t)^2 + \epsilon O(\sqrt{s_{t+1} - s_t}). \quad (49)$$

By Markov's inequality, we have

$$\begin{aligned} \mathbb{P} \left( \left\| \hat{\Psi}_{s_t}(\mathbf{z}_{s_t}, \epsilon_\phi) - \Psi_{s_t}(\mathbf{z}_{s_t}) \right\| \geq \eta \right) &\leq \frac{\mathbb{E} \left[ \left\| \hat{\Psi}_{s_t}(\mathbf{z}_{s_t}, \epsilon_\phi) - \Psi_{s_t}(\mathbf{z}_{s_t}) \right\| \right]}{\eta} \\ &\leq \frac{1}{\eta} [C(s_{t+1} - s_t)^2 + \epsilon O(\sqrt{s_{t+1} - s_t})]. \end{aligned} \quad (50)$$

Thus,

$$\begin{aligned}
& \mathbb{E} \left[ \|f_{\theta}(\mathbf{x}_{t+1}, t+1) - f_{\theta}(\Phi_t(\mathbf{x}_{t+1}), t)\|^2 \right] \\
&= \mathbb{E} \left[ \|f_{\theta}(\mathbf{z}_{s_{t+1}}, t+1) - f_{\theta}(\Psi_{s_t}(\mathbf{z}_{s_{t+1}}), t)\|^2 \right] \\
&= \mathbb{E} \left[ \left\| f_{\theta}(\mathbf{z}_{s_{t+1}}, t+1) - f_{\theta}(\hat{\Psi}_{s_t}(\mathbf{z}_{s_{t+1}} + \delta_{s_t}, \epsilon_{\phi}), t) \right\|^2 \right] \\
&= \mathbb{E} \left[ \left( \mathbf{1}_{\|\delta_{s_t}\| > \eta} + \mathbf{1}_{\|\delta_{s_t}\| \leq \eta} \right) \left\| f_{\theta}(\mathbf{z}_{s_{t+1}}, t+1) - f_{\theta}(\hat{\Psi}_{s_t}(\mathbf{z}_{s_{t+1}} + \delta_{s_t}, \epsilon_{\phi}), t) \right\|^2 \right] \\
&\leq \mathbb{E} \left[ \sup_{\|\delta\| \leq \eta} \left\| f_{\theta}(\mathbf{z}_{s_{t+1}}, t+1) - f_{\theta}(\hat{\Psi}_{s_t}(\mathbf{z}_{s_{t+1}} + \delta_{s_t}, \epsilon_{\phi}), t) \right\|^2 \right] + 4D^2 \mathbb{P}(\|\delta_{s_t}\|^2 \geq \eta) \\
&\leq \mathbb{E} \left[ \sup_{\|\delta\| \leq \eta} \left\| f_{\theta}(\mathbf{z}_{s_{t+1}}, t+1) - f_{\theta}(\hat{\Psi}_{s_t}(\mathbf{z}_{s_{t+1}} + \delta, \epsilon_{\delta}), t) \right\|^2 \right] + \frac{4D^2}{\eta} [C(s_{t+1} - s_t)^2 + \epsilon O(\sqrt{s_{t+1} - s_t})].
\end{aligned} \tag{51}$$

Taking sum over  $t$  and combining Theorem 2, we prove our conclusion.  $\square$

Therefore, in this theorem, by taking  $\Delta_{s_t} = s_{t+1} - s_t$  close to zero, we get the results in Theorem 3.

## C THE CONNECTION TO STANDARD ADVERSARIAL TRAINING

In this section, we clarify why the proposed AT objective (14) is a general version of the standard AT objective proposed in (Madry et al., 2018) used for classification problems.

For classification problem, given model  $f_{\theta}(\mathbf{x})$ , data  $\mathbf{x}$ , and label  $y$ , it aims to minimize the adversarial training objective

$$\min_{\theta} \mathbb{E}_{(\mathbf{x}, y)} \left[ \sup_{\delta: \|\delta\| \leq \eta_0} \ell(f_{\theta}(\mathbf{x} + \delta), y) \right], \tag{52}$$

for some loss function  $\ell$  (e.g. cross entropy) and adversarial radius  $\eta_0$ . However, the objective is not directly generalized to the diffusion model, as its training objective is a regression problem instead of classification (52). Thus, we should refer to the general version of adversarial training as in (Yi et al., 2021; Sinha et al., 2018), where the training objective is  $\min_{\theta} \mathbb{E}_{\mathbf{x}}[\ell_{\theta}(\mathbf{x})]$ , and the adversarial training objective becomes

$$\min_{\theta} \mathbb{E}_{\mathbf{x}} \left[ \sup_{\delta: \|\delta\| \leq \eta_0} \ell_{\theta}(\mathbf{x} + \delta) \right], \tag{53}$$

where  $\ell_{\theta}$  is the parameterized loss function, and  $\mathbf{x}$  is data. Then, we can conclude our objective (14) follows the above formulation, such that the goal is represented as

$$\min_{\theta} \sum_{t=0}^{T-1} \mathbb{E}_{\mathbf{x}_0} \left[ \mathbb{E}_{\mathbf{x}_t | \mathbf{x}_0} \left[ \sup_{\delta: \|\delta\| \leq \eta_0} \ell_{\theta}^{\mathbf{x}_0}(\mathbf{x}_t + \delta) \right] \right], \tag{54}$$

compared with the original noise prediction objective  $\min_{\theta} \sum_{t=0}^{T-1} \mathbb{E}_{\mathbf{x}_0} [\mathbb{E}_{\mathbf{x}_t | \mathbf{x}_0} [\ell_{\theta}^{\mathbf{x}_0}(\mathbf{x}_t)]]$  (5), such that the loss function

$$\ell_{\theta}^{\mathbf{x}_0}(\mathbf{x}_t) = \left\| \epsilon_{\theta}(t, \mathbf{x}_t) - \frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t} \mathbf{x}_0}{\sqrt{1 - \bar{\alpha}_t}} \right\|^2. \tag{55}$$

This clarifies the equivalence of our objective (14) to general adversarial training.

## D ADVERSARIAL TRAINING ON CONSISTENCY TRAINING MODEL

In (Song et al., 2023), the consistency model can be even trained without estimator  $\hat{\phi}_s$ . They prove that the empirical consistency distillation loss  $\hat{\mathcal{L}}_{CD}(\theta)$  can be approximated by the following  $\mathcal{L}_{CT}(\theta)$

$$\mathcal{L}_{CT}(\theta) = \sum_{t=0}^{T-1} \mathbb{E}_{\mathbf{x}_{t+1} \sim q(\mathbf{x}_{t+1})} [\|f_{\theta}(\mathbf{x}_t, t) - f_{\theta}(\mathbf{x}_{t+1}, t+1)\|^2]. \tag{56}$$

In our adversarial regime, we can also prove that the desired  $\hat{\mathcal{L}}_{CD}^{Adv}(\theta)$  can be approximated by the following  $\mathcal{L}_{CT}^{Adv}(\theta)$  with adversarial perturbation

$$\mathcal{L}_{CT}^{Adv}(\theta) = \sum_{t=0}^{T-1} \mathbb{E}_{\mathbf{x}_{t+1} \sim q(\mathbf{x}_{t+1})} \left[ \sup_{\|\delta\| \leq \eta} \|f_{\theta}(\mathbf{x}_t + \delta, t) - f_{\theta}(\mathbf{x}_{t+1}, t+1)\|^2 \right]. \quad (57)$$

The results can be checked by the following theorem.

**Theorem 5.** *Suppose  $f_{\theta}(\mathbf{x}_t, t)$  is twice continuously differentiable with a bounded second derivative. Then*

$$\hat{\mathcal{L}}_{CD}^{Adv}(\theta) \lesssim \mathcal{L}_{CT}^{Adv}(\theta) + O\left(T - \sum_{t=1}^T \sqrt{\alpha_t} + T\eta^2\right), \quad (58)$$

where “ $\lesssim$ ” means approximately less than or equal.

*Proof.* Due to the continuity of  $f_{\theta}(\mathbf{x}, t)$ , for any  $\delta$  with  $\|\delta\| \leq \eta$ , by Taylor’s expansion on  $\mathbf{x}_{t+1}$  from  $\mathbf{x}_t + \delta$ , we have

$$\begin{aligned} \mathbb{E} [\|f_{\theta}(\mathbf{x}_t + \delta, t) - f_{\theta}(\mathbf{x}_{t+1}, t+1)\|^2] &= \mathbb{E} [\|f_{\theta}(\mathbf{x}_{t+1}, t) - f_{\theta}(\mathbf{x}_{t+1}, t+1)\|^2] \\ &+ \mathbb{E} \left[ (f_{\theta}(\mathbf{x}_{t+1}, t) - f_{\theta}(\mathbf{x}_{t+1}, t+1))^{\top} \nabla f_{\theta}(\mathbf{x}_{t+1}, t) (\mathbf{x}_t + \delta - \mathbf{x}_{t+1}) \right] + O(\mathbb{E} [\|\mathbf{x}_{t+1} - \mathbf{x}_t - \delta\|^2]). \end{aligned} \quad (59)$$

Due to the Taylor’s expansion  $f_{\theta}(\mathbf{x}_t + \delta, t) = f_{\theta}(\mathbf{x}_{t+1}, t) + \nabla f_{\theta}(\mathbf{x}_{t+1}, t)(\mathbf{x}_t + \delta - \mathbf{x}_{t+1}) + \mathcal{O}(\|\mathbf{x}_{t+1} - \mathbf{x}_t - \delta\|^2)$ . Then, from the formulation of  $\mathbf{x}_t$ , we know  $\mathbb{E} [\|\mathbf{x}_{t+1} - \mathbf{x}_t - \delta\|^2] = O(1 - \sqrt{\alpha_t} + \eta^2)$ . Noting that due to definition of  $s_t$ , we have

$$\begin{aligned} \mathbb{E}[\mathbf{x}_t \mid \mathbf{x}_{t+1} = \mathbf{z}_{s_{t+1}}] &= \mathbb{E}[\mathbf{z}_{s_t} \mid \mathbf{z}_{s_{t+1}}] \\ &= \frac{1}{\sqrt{\alpha_{t+1}}} (\mathbf{z}_{s_{t+1}} - (1 - \alpha_{t+1}) \nabla_{\mathbf{x}} \log q_{s_{t+1}}(\mathbf{z}_{s_{t+1}})) \\ &= \exp\left(\frac{1}{2} \int_{s_t}^{s_{t+1}} \beta_s ds\right) \left( \mathbf{z}_{s_{t+1}} - \left(1 - e^{\int_{s_t}^{s_{t+1}} \beta_s ds}\right) \nabla_{\mathbf{z}} \log q_{s_{t+1}}(\mathbf{z}_{s_{t+1}}) \right) \quad (60) \\ &\approx \left(1 + \frac{1}{2} \int_{s_t}^{s_{t+1}} \beta_s ds\right) \mathbf{z}_{s_{t+1}} + \frac{1}{2} \int_{s_t}^{s_{t+1}} \beta_s ds \nabla_{\mathbf{z}} \log q_{s_{t+1}}(\mathbf{z}_{s_{t+1}}) \\ &\approx \hat{\Psi}_{s_t}(\mathbf{z}_{s_{t+1}}, \sqrt{1 - \bar{\alpha}_{t+1}} \nabla_{\mathbf{z}} \log q_{s_{t+1}}), \end{aligned}$$

where the first equality is due to Tweedie’s formula i.e., Lemma 11 in (Bao et al., 2022), the “ $\approx$ ” is due to  $e^a \approx 1 + a$  when  $a \rightarrow 0$ , and the last  $\approx$  is due to Euler-Mayaruma discretion. Due to this, we notice that

$$\begin{aligned} &\mathbb{E} \left[ (f_{\theta}(\mathbf{x}_{t+1}, t) - f_{\theta}(\mathbf{x}_{t+1}, t+1))^{\top} \nabla f_{\theta}(\mathbf{x}_{t+1}, t) (\mathbf{x}_t + \delta - \mathbf{x}_{t+1}) \mid \mathbf{x}_{t+1} = \mathbf{z}_{s_{t+1}} \right] \\ &= \mathbb{E} \left[ (f_{\theta}(\mathbf{x}_{t+1}, t) - f_{\theta}(\mathbf{x}_{t+1}, t+1))^{\top} \nabla f_{\theta}(\mathbf{x}_{t+1}, t) (\mathbb{E}[\mathbf{x}_t + \delta \mid \mathbf{x}_{t+1} = \mathbf{z}_{s_{t+1}}] - \mathbf{x}_{t+1}) \mid \mathbf{x}_{t+1} = \mathbf{z}_{s_{t+1}} \right] \\ &\approx \mathbb{E} \left[ (f_{\theta}(\mathbf{z}_{s_{t+1}}, t) - f_{\theta}(\mathbf{z}_{s_{t+1}}, t+1))^{\top} \nabla f_{\theta}(\mathbf{z}_{s_{t+1}}, t) \left( \hat{\Psi}_{s_t}(\mathbf{z}_{s_{t+1}}, \nabla_{\mathbf{z}} \log q_{s_{t+1}}) + \mathbb{E}[\delta \mid \mathbf{z}_{s_{t+1}}] - \mathbf{z}_{s_{t+1}} \right) \right], \end{aligned} \quad (61)$$

where the first equality is due to the property of conditional expectation, and the second “ $\approx$ ” is due to (60). Combining this with (59), we have

$$\begin{aligned} &\mathbb{E} [\|f_{\theta}(\mathbf{x}_t + \delta, t) - f_{\theta}(\mathbf{x}_{t+1}, t+1)\|^2 \mid \mathbf{x}_{t+1} = \mathbf{z}_{s_{t+1}}] \\ &= \mathbb{E} [\|f_{\theta}(\mathbf{z}_{s_t} + \delta, t) - f_{\theta}(\mathbf{z}_{s_{t+1}}, t+1)\|^2 \mid \mathbf{z}_{s_{t+1}}] \\ &= \mathbb{E} [\|f_{\theta}(\mathbf{z}_{s_{t+1}}, t) - f_{\theta}(\mathbf{z}_{s_{t+1}}, t+1)\|^2] \\ &+ \mathbb{E} \left[ (f_{\theta}(\mathbf{z}_{s_{t+1}}, t) - f_{\theta}(\mathbf{z}_{s_{t+1}}, t+1))^{\top} \nabla f_{\theta}(\mathbf{z}_{s_{t+1}}, t) \left( \hat{\Psi}_{s_t}(\mathbf{z}_{s_{t+1}}, \nabla_{\mathbf{z}} \log q_{s_{t+1}}) + \mathbb{E}[\delta \mid \mathbf{z}_{s_{t+1}}] - \mathbf{z}_{s_{t+1}} \right) \right] \\ &+ O(1 - \sqrt{\alpha_t} + \eta^2) \\ &= \mathbb{E} \left[ \left\| f_{\theta}(\hat{\Psi}_{s_t}(\mathbf{z}_{s_{t+1}}, \nabla_{\mathbf{z}} \log q_{s_{t+1}}) + \delta, t) - f_{\theta}(\mathbf{z}_{s_{t+1}}, t+1) \right\|^2 \right] + O(1 - \sqrt{\alpha_t} + \eta^2), \end{aligned} \quad (62)$$

where the last equality is due to Taylor’s expansion from  $f_{\theta}(\hat{\Psi}_{s_t}(\mathbf{z}_{s_{t+1}}, \nabla_{\mathbf{z}} \log q_{s_{t+1}}) + \delta, t)$  to  $f_{\theta}(\mathbf{z}_{s_{t+1}}, t)$ . Due to the arbitrariness of  $\delta$ , we prove our conclusion.  $\square$

## E IMPLEMENTATION DETAILS

### E.1 HYPERPARAMETERS OF DIFFUSION MODELS

For the diffusion models, all methods adopt the ADM model (Dhariwal & Nichol, 2021) as the backbone and follow the same training pipeline. Following existing work (Dhariwal & Nichol, 2021; Ning et al., 2023), we train models using the AdamW optimizer (Loshchilov & Hutter, 2019) with mixed precision training and the EMA rate is set to 0.9999. For CIFAR-10, the pretrained ADM is trained using a batch size of 128 for 250K iterations with a learning rate set to 1e-4. For ImageNet, the pretrained model is trained with a batch size of 1024 for 400K iterations, employing a learning rate of 3e-4. The models are trained in a cluster of NVIDIA Tesla V100s. More hyperparameters are reported in Table 4.

Table 4: Hyperparameters of diffusion model on each datasets.

Hyperparameters	CIFAR10 32 × 32	ImageNet 64 × 64
Channels	128	192
Batch size	128	1024
Learning rate	1e-4	3e-4
Fine-tuning iterations	200K	200K
Dropout	0.3	0.1
Noise schedule	Cosine	Cosine

### E.2 HYPERPARAMETERS OF LATENT CONSISTENCY MODELS

For experiments on Latent Consistency Models (LCM) (Luo et al., 2023), we train models on LAIOIN-Aesthetic-6.5+ (Schuhmann et al., 2022) at the resolution of 512×512, comprising 650K text-image pairs with predicted aesthetic scores higher than 6.5. Stable Diffusion v1.5 (Rombach et al., 2022) is adopted as the teacher model and initialized the student and target models in the latent consistency distillation framework. We set the range of the guidance scale  $[w_{min}, w_{max}] = [3, 5]$  during training and use  $w = 4$  in sampling because it performs better in our preliminary experiments, which is similar to DMD (Yin et al., 2024). The models are trained in a cluster of NVIDIA Tesla V100s. Both models of our AT and the original LCM training are trained from scratch with the same hyperparameters. We select the adversarial learning rate  $\alpha$  from  $\{0.02, 0.05\}$  and adversarial step  $K$  from  $\{2, 3\}$ . More details of hyperparameters are shown in Table 5 and other details of implementations can be found in the original LCM paper (Luo et al., 2023).

Table 5: Hyperparameters of latent consistency model.

Hyperparameters	LAIOIN-Aesthetic-6.5+
Batch size	64
Learning rate	8e-6
Training iterations	100K
EMA rate of target model	0.95
Conditional guidance scale $[w_{min}, w_{max}]$	$[3, 5]$

## F ADDITIONAL RESULTS

### F.1 RESULTS OF CLASSIFICATION ACCURACY SCORE

Classification Accuracy Score (CAS) (Ravuri & Vinyals, 2019) is proposed to evaluate the utility of the images produced by the generative model for downstream classification tasks. The underlying motivation for this metric is that if the generative model captures the real data distribution, the real data distribution can be replaced by the model-generated data and achieve similar results on downstream tasks like image classification.

Table 6: Comparison of CAS of different methods on CIFAR-10 32×32 dataset.

Methods	CAS
Real	92.5
<i>only using the synthetic data.</i>	
ADM	91.0
ADM-IP	89.2
ADM-AT (Ours)	<b>91.6</b>
<i>using the synthetic data with real data.</i>	
ADM	95.0
ADM-IP	94.9
ADM-AT (Ours)	<b>95.4</b>

Following the evaluation pipeline in Ravuri & Vinyals (2019), we train the image classifier in two settings: only on synthetic data or real data augmented with synthetic data, and use the classifier to predict labels on the test set of real data. Synthetic images are generated with a DDIM sampler under 20 NFEs. We use ResNet-18 (He et al., 2016) as the image classifier and train it for 200 epochs with a learning rate of 0.1 and a batch size of 128. We report CAS in the CIFAR-10 dataset at a resolution of 32×32 in Table 6. The results indicate that our method consistently performs better than other baseline methods on CAS metric in both settings. Although CAS with synthetic data cannot surpass real data, it demonstrates significant potential for enhancing classifier accuracy when employed as an augmentation technique alongside real data.

Table 7: Comparison of AT with TS-DDIM on CIFAR10 32×32. Both models are based on the ADM backbone. The results of TS are taken directly from the original paper.

Methods \ NFEs	50	20	10	5
ADM-TS-DDIM	3.52	5.35	10.73	26.94
ADM-AT (Ours)	<b>3.07</b>	<b>4.40</b>	<b>9.30</b>	<b>26.38</b>

## F.2 COMPARISON TO TS-DDIM

Li et al. (2024) introduces another approach named Time-Shift (TS) to alleviate the DPM distribution mismatch by searching for coupled time steps in sampling. Table 7 shows the comparison between our AT method with TS on CIFAR-10 with the DDIM Sampler. Both methods are based on the ADM pretrained model (Dhariwal & Nichol, 2021) as a backbone, which is the same as Section 6.2. We observe our method consistently better than the TS method across various sampling steps.

## F.3 RESULTS OF MORE NFES

We present results obtained with various samplers under 100 or 200 NFEs on CIFAR10 32x32 and ImageNet 64x64 in Table 8 and Table 9, respectively. The results show that our method is still effective for samplers under hundreds of NFEs.

## F.4 RESULTS OF MORE METRICS

We present the results of more generation quality metrics, including sFID, Inception Score (IS), Precision, and Recall, on CIFAR10 32x32 (Table 10 and Table 11) and ImageNet 64x64 (Table 12 and Table 13). The evaluation is performed following Dhariwal & Nichol (2021). We observe that our method shows effectiveness across these metrics.



Table 8: Sample quality measured by FID  $\downarrow$  of various sampling methods of DPM under 100 or 200 NFEs on CIFAR10 32x32.

Methods	IDDPM		DDIM		ES		DPM-Solver	
	100	200	100	200	100	200	100	200
ADM-FT	3.34	3.02	4.02	4.22	2.38	2.45	2.97	<b>2.97</b>
ADM-IP	2.83	2.73	6.69	8.44	2.97	3.12	10.10	10.11
ADM-AT (Ours)	<b>2.52</b>	<b>2.46</b>	<b>3.19</b>	<b>3.23</b>	<b>2.18</b>	<b>2.35</b>	<b>2.83</b>	3.00

Table 9: Sample quality measured by FID  $\downarrow$  of various sampling methods of DPM under 100 or 200 NFEs on ImageNet 64x64.

Methods	IDDPM		DDIM		ES		DPM-Solver	
	100	200	100	200	100	200	100	200
ADM-FT	3.88	3.48	4.71	4.38	3.07	<b>2.98</b>	<b>4.20</b>	4.13
ADM-IP	3.55	<b>3.08</b>	8.53	10.43	3.36	3.31	9.75	9.77
ADM-AT (Ours)	<b>3.35</b>	3.16	<b>4.58</b>	<b>4.34</b>	<b>3.05</b>	3.10	4.31	<b>4.10</b>

Table 10: Comparison of sFID  $\downarrow$  and IS  $\uparrow$  on CIFAR10 32x32.

(a) IDDPM

	5		8		10		20		50	
	sFID	IS	sFID	IS	sFID	IS	sFID	IS	sFID	IS
ADM	20.95	8.25	25.03	8.51	23.56	8.50	16.01	9.14	6.81	9.49
ADM-IP	25.81	7.02	<b>24.51</b>	8.04	<b>19.02</b>	8.50	8.99	9.28	5.32	<b>9.66</b>
ADM-AT	<b>19.78</b>	<b>8.71</b>	25.67	<b>8.66</b>	23.09	<b>8.77</b>	<b>6.01</b>	<b>9.30</b>	<b>5.04</b>	9.65

(b) DDIM

	5		8		10		20		50	
	sFID	IS	sFID	IS	sFID	IS	sFID	IS	sFID	IS
ADM	12.75	7.76	8.53	8.62	8.39	8.70	6.19	9.08	4.99	9.19
ADM-IP	15.53	7.55	8.00	<b>8.98</b>	7.12	<b>9.15</b>	<b>5.30</b>	<b>9.41</b>	5.64	9.49
ADM-AT	<b>12.56</b>	<b>7.97</b>	<b>7.93</b>	8.90	<b>7.08</b>	8.90	5.37	9.17	<b>4.66</b>	<b>9.51</b>

(c) ES

	5		8		10		20		50	
	sFID	IS	sFID	IS	sFID	IS	sFID	IS	sFID	IS
ADM	27.39	6.14	14.91	8.33	10.04	8.79	5.45	9.55	4.12	9.62
ADM-IP	34.70	5.73	16.84	8.23	10.89	8.88	4.94	9.59	4.08	9.70
ADM-AT	<b>16.84</b>	<b>6.97</b>	<b>10.33</b>	<b>8.60</b>	<b>8.00</b>	<b>8.95</b>	<b>4.78</b>	<b>9.65</b>	<b>4.04</b>	<b>9.77</b>

(d) DPM-Solver

	5		8		10		20		50	
	sFID	IS	sFID	IS	sFID	IS	sFID	IS	sFID	IS
ADM	11.82	8.00	5.79	9.12	<b>5.05</b>	9.41	<b>4.43</b>	9.78	<b>4.32</b>	9.82
ADM-IP	26.46	7.09	5.93	9.19	5.49	9.45	7.53	9.66	8.37	9.75
ADM-AT	<b>11.19</b>	<b>8.43</b>	<b>5.10</b>	<b>9.35</b>	5.29	<b>9.65</b>	4.75	<b>10.03</b>	4.59	<b>9.93</b>

## G MORE ANALYSIS

### G.1 EFFICIENT AT VS STANDARD AT

In this section, we conduct an ablation of the AT method in diffusion model training. We compare the performance of our used efficient AT and a standard AT method PGD on CIFAR-10 dataset at the resolution of 32 $\times$ 32. The adversarial step  $K$  is set to be 3 for both methods. We fine-tune both

Table 11: Comparison of Precision (P)  $\uparrow$  and Recall (R)  $\uparrow$  on CIFAR10 32x32.

(a) IDDPM										
	5		8		10		20		50	
	P	R	P	R	P	R	P	R	P	R
ADM	<b>0.54</b>	<b>0.47</b>	<b>0.59</b>	<b>0.45</b>	0.61	<b>0.46</b>	0.64	0.52	0.68	0.58
ADM-IP	<b>0.54</b>	0.39	<b>0.59</b>	0.43	0.61	<b>0.46</b>	0.66	0.54	0.68	0.59
ADM-AT	0.52	<b>0.47</b>	0.57	<b>0.45</b>	<b>0.62</b>	<b>0.46</b>	<b>0.68</b>	<b>0.55</b>	<b>0.69</b>	<b>0.59</b>
(b) DDIM										
	5		8		10		20		50	
	P	R	P	R	P	R	P	R	P	R
ADM	0.57	<b>0.47</b>	0.59	0.52	0.61	0.52	0.64	0.52	0.63	0.60
ADM-IP	0.57	0.44	<b>0.62</b>	<b>0.53</b>	<b>0.63</b>	<b>0.56</b>	0.65	<b>0.60</b>	0.65	<b>0.61</b>
ADM-AT	<b>0.59</b>	0.46	<b>0.62</b>	0.52	<b>0.63</b>	0.54	<b>0.65</b>	0.58	<b>0.66</b>	<b>0.61</b>
(c) ES										
	5		8		10		20		50	
	P	R	P	R	P	R	P	R	P	R
ADM	0.54	0.37	0.60	0.48	0.61	0.52	0.64	0.52	0.63	0.60
ADM-IP	0.46	0.32	0.58	0.45	0.62	0.51	<b>0.67</b>	<b>0.58</b>	<b>0.68</b>	0.60
ADM-AT	<b>0.61</b>	<b>0.45</b>	<b>0.64</b>	<b>0.51</b>	<b>0.65</b>	<b>0.54</b>	0.65	<b>0.58</b>	0.66	<b>0.61</b>
(d) DPM-Solver										
	5		8		10		20		50	
	P	R	P	R	P	R	P	R	P	R
ADM	0.61	0.47	<b>0.65</b>	0.58	<b>0.65</b>	0.59	0.66	0.61	0.63	<b>0.62</b>
ADM-IP	0.49	0.32	<b>0.65</b>	0.58	<b>0.65</b>	0.59	0.62	0.58	0.61	0.56
ADM-AT	<b>0.62</b>	<b>0.49</b>	<b>0.65</b>	<b>0.59</b>	<b>0.65</b>	<b>0.61</b>	<b>0.67</b>	<b>0.62</b>	<b>0.65</b>	0.61

models from the same pretrained ADM model with 100K update iterations of the model. The results are shown in Table 14. We report the results of 4 sampler settings (method-NFEs): IDDPM-50, DDIM-50, ES-20, and DPM-Solver-10.

We observe that efficient AT achieves performance comparable to or even better than PGD with the same model update iterations while accelerating the training ( $2.6\times$  speed-up). Thus, we propose applying the efficient AT method for our adversarial training framework.

## G.2 CONVERGENCE OF AT ON DIFFUSION MODELS

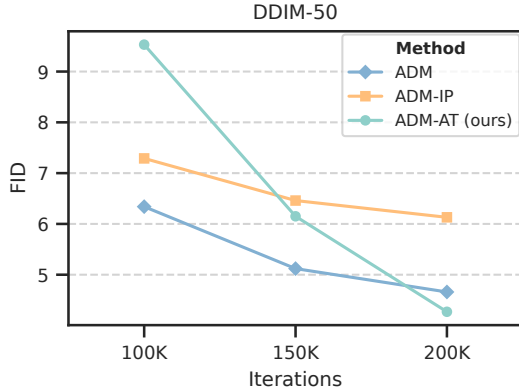
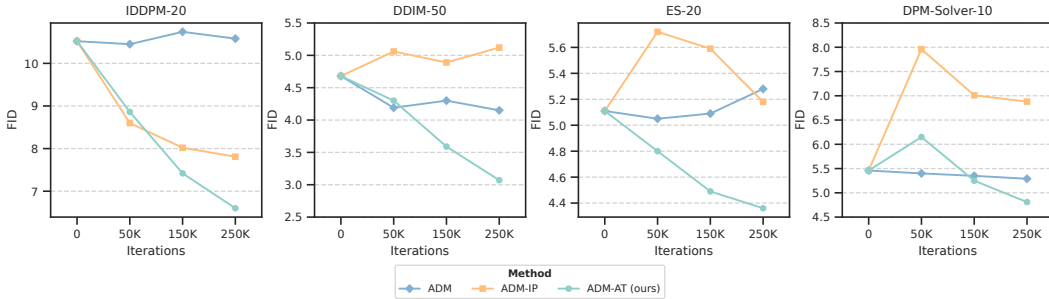


Figure 2: The convergence of methods trained from scratch on CIFAR-10  $32 \times 32$ . We use the DDIM sampler with 50 NFEs for sampling.

Table 12: Comparison of sFID  $\downarrow$  and IS  $\uparrow$  on ImageNet 64x64.

(a) IDDPM										
	5		8		10		20		50	
	sFID	IS	sFID	IS	sFID	IS	sFID	IS	sFID	IS
ADM	26.17	12.55	<b>36.34</b>	22.61	40.52	26.55	26.08	39.10	11.35	45.68
ADM-IP	40.90	12.19	47.98	23.47	37.72	27.86	25.06	<b>39.40</b>	6.75	44.87
ADM-AT	<b>24.82</b>	<b>14.50</b>	37.04	<b>23.84</b>	<b>36.50</b>	<b>30.03</b>	<b>22.83</b>	39.12	<b>5.69</b>	<b>46.25</b>
(b) DDIM										
	5		8		10		20		50	
	sFID	IS	sFID	IS	sFID	IS	sFID	IS	sFID	IS
ADM	27.74	14.30	14.27	25.88	12.78	28.29	8.84	33.54	6.31	38.08
ADM-IP	52.08	10.21	16.40	22.03	11.70	25.94	9.09	32.04	15.14	31.62
ADM-AT	<b>25.49</b>	<b>14.82</b>	<b>10.68</b>	<b>26.62</b>	<b>9.22</b>	<b>29.29</b>	<b>6.41</b>	<b>34.33</b>	<b>4.66</b>	<b>39.36</b>
(c) ES										
	5		8		10		20		50	
	sFID	IS	sFID	IS	sFID	IS	sFID	IS	sFID	IS
ADM	34.55	13.29	42.32	24.98	34.44	29.36	14.44	40.45	6.41	45.36
ADM-IP	44.81	10.07	41.01	22.44	30.12	27.66	<b>10.13</b>	39.50	<b>4.67</b>	44.69
ADM-AT	<b>29.72</b>	<b>16.49</b>	<b>33.58</b>	<b>27.85</b>	<b>27.64</b>	<b>31.94</b>	10.22	<b>42.18</b>	5.10	<b>45.59</b>
(d) DPM-Solver										
	5		8		10		20		50	
	sFID	IS	sFID	IS	sFID	IS	sFID	IS	sFID	IS
ADM	25.70	24.34	11.08	34.77	8.05	<b>37.45</b>	5.35	<b>40.54</b>	4.69	<b>41.31</b>
ADM-IP	42.68	16.93	<b>7.47</b>	33.85	7.22	33.57	14.74	31.29	18.99	30.32
ADM-AT	<b>20.79</b>	<b>26.32</b>	7.60	<b>34.89</b>	<b>6.36</b>	36.51	<b>4.51</b>	38.79	<b>4.22</b>	39.10

Figure 3: The convergence of methods fine-tuned from a same pretrained model on CIFAR-10  $32 \times 32$ . We compare the performance of methods on various samplers.

In classification tasks, adding adversarial perturbations usually slows the convergence of model training (Zhu et al., 2020). We are interested to see whether AT also affects the convergence of the diffusion training process.

Firstly, we explore the convergence of models trained from scratch. We utilize DDIM as the sampler with 50 NFEs and the results are shown in Figure 2. We observe that our AT method and ADM-IP exhibit slower convergence compared to ADM at the beginning (before 100K iterations), while as training more iterations (200K), our AT method shows a notable advantage.

Moreover, we explore the convergence of models under fine-tuning setting and the results are shown in Figure 3. We observe under this setting, when given a pretrained diffusion model like ADM, fine-tuning it with our proposed AT improves performance faster than other baselines. Overall, we observe that incorporating AT with a diffusion framework does not affect the convergence of the model much, especially in the fine-tuning setting.

Table 13: Comparison of Precision (P)  $\uparrow$  and Recall (R)  $\uparrow$  on ImageNet 64x64.

(a) IDDPM										
	5		8		10		20		50	
	P	R	P	R	P	R	P	R	P	R
ADM	0.34	0.48	0.46	<b>0.50</b>	0.51	0.48	0.65	0.52	0.73	0.57
ADM-IP	0.39	0.39	<b>0.50</b>	0.45	<b>0.56</b>	0.48	0.68	<b>0.55</b>	0.73	<b>0.60</b>
ADM-AT	<b>0.40</b>	<b>0.50</b>	<b>0.50</b>	<b>0.50</b>	0.55	<b>0.49</b>	<b>0.69</b>	0.52	<b>0.77</b>	0.59
(b) DDIM										
	5		8		10		20		50	
	P	R	P	R	P	R	P	R	P	R
ADM	0.42	<b>0.47</b>	0.54	<b>0.56</b>	0.58	<b>0.58</b>	0.65	0.60	0.69	<b>0.61</b>
ADM-IP	0.38	0.40	0.51	0.53	0.55	0.57	0.63	<b>0.61</b>	0.62	<b>0.61</b>
ADM-AT	<b>0.44</b>	0.43	<b>0.58</b>	0.55	<b>0.62</b>	0.56	<b>0.69</b>	0.59	<b>0.72</b>	<b>0.61</b>
(c) ES										
	5		8		10		20		50	
	P	R	P	R	P	R	P	R	P	R
ADM	0.40	0.44	0.52	0.47	0.58	0.48	0.69	0.55	0.73	0.59
ADM-IP	0.37	0.35	0.49	0.44	0.56	<b>0.49</b>	0.68	<b>0.57</b>	0.72	<b>0.60</b>
ADM-AT	<b>0.44</b>	<b>0.46</b>	<b>0.58</b>	<b>0.48</b>	<b>0.63</b>	<b>0.49</b>	<b>0.73</b>	0.55	<b>0.76</b>	0.59
(d) DPM-Solver										
	5		8		10		20		50	
	P	R	P	R	P	R	P	R	P	R
ADM	0.51	0.49	0.65	0.58	0.67	<b>0.60</b>	0.69	<b>0.62</b>	0.69	<b>0.62</b>
ADM-IP	0.39	0.44	0.64	<b>0.60</b>	0.64	<b>0.60</b>	0.59	0.60	0.57	0.59
ADM-AT	<b>0.56</b>	<b>0.50</b>	<b>0.68</b>	0.57	<b>0.69</b>	0.59	<b>0.72</b>	0.60	<b>0.71</b>	0.61

Table 14: Comparison of different AT methods used in our AT framework. All models are trained with the same model-updating iterations while the efficient AT has less training time.

Methods	FID				Training Time Speedup
	IDDPM-50	DDIM-50	ES-20	DPM-Solver-10	
Standard AT PGD-3	4.02	<b>3.37</b>	6.42	7.60	1.0 $\times$
Efficient AT (Ours)	<b>3.97</b>	3.42	<b>5.98</b>	<b>6.05</b>	<b>2.6<math>\times</math></b>

Table 15: Comparison of different adversarial learning rate  $\alpha$  of our AT framework on CIFAR10 32x32. IDDPM is adopted as the inference sampler.

$\alpha \setminus$ NFEs	5	8	10	20	50
$\alpha = 0.05$	51.72	32.09	25.48	10.38	4.36
$\alpha = 0.1$	37.15	23.59	15.88	6.60	3.34
$\alpha = 0.5$	63.73	40.08	27.57	7.23	3.42

Table 16: Comparison of different adversarial learning rate  $\alpha$  of our AT framework on ImageNet 64x64. IDDPM is adopted as the inference sampler.

$\alpha \setminus$ NFEs	5	8	10	20	50
$\alpha = 0.1$	56.92	27.39	24.06	10.17	5.82
$\alpha = 0.5$	45.65	23.79	19.18	8.28	4.01
$\alpha = 0.8$	46.92	28.46	22.47	9.70	4.25

Table 17: Comparison of different perturbation norms ( $l_1, l_2, l_\infty$ ) of our AT framework on CIFAR10 32x32.

Perturbation Norm	IDDPM-50	DDIM-50	ES-20	DPM-Solver-10
$l_1$	4.45	4.91	4.72	5.05
$l_2$	3.34	3.07	4.36	4.81
$l_\infty$	3.87	3.63	4.48	5.32

### G.3 MORE ABLATION STUDY

**Ablation on  $\alpha$**  We investigate the impact of adversarial learning rate  $\alpha$  in our framework. The results of various  $\alpha$  on CIFAR10 32x32 and ImageNet 64x64 are shown in Table 15 and Table 16, respectively. We observe that  $\alpha$  set to 0.1 is better on CIFAR10 32x32 and  $\alpha = 0.5$  is better for ImageNet 64x64. That says, the image in larger size corresponds to larger optimal perturbation level  $\alpha$ . We speculate this is because we use the perturbation measured under  $l_2$ -norm, where the  $l_2$ -norm of vector will increase with its dimension.

**Ablation on perturbation norm** During our experiments, we adopt  $l_2$ -adversarial perturbation. Actually, perturbations in Euclidean space under different  $l_p$  norm are equivalent with each other, e.g., for vector  $\delta \in \mathbb{R}^d$ , it holds  $\|\delta\|_\infty \leq \|\delta\|_2 \leq \sqrt{d}\|\delta\|_\infty$ . Therefore, we select  $\|\cdot\|_2$  as representation in our paper. Next, we explore the proposed ADM-AT under different adversarial perturbations.

The results are in Table 17. We found that our method under  $l_2$ -perturbation is more stable and indeed has better performance, thus we suggest to use  $l_2$ -perturbation as in the main body of this paper.

### G.4 QUALITATIVE COMPARISONS

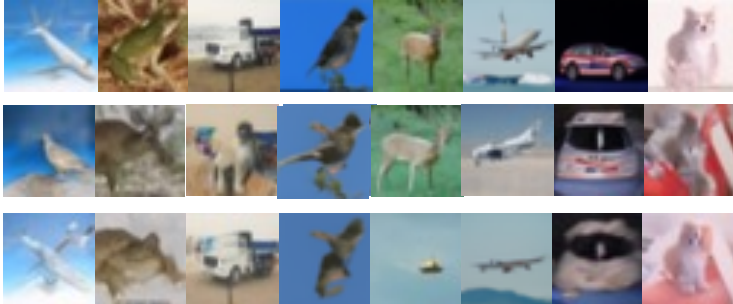


Figure 4: The qualitative comparisons of ADM-AT (top, FID 6.60), ADM-IP (middle, FID 7.81), and ADM (bottom, FID 10.58) on CIFAR10  $32 \times 32$ . We use the IDDPM sampler with 20 NFEs for sampling.

Figure 4, 5, 6, 7 show the qualitative comparisons between our proposed AT method and baselines. Our proposed AT method generates more realistic and higher-fidelity samples. We attribute this to our AT algorithm mitigates the distribution mismatch problem.



Figure 5: The qualitative comparisons of ADM-AT (top, FID 6.20), ADM-IP (middle, FID 8.40) and ADM (bottom, FID 8.32) on ImageNet  $64 \times 64$ . We use the DDIM sampler with 20 NFEs for sampling.



Figure 6: The qualitative comparisons of LCM (left) and LCM-AT (right) with one-step generation. The text prompt is *A photo of beautiful mountain with realistic sunset and blue lake, highly detailed, masterpiece.*

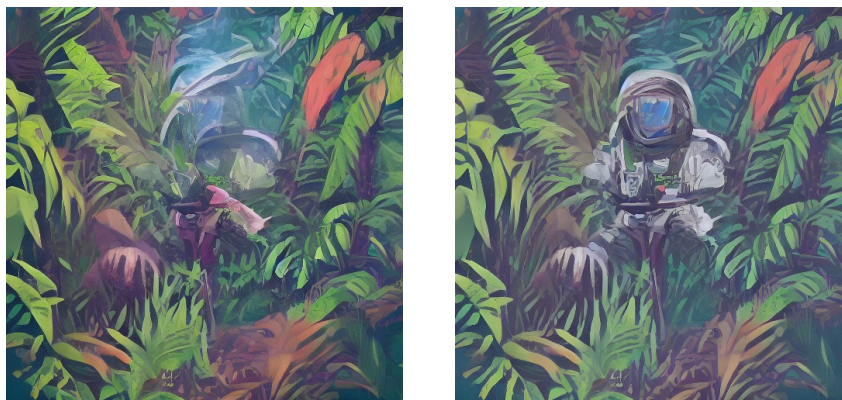


Figure 7: The qualitative comparisons of LCM (left) and LCM-AT (right) with one-step generation. The text prompt is *Astronaut in a jungle, cold color palette, muted colors, detailed, 8k.*