DIFF-INSTRUCT*: TOWARDS HUMAN-PREFERRED ONE-STEP TEXT-TO-IMAGE GENERATIVE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

In this paper, we introduce the Diff-Instruct* (DI*), an image data-free approach for building one-step text-to-image generative models that align with human preference while maintaining the ability to generate highly realistic images. We frame human preference alignment as online reinforcement learning using human feedback (RLHF), where the goal is to maximize the reward function while regularizing the generator distribution to remain close to a reference diffusion process. Unlike traditional RLHF approaches, which rely on the KL divergence for regularization, we introduce a novel score-based divergence regularization, which leads to significantly better performances. Although the direct calculation of this divergence remains intractable, we demonstrate that we can efficiently compute its gradient by deriving an equivalent yet tractable loss function. Remarkably, with Stable Diffusion V1.5 as the reference diffusion model, DI* outperforms all previously leading models by a large margin. When using the 2.6B Stable Diffusion XL architecture, the DI* results in a solid human-preferred one-step model that is able to generate aesthetic images of 1024×1024 resolutions. When using the 0.6B PixelArt- α model as the reference diffusion, DI* achieves a new record Aesthetic Score of 6.30 and an Image Reward of 1.31 with only a single generation step, almost doubling the scores of the rest of the models with similar sizes. It also achieves an HPSv2.0 score of 28.70, establishing a new state-of-the-art benchmark, with a better layout, richer details, and aesthetic colors.

029 030 031

032

004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

1 INTRODUCTIONS

Deep generative models have made substantial progress these years, largely transforming the content creation and editing across various domains (Karras et al., 2020; Nichol & Dhariwal, 2021; Poole et al., 2022; Kim et al., 2022; Tashiro et al., 2021; Meng et al., 2021; Couairon et al., 2022; Ramesh et al., 2022; Esser et al., 2024). These models demonstrated exceptional capabilities in generating high-resolution outputs, such as photorealistic images, videos, audio, and 3D assets (Oord et al., 2016; Ho et al., 2022; Poole et al., 2022; Brooks et al., 2024), and others (Zhang et al., 2023; Xue et al., 2023; Luo & Zhang, 2024; Luo et al., 2023b; Zhang et al., 2023b; Feng et al., 2023; Deng et al., 2024; Luo et al., 2024; Geng et al., 2024; Wang et al., 2024; Pokle et al., 2022).

Within this field, two types of generative models have gained significant attention, diffusion models
and one-step generators. Diffusion models (DMs) (Sohl-Dickstein et al., 2015; Ho et al., 2020),
or score-based generative models (Song et al., 2020), first progressively corrupt data with diffusion
processes and then train models to approximate the score functions of the noisy data distributions
across varying noise levels. The learned score functions can be used in the reverse process to generate high-quality samples by iterative denoising the noisy samples through stochastic differential
equations. While DMs can produce high-quality outputs, they often require a large number of model
evaluations, which limits their efficiency in applications.

Different from diffusion models, one-step generators (Zheng & Yang, 2024; Kang et al., 2023a;
Sauer et al., 2023a; Yin et al., 2024; Zhou et al., 2024a; Luo et al., 2024c) have emerged as a
highly efficient alternative to multi-step diffusion models. Unlike DMs, one-step generative models
directly transform latent noises to output samples with a single neural network forward pass. This
mechanism significantly reduced the inference cost, making them ideal for real-time applications
such as text-to-image and text-to-video generations. Many existing works have demonstrated the



Figure 1: None cherry-picked generated images from one-step 0.6B DiT-DI* model with an recordbreaking HPSv2.0 score of 28.70. After being trained with Diff-Instruct*, the images show better layouts, rich colors, vivid details, and aesthetic appearance, making them favored in terms of human preferences. Refer to the Appendix B.1 for the prompts used in comparison.



Figure 2: A visual comparison of our 0.6B DiT-DI* and 0.86B SD1.5-DI* models against other text-t0-image models. Refer to the Appendix B.1 for the prompts used in comparison.

108 leading performances of one-step text-to-image generators (Zhou et al., 2024a; Yin et al., 2024; 109 Luo et al., 2024c) by employing diffusion distillation (Luo, 2023). However, these works focus 110 on matching the generator distributions with pretrained diffusion models, without considering the 111 critical challenge of teaching one-step text-to-image models to satisfy human preferences, which is 112 one of the most important needs in the age of human-centric AI.

113 To close this gap, we introduce Diff-Instruct* (DI*), a novel approach to train human-preferred one-114 step generators. Inspired by the success of reinforcement learning using human feedback (RLHF) in 115 training large language models (Christiano et al., 2017; Ouyang et al., 2022), we frame the human-116 preference alignment problem as maximizing the rewards with a score-based divergence constraint. 117 This yields generated samples that not only adhere to user prompts but also improve in aesthetic 118 quality. Our approach differs from traditional RLHF methods, which rely on Kullback-Leibler (KL) divergence for distribution regularization. Instead, we propose score-based divergences, which offer 119 more stable training dynamics and better final performance. While the direct computation of the 120 score-based divergence is intractable, we develop a novel solution by deriving a tractable gradient 121 computation method, leading to a tractable pseudo-loss that stands equivalent to the intractable one 122 in the gradient perspective. With DI*, we are able to train large-scale human-preferred one-step 123 text-to-image generative models. 124

In our evaluation, our best 0.6B DiT-DI*-1step model, using diffusion transformer (DiT) (Peebles & 125 Xie, 2022) as generator architecture and PixelArt- α (Chen et al., 2023) as reference diffusion model, 126 has set a new benchmark on the COCO-2017 validation prompts dataset. It achieves an ImageRe-127 ward (Xu et al., 2023a) of 1.31 and an Aesthetic Score(Schuhmann, 2022) of 6.30, significantly out-128 performing the previous leading models such as 2.6B SDXL-DMD2 4-step models (Yin et al., 2024) 129 and Stable Diffusion XL (Podell et al., 2023) with 25 model steps by 50.5% and 77.0%, respec-130 tively. Moreover, this model also reaches a new high with Human Preference Score V2.0 (Wu et al., 131 2023) of **28.70**, surpassing leading models like SDXL-DMD2, the PixelArt- α , DALL-E 2 (Ramesh 132 et al., 2022) and the Stable Diffusion XL with Refiners. When using Stable Diffusion V1.5 as the 133 reference model and the same UNet architecture for the generator, our SD1.5-DI*-1step model with 134 only 0.86B parameters exceeds the performance of 2.6B SDXL-DMD2 (Yin et al., 2024) and SD1.5-135 DPO (Wallace et al., 2024). These results establish Diff-Instruct* as a versatile, architecture-flexible 136 approach for aligning one-step text-to-image generators with human preferences.

137 138 139

145

151 152

2 PRELIMINARY

140 **Diffusion Models.** In this section, we introduce preliminary knowledge and notations about dif-141 fusion models. Assume we observe data from the underlying distribution $q_d(x)$. The goal of gener-142 ative modeling is to train models to generate new samples $x \sim q_0(x)$. Under mild conditions, the 143 forward diffusion process of a diffusion model can transform initial distribution q_0 towards some simple noise distribution, 144

$$d\boldsymbol{x}_t = \boldsymbol{F}(\boldsymbol{x}_t, t)dt + G(t)d\boldsymbol{w}_t, \qquad (2.1)$$

146 where F is a pre-defined vector-valued drift function, G(t) is a pre-defined scalar-value diffusion 147 coefficient, and w_t denotes an independent Wiener process. A continuous-indexed score network 148 $s_{\varphi}(x,t)$ is employed to approximate marginal score functions of the forward diffusion process (2.1). 149 The learning of score networks is achieved by minimizing a weighted denoising score matching 150 objective (Vincent, 2011; Song et al., 2020),

$$\mathcal{L}_{DSM}(\varphi) = \int_{t=0}^{T} \lambda(t) \mathbb{E}_{\substack{\boldsymbol{x}_{0} \sim q_{0}, \\ \boldsymbol{x}_{t} \mid \boldsymbol{x}_{0} \sim p_{t}(\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0})}} \|\boldsymbol{s}_{\varphi}(\boldsymbol{x}_{t}, t) - \nabla_{\boldsymbol{x}_{t}} \log p_{t}(\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0})\|_{2}^{2} \mathrm{d}t.$$
(2.2)

153 Here the weighting function $\lambda(t)$ controls the importance of the learning at different time levels and 154 $p_t(\boldsymbol{x}_t|\boldsymbol{x}_0)$ denotes the conditional transition of the forward diffusion (2.1). After training, the score 155 network $s_{\varphi}(\boldsymbol{x}_t, t) \approx \nabla_{\boldsymbol{x}_t} \log q_t(\boldsymbol{x}_t)$ is a good approximation of the marginal score function of the 156 diffused data distribution. 157

3 **PREFERENCE ALIGNMENT USING DIFF-INSTRUCT***

159 160

158

In this section, we introduce Diff-Instruct*, a general method tailored for training one-step text-161 to-image generator according to human preference. We begin with outlining the problem setup and defining notations. We frame the training objective as a special form of reinforcement learning using
 human feedback (RLHF). We then introduce a family of score-based probability divergences and
 show how these divergences effectively regularize the RLHF process to align closely with human
 aesthetics and content preferences. Next, we introduce the classifier-free reward and discuss how to
 balance the explicit human reward model and the implicit classifier-free reward.

168 **Problem Setup.** Assume we have a human reward function $r(x_0, c)$, which encodes the human 169 preference for an image x_0 and corresponding text description c. Besides, we also have a pre-trained 170 diffusion model which will later act as a reference distribution $p_{ref}(x_0) = q_0(x_0)$. We assume the 171 reference diffusion model is specified by the score function

$$\boldsymbol{s}_{q_t}(\boldsymbol{x}_t) \coloneqq \nabla_{\boldsymbol{x}_t} \log q_t(\boldsymbol{x}_t) \tag{3.1}$$

where $q_t(x_t)$'s is the underlying distribution diffused at time t according to (2.1). We assume that the pre-trained diffusion model well captures the ground-truth data distribution, and thus will be the only item of consideration for our approach.

177 Our goal is to train a human-preferred one-step generator model g_{θ} that generates images by directly 178 mapping a random noise $z \sim p_z$ to obtain $x_0 = g_{\theta}(z, c)$, conditioned on the input text $c \sim C$. 179 The generator's output distribution, $p_{\theta}(x_0|c)$, should maximize the expected human rewards while 180 adhering to the constraints on the divergence to the reference distribution $p_{ref}(\cdot)$. Let $D(\cdot, \cdot)$ be a 181 distribution divergence. For any fixed prompt c, the training objective is defined as:

$$\theta^* = \operatorname*{argmax}_{\theta} \mathbb{E}_{\boldsymbol{x}_0 \sim p_{\theta}(\boldsymbol{x}_0 | \boldsymbol{c})} \big[r(\boldsymbol{x}_0, \boldsymbol{c}) \big], \quad \text{s.t. } \boldsymbol{D}(p_{\theta}(\cdot | \boldsymbol{c}), p_{ref}(\cdot | \boldsymbol{c})) \le \delta$$
(3.2)

The constraint $D(p_{\theta}(\cdot|c), p_{ref}(\cdot|c)) \le \delta$ is often referred to as a *trust-region* in literature of reinforcement learning (Schulman, 2015; Schulman et al., 2017). This objective (3.2) is equivalent to the minimization of objective (3.3) that adds the constraining term onto the negative reward:

$$\theta^* = \operatorname*{argmin}_{\theta} \mathbb{E}_{\boldsymbol{x}_0 \sim p_{\theta}(\boldsymbol{x}_0 | \boldsymbol{c})} \left[-\alpha r(\boldsymbol{x}_0, \boldsymbol{c}) \right] + \boldsymbol{D}(p_{\theta}, p_{ref})$$
(3.3)

190 Here α is a coefficient that balances reward influences and $D(\cdot)$ acts as a regularization term. 191 Traditional reinforcement learning from human feedback (RLHF) methods in large language mod-192 els (Ouyang et al., 2022) use the Kullback-Leibler divergences for regularization. Recent work (Luo, 193 2024) has studied using the integral of KL divergence in objective (3.3) to train one-step text-to-194 image generators. However, since the KL divergences are defined with the density ratio of two 195 distributions, any misalignment of density supports will lead to severe numerical instability, poten-196 tially resulting in annoying mode-seeking behavior (Bishop, 2006). Besides, some recent works have shown that score-based divergences (Zhou et al., 2024b; Luo et al., 2024c) result in better and 197 more stable performance than KL divergences (Luo et al., 2024b; Yin et al., 2023; Nguyen & Tran, 198 2023) in the literature of diffusion distillation. Motivated by the above two inspirations, we propose 199 a novel score-based online PPO algorithm that uses a general family of score-based divergence as 200 regularization in this paper. We provide solid theoretical foundations of score-based regularization 201 and its better empirical performances through large-scale text-to-image models. 202

203 204

211 212 213

167

172

173

182

183 184

188 189

3.1 GENERAL SCORE-BASED DIVERGENCES

Different from KL divergence, we can define the regularization term $D(p_{\theta}, p_{ref})$ via the following general score-based divergence. Assume $\mathbf{d} : \mathbb{R}^d \to \mathbb{R}$ is a scalar-valued proper distance function (i.e., a non-negative function that satisfies $\forall x, \mathbf{d}(x) \ge 0$ and $\mathbf{d}(x) = 0$ if and only if x = 0). Given a parameter-independent sampling distribution π_t that has large distribution support, we can formally define a time-integral score divergence as

$$\mathbf{D}^{[0,T]}(p_{\theta}, p_{ref}) \coloneqq \int_{t=0}^{T} w(t) \mathbb{E}_{\boldsymbol{x}_{t} \sim \pi_{t}} \left\{ \mathbf{d}(\boldsymbol{s}_{p_{\theta,t}}(\boldsymbol{x}_{t}) - \boldsymbol{s}_{q_{t}}(\boldsymbol{x}_{t})) \right\} \mathrm{d}t,$$
(3.4)

where $p_{\theta,t}$ and q_t denote the marginal densities of the diffusion process (2.1) at time t initialized with $p_{\theta,0} = p_{\theta}$ and $q_0 = p_{ref}$ respectively. w(t) is an integral weighting function. Clearly, we have $\mathbf{D}^{[0,T]}(p_{\theta}, p_{ref}) = 0$ if and only if $p_{\theta}(\mathbf{x}_0) = p_{ref}(\mathbf{x}_0)$, a.s. π_0 .

Algorithm 1: Diff-Instruct* for training human-preferred one-step text-to-image generators.
Input: prompt dataset C , generator $g_{\theta}(\boldsymbol{x}_0 \boldsymbol{z}, \boldsymbol{c})$, prior distribution p_z , reward model $r(\boldsymbol{x}, \boldsymbol{c})$,
reward model scale α_{rew} , CFG reward scale α_{cfg} , reference diffusion model
$s_{ref}(x_t c,c)$, assistant diffusion $s_{\psi}(x_t t,c)$, forward diffusion $p_t(x_t x_0)$ (2.1),
assistant diffusion updates rounds K_{TA} , time distribution $\pi(t)$, diffusion model
weighting $\lambda(t)$, generator loss time weighting $w(t)$.
while not converge do
freeze θ , update ψ for K_{TA} rounds using SGD by minimizing
$\mathcal{L}(a b) = \mathbb{E}$ $\sum_{n=1}^{\infty} \sum_{j=1}^{n} a_{nj}(a_{j} +a_{j}) - \sum_{j=1}^{\infty} a_{nj}(a_{j} +a_{j}) ^{2} dt$
$\mathcal{L}(\psi) = \mathbb{E}_{\substack{\mathbf{c} \sim \mathcal{C}, \mathbf{z} \sim p_z, t \sim \pi(t) \\ \mathbf{x}_0 = g_\theta(\mathbf{z} \mathbf{c}), \mathbf{x}_t \mid \mathbf{x}_0 \sim p_t(\mathbf{x}_t \mathbf{x}_0)}} \Lambda(t) \ \mathbf{s}_{\psi}(\mathbf{x}_t t, \mathbf{c}) - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t \mathbf{x}_0)\ _2 \mathrm{d}t.$
for a lot of the SCD to activity in the
There ψ , update θ using SGD by minimizing loss
$ \mathcal{L}_{DI*}(\theta) = \mathbb{E}_{\substack{\mathbf{c} \sim \mathcal{C}, \mathbf{z} \sim p_z, \\ \mathbf{x}_0 = g_{\theta}(\mathbf{z}, \mathbf{c})}} \left\{ -\alpha_{rew} \cdot r(\mathbf{x}_0, \mathbf{c}) + \mathbb{E}_{\substack{t \sim \pi(t), \\ \mathbf{x}_t \mid \mathbf{x}_0 \sim p_t(\mathbf{x}_t \mid \mathbf{x}_0)}} \right\} $
(1) (1) (1) (1) (1) (1) (1) (1) (1) (1)
$-w(t)\{\mathbf{d}'(s_{\psi}(\boldsymbol{x}_{t} t,\boldsymbol{c})-s_{ref}(\boldsymbol{x}_{t} t,\boldsymbol{c}))\} \{s_{\psi}(\boldsymbol{x}_{t} t,\boldsymbol{c})-\nabla_{\boldsymbol{x}_{t}}\log p_{t}(\boldsymbol{x}_{t} \boldsymbol{x}_{0})\}$
$\left[\begin{array}{c} (1) \left(\begin{array}{c} (1) \left(\begin{array}{c} 1 \end{array}\right) \right) \right] \right] $
$+ \alpha_{cfg} \cdot w(t) \{ \boldsymbol{s}_{ref}(\operatorname{sg}[\boldsymbol{x}_t] t, \boldsymbol{c}) - \boldsymbol{s}_{ref}(\operatorname{sg}[\boldsymbol{x}_t] t, \boldsymbol{\omega}) \} \boldsymbol{x}_t \} $ (3.7)

end

 return θ, ψ .

3.2 DIFF-INSTRUCT*

Recall that g_{θ} is a one-step model, therefore samples from p_{θ} can be implemented through a direct mapping $\boldsymbol{x}_0 = g_{\theta}(\boldsymbol{z}|\boldsymbol{c})$. With the score-based regularization term (3.4), for each given text prompt \boldsymbol{c} , we can formally write down our training objective to minimize as:

$$\mathcal{L}_{Orig}(\theta) = \mathbb{E}_{\substack{\boldsymbol{x}_{0} = g_{\theta}(\boldsymbol{x}, \boldsymbol{c})}} \left[-\alpha r(\boldsymbol{x}_{0}, \boldsymbol{c}) \right] + \mathbf{D}^{[0,T]}(p_{\theta}, p_{ref})$$
(3.5)

Now we are ready to reveal the objective of Diff-Instruct* that we use to train human-preferred onestep generator g_{θ} . Notice that directly minimizing objective (3.5) is intractable because we do not know the relationship between θ and corresponding $p_{\theta,t}$. However, we show in Theorem 3.1 that an equivalent tractable loss (3.6) will have the same θ gradient as the intractable loss function (3.5):

$$\mathcal{L}_{DI*}(\theta) = \mathbb{E}_{\substack{\boldsymbol{x} \sim p_{z}, \\ \boldsymbol{x}_{0} = g_{\theta}(\boldsymbol{x})}} \left[-\alpha r(\boldsymbol{x}_{0}, \boldsymbol{c}) + \int_{t=0}^{T} w(t) \mathbb{E}_{\substack{\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0} \\ \sim p_{t}(\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0})}} \left\{ -\mathbf{d}'(\boldsymbol{y}_{t}) \right\}^{T} \left\{ \boldsymbol{s}_{p_{\mathrm{sg}[\theta], t}}(\boldsymbol{x}_{t}) - \nabla_{\boldsymbol{x}_{t}} \log p_{t}(\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0}) \right\} \mathrm{d}t \right]$$
(3.6)

with $oldsymbol{y}_t\coloneqqoldsymbol{s}_{p_{ ext{sg}[m{ heta}],t}}(oldsymbol{x}_t)-oldsymbol{s}_{q_t}(oldsymbol{x}_t).$

Theorem 3.1. Under mild assumptions, if we take the sampling distribution in (3.4) as $\pi_t = p_{sg[\theta],t}$, then the θ gradient of (3.5) is the same as the objective (3.6): $\frac{\partial}{\partial \theta} \mathcal{L}_{Orig}(\theta) = \frac{\partial}{\partial \theta} \mathcal{L}_{DI*}(\theta)$.

We will give the proof in Appendix A.1. In practice, we can use another assistant diffusion model $s_{\psi}(\boldsymbol{x}_t, t)$ to approximate the generator model's score function $s_{p_{sg[\theta],t}}(\boldsymbol{x}_t)$ pointwise, which was also done in the literature of diffusion distillations works such as Zhou et al. (2024a;b); Luo et al. (2024b;c); Yin et al. (2023; 2024). Therefore, we can alternate between 1) updating the assistant diffusion $s_{\psi}(x_t, t)$ using generator-generated samples (which are efficient) and 2) updating the gen-erator by minimizing the tractable objective (3.6). We name our training method that minimizes the objective $\mathcal{L}_{DI*}(\theta)$ in (3.6) the Diff-Instruct* because it is inspired by Diff-Instruct(Luo et al., 2024b) and Diff-Instruct++(Luo, 2024) that involves an additional diffusion model and a reward model to train one-step generators.

3.3 DECOUPLING THE EXPLICIT REWARD AND IMPLICIT GUIDANCE-BASED REWARDS

Classifier-free Guidance Corresponds to Implicit Reward. In previous sections, we have shown in theory that with explicitly available reward models, we can readily train the one-step generator to



Figure 3: Comparison of Aesthetic Scores and Image Reward on 1K MSCOCO 2017 validation prompts of Score-based (DI*) and KL divergence (DI++(Luo, 2024)) for alignment with (3.5).

align with human preference. In this section, we enhance the DI* by incorporating the classifier-free reward that is implied by the classifier-free guidance of diffusion models.

The classifier-free guidance (Ho & Salimans, 2022) (CFG) uses a modified score function of a form $\tilde{\alpha}_{-}(m,t|\alpha) := \alpha_{-}(m,t|\alpha) + (\sqrt{\alpha_{-}(m,t|\alpha)}) + (\sqrt{\alpha_{-}(m,t|\alpha)})$

$$\boldsymbol{s}_{ref}(\boldsymbol{x}_t, t | \boldsymbol{c}) \coloneqq \boldsymbol{s}_{ref}(\boldsymbol{x}_t, t | \boldsymbol{\varnothing}) + \omega \{ \boldsymbol{s}_{ref}(\boldsymbol{x}_t, t | \boldsymbol{c}) - \boldsymbol{s}_{ref}(\boldsymbol{x}_t, t | \boldsymbol{\varnothing}) \}$$

to replace the original conditions score function $s_{ref}(x_t, t|c)$. Using CFG for diffusion models empirically leads to better sampling quality.

As is first pointed out by Luo (2024), the classifier-free guidance is related to an implicit reward function. In this part, we derive a tractable loss function that minimizes the so-called classifier-free reward, which we use together with the explicit reward $r(\cdot, \cdot)$ in DI*.

Theorem 3.2. Under mild conditions, if we set an implicit reward function as (3.9), the loss (3.8)

$$\mathcal{L}_{cfg}(\theta) = \int_{t=0}^{T} \mathbb{E}_{\substack{\boldsymbol{z} \sim p_{z}, \boldsymbol{x}_{0} = g_{\theta}(\boldsymbol{z}, \boldsymbol{c}) \\ \boldsymbol{x}_{t} \mid \boldsymbol{x}_{0} \sim p(\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0})}} w(t) \left\{ \boldsymbol{s}_{ref}(\mathrm{sg}[\boldsymbol{x}_{t}] \mid t, \boldsymbol{c}) - \boldsymbol{s}_{ref}(\mathrm{sg}[\boldsymbol{x}_{t}] \mid t, \boldsymbol{\varnothing}) \right\}^{T} \boldsymbol{x}_{t} \mathrm{d}t \qquad (3.8)$$

has the same gradient as the negative implicit reward function (3.9)

$$-r(\boldsymbol{x}_{0},\boldsymbol{c}) = -\int_{t=0}^{T} \mathbb{E}_{\boldsymbol{x}_{t} \sim p_{\theta,t}} w(t) \log \frac{p_{ref}(\boldsymbol{x}_{t}|t,\boldsymbol{c})}{p_{ref}(\boldsymbol{x}_{t}|t)} \mathrm{d}t.$$
(3.9)

The notation $sg[x_t]$ means detaching the θ gradient on x_t . We give the proof of Theorem 3.2 in Appendix A.2. Theorem 3.2 gives a tractable loss function (3.8) aiming to minimize the negative classifier-free reward function. Therefore, we can scale and add this loss $\mathcal{L}_{cfg}(\theta)$ (3.8) to the DI* loss (3.6) to balance the effects of explicit reward and implicit CFG reward.

305 3.4 THE PRACTICAL ALGORITHM

306 Now it is time for us to introduce the practical algorithm. As Algorithm 1 (and a more executable 307 version in Algorithm 2) shows, the DI* involves three models, with one generator model g_{θ} , one 308 reference diffusion model s_{ref} and one assistant diffusion model s_{ψ} . The reference diffusion does 309 not need to be trained, while the generator and the assistant diffusion are updated alternatively. Two 310 hyper-parameters, the α_{rew} and α_{cfg} control the strength of the explicit reward and the implicit CFG reward during training. The explicit reward model can either be an off-the-shelf reward model, 311 such as the CLIP similarity score (Radford et al., 2021), or the Image Reward (Xu et al., 2023a) or 312 trained in-house with researchers' internal human feedback data. Due to page limitations, we put a 313 discussion about the meanings of hyper-parameters in Appendix B.2. 314

Flexibility in Distance Functions. Clearly, various choices of distance function d(.) result in different training algorithms. For instance, $\mathbf{d}(y_t) = \|y_t\|_2^2$ is a naive choice. Interestingly, such a distance function has been studied in pure diffusion distillation literature in Zhou et al. (2024b;a); Luo et al. (2024c). In this paper, we draw inspiration from (Luo et al., 2024c) and find that using the so-called pseudo-Huber distance leads to better performance. The distance and corresponding loss writes $\mathbf{d}(\mathbf{y}) \coloneqq \sqrt{\|\mathbf{y}_t\|_2^2 + c^2} - c$, and

321 322

323

279

281 282

283

284

285

286 287

289

290

291

292

293

295

304

$$\mathbf{D}^{[0,T]}(p_{\theta}, p_{ref}) = -\left\{\frac{\boldsymbol{y}_t}{\sqrt{\|\boldsymbol{y}_t\|_2^2 + c^2}}\right\}^T \left\{\boldsymbol{s}_{\psi}(\boldsymbol{x}_t, t) - \nabla_{\boldsymbol{x}_t} \log p_t(\boldsymbol{x}_t | \boldsymbol{x}_0)\right\}.$$
(3.10)

Here $\boldsymbol{y}_t\coloneqq \boldsymbol{s}_{p_{\mathrm{sg}[\theta],t}}(\boldsymbol{x}_t) - \boldsymbol{s}_{q_t}(\boldsymbol{x}_t).$

³²⁴ 4 RELATED WORKS

326 Diffusion Distillation Through Divergence Minimization. Diff-Instruct* is inspired by research 327 on diffusion distillation (Luo, 2023) which aims to minimize certain distribution divergence to train 328 one-step generators. Luo et al. (2024b) first study the diffusion distillation by minimizing the Inte-329 gral KL divergence. Yin et al. (2023) generalize such a concept and add a data regression loss for better performance. Zhou et al. (2024b) study the distillation by minimizing the Fisher divergence. 330 Luo et al. (2024c) study the distillation using the general score-based divergence. Many other works 331 also introduced additional techniques and improved the performance (Geng et al., 2023; Kim et al., 332 2023; Song et al., 2023; Song & Dhariwal; Nguyen & Tran, 2023; Song et al., 2024; Yin et al., 2024; 333 Zhou et al., 2024a; Heek et al., 2024; Xie et al., 2024; Salimans et al., 2024; Geng et al., 2024). 334

335 Preference Alignment for Diffusion Models and One-step Generators. In recent years, many 336 works have emerged trying to align diffusion models with human preferences. There are three main 337 lines of alignment methods for diffusion models. 1) The first kind of method fine-tunes the diffusion 338 model over a specifically curated image-prompt dataset (Dai et al., 2023; Podell et al., 2023). 2) 339 the second line of methods tries to maximize some reward functions either through the multi-step 340 diffusion generation output (Prabhudesai et al., 2023; Clark et al., 2023; Lee et al., 2023) or through 341 policy gradient-based RL approaches (Fan et al., 2024; Black et al., 2023). For these methods, the 342 backpropagation through the multi-step diffusion generation output is expensive and hard to scale. 343 3) the third line, such as Diffusion-DPO (Wallace et al., 2024), Diffusion-KTO (Yang et al., 2024), tries to directly improve the diffusion model's human preference property with raw collected data 344 instead of reward functions. Besides the human preference alignment of diffusion models, Diff-345 Instruct++(Luo, 2024) recently arose as the first attempt to improve human preferences for one-346 step generators. Though inspired by DI++, DI* uses score-based divergences which are technically 347 different from the KL divergence used in DI++. Besides, as we show in Section 5.2, DI* archives 348 better performances than DI++ in Luo (2024).

349 350

351

5 EXPERIMENTS

5.1 EXPERIMENT SETTTINGS

352 353

Experiment Settings for SD1.5 and SDXL Experiment. For experiments of SD1.5, we use the 354 open-sourced SD1.5 of a resolution of 512×512 as our reference diffusion in Algorithm 1. We 355 implement our experiments based on SiD-LSG (Zhou et al., 2024a) codebase. We construct the 356 one-step generator with the same architecture as the reference SD1.5 model, following the same 357 configuration of SiD-LSG. We use the prompts of the LAION-AESTHETIC dataset with an aes-358 thetic score larger than 6.25, which resulting a total of 3M text prompts. To better explore the ad-359 vantages of our score-based divergence over KL divergence, we refer to a recent work (Luo, 2024) 360 that uses KL divergence for training and conducts a detailed comparison between two divergences. 361 We explore different combinations and find that ($\alpha_{rew} = 1000, \alpha_{cfg} = 1.5$) in Algorithm 1 is the 362 best. For SD1.5 experiments, we find that training one-step generators from scratch leads to longer training, therefore we initialize our generator with the weights of the SiD-LSG pre-trained one-step 363 model. For experiments of SDXL, we use the SDXL(Podell et al., 2023) as the teacher model and 364 its architecture as the one-step student model. We initialize the one-step generator with the pretrained DMD2-SDXL-1step model (Yin et al., 2024), which is a pretty good initialization based on 366 SDXL architectures. We follow similar settings as the SD1.5 experiment: using a scale for explicit 367 ImageReward of 1000 and a scale for implicit CFG reward of 8.0.

368

369 **Experiment Settings for PixelArt**- α **Experiment.** PixArt- α is a high-quality DiT-based(Peebles 370 & Xie, 2022) open-sourced text-to-image diffusion model. We use the DiT architecture and the 371 0.6B PixArt- α of a resolution of 512×512 as our reference diffusion to demonstrate the compatibil-372 ity of DI* for different neural network architectures. We use the prompts from the training dataset 373 of PixArt- α (the SAM-Recaptioned Dataset), resulting in a total of 10M prompts. After exploring 374 different combinations, we find that ($\alpha_{rew} = 10, \alpha_{cfq} = 4.5$) in Algorithm 1 gives the best per-375 formances. For DiT generators with a total of 0.6B parameters, we find that training from scratch without initialization results in strong results. Therefore we do not apply special initialization for 376 DiT-based one-step generators. Our best DiT-DI* model in Table 2 is trained from scratch with 8 377 H100-80G GPUs for 200000 iterations with a batch size of 128. The total wall-clock training costs

Table 1: Quantitative comparisons of text-to-image models on MSCOCO-2017 validation prompts (the upper part) and Parti(Yu et al., 2022) Prompts (the under part). DI* is short for Diff-Instruct*. α_r and α_c are short for α_{rew} and α_{cfg} in Algorithm 1. † means our implementation. Data means the model needs image data for training. Sampling means the model needs to draw samples from reference diffusion models. Reward means the model needs a human reward model for training. † indicates our implementation. ‡ indicates the same 4-step model of DMD2 but with different inference steps.

MODEL	STEPS	Туре	PARAMS	IMAGE Reward	AES SCORE	PICK Score	CLIP Score	ADDITIONAL REQUIREMENTS
SD15-DPO(WALLACE ET AL., 2024)	15	UNET	0.86B	0.20	5.29	0.214	31.07	DATA, SAMPLING
SD15-DPO(WALLACE ET AL., 2024)	25	UNET	0.86B	0.28	5.37	0.218	31.25	PREFERENCE DATA
SD15-LCM(LUO ET AL., 2023A)	1	UNET	0.86B	-1.58	5.04	0.194	27.20	DATA, SAMPLING
SD15-LCM(LUO ET AL., 2023A)	4	UNET	0.86B	-0.23	5.40	0.214	30.11	DATA, SAMPLING
SD15-TCD(ZHENG ET AL., 2024)	1	UNET	0.86B	-1.49	5.10	0.196	28.30	DATA, SAMPLING
SD15-TCD(ZHENG ET AL., 2024)	4	UNET	0.86B	-0.04	5.28	0.212	30.43	DATA, SAMPLING
PERFLOW(YAN ET AL., 2024)	4	UNET	0.86B	-0.20	5.51	0.211	29.54	DATA, SAMPLING
SD15-Hyper(Ren et al., 2024)	1	UNET	0.86B	0.28	5.49	0.214	30.82	DATA, SAMPLING
SD15-Hyper(Ren et al., 2024)	4	UNET	0.86B	0.42	5.41	0.217	31.03	REWARD , SEG-MODEL
SD15-INSTAFLOW(LIU ET AL., 2023)	1	UNET	0.86B	-0.16	5.03	0.207	30.68	DATA, SAMPLING
SDXL-BASE(ROMBACH ET AL., 2022)	25	UNET	2.6B	0.74	5.57	0.226	31.83	IMAGE-TEXT
SDXL-BASE(ROMBACH ET AL., 2022)	15	UNET	2.6B	0.68	5.56	0.224	31.99	IMAGE-TEXT
SDXL-DMD2 [‡] -1024(Yin et al., 2024)	1	UNET	2.6B	0.82	5.45	0.224	31.78	IMAGE-TEXT
SDXL-DMD2 [‡] -1024(Yin et al., 2024)	4	UNET	2.6B	0.87	5.52	0.231	31.50	IMAGE-TEXT
SDXL-DMD2‡-512(YIN ET AL., 2024)	1	UNET	2.6B	0.36	5.03	0.215	31.54	IMAGE-TEXT
SDXL-DMD2‡-512(YIN ET AL., 2024)	4	UNET	2.6B	-0.18	5.17	0.206	29.28	IMAGE-TEXT
SD15-DMD2-512(YIN ET AL., 2024)	1	UNET	2.6B	-0.12	5.24	0.211	30.00	IMAGE-TEXT
SD21-TURBO(SAUER ET AL., 2023B)	1	UNET	0.86B	0.56	5.47	0.225	31.50	IMAGE-TEXT
SD15-BASE(ROMBACH ET AL., 2022)	15	UNET	0.86B	0.08	5.25	0.212	30.99	IMAGE-TEXT
SD15-BASE(ROMBACH ET AL., 2022)	25	UNET	0.86B	0.22	5.32	0.216	31.13	IMAGE-TEXT
PIXELART- α -512(Chen et al., 2023)	25	DIT	0.6B	0.82	6.01	0.227	31.20	IMAGE-TEXT
PIXELART- α -512(Chen et al., 2023)	15	DIT	0.6B	0.82	6.03	0.226	31.16	IMAGE-TEXT
SD15-SIDLSG(ZHOU ET AL., 2024A)	1	UNET	0.86B	-0.18	5.16	0.210	30.04	Text
SDXL-DMD2-1024(YIN ET AL., 2024)	1	UNET	2.6B	0.85	5.46	0.225	31.86	IMAGE-TEXT
SD15-DI++(Luo, 2024) [†]	1	UNET	0.86B	0.82	5.78	0.219	30.30	REWARD
DIT-DI++(Luo, 2024) [†]	1	DIT	0.6B	1.24	6.19	0.225	30.80	REWARD
SD15-DI *($\alpha_r = 0, \alpha_c = 1.5$)	1	UNET	0.86B	0.34	5.27	0.217	30.83	REWARD
SD15-DI *($\alpha_r = 100, \alpha_c = 1.5$)	1	UNET	0.86B	0.62	5.44	0.218	30.76	REWARD
SD15-DI *($\alpha_r = 1000, \alpha_c = 4.5$)	1	UNET	0.86B	0.73	5.56	0.219	30.71	REWARD
SD15-DI *($\alpha_r = 1000, \alpha_c = 1.5$)	1	UNET	0.86B	0.94	5.83	0.220	30.49	REWARD
SDXL-DI*-1024 ($\alpha_r = 1000, \alpha_c = 8.0$) 1	UNET	2.6B	0.88	5.56	0.225	32.07	REWARD
DIT-DI *($\alpha_r = 1, \alpha_c = 4.5$)	1	DIT	0.6B	0.98	6.02	0.225	31.00	REWARD
DIT-DI *($\alpha_r = 10, \alpha_c = 4.5$)	1	DIT	0.6B	1.31	6.30	0.225	30.84	Reward
SDXL-BASE(ROMBACH ET AL., 2022)	15	UNET	2.6B	0.69	5.68	0.224	32.76	IMAGE-TEXT
PIXELART- α -512(CHEN ET AL., 2023)	15	DIT	0.6B	0.96	6.00	0.227	31.76	IMAGE-TEXT
SDXL-DMD2-1024(YIN ET AL., 2024)	1	UNET	2.6B	0.94	5.53	0.225	33.00	IMAGE-TEXT
SDXL-DI*-1024 ($\alpha_r = 1000, \alpha_c = 8.0$) 1	UNET	2.6B	1.06	5.61	0.225	33.27	REWARD
DIT-DI *($\alpha_r = 1, \alpha_c = 4.5$)	1	DIT	0.6B	0.96	6.06	0.224	31.11	REWARD
DIT-DI *($\alpha_r = 10, \alpha_c = 4.5$)	1	DIT	0.6B	1.26	6.23	0.225	30.64	REWARD

+14

415 416 417

418

419

is less than 30 hours. As a comparison, other industry models in Table 1 usually require hundreds of A100 GPU days for training. Such a low cost might benefit from the property that DI* needs neither real image data nor samples from reference diffusion models.

420 421

422 423

Quantitative Evaluations Metrics. We compare our generators with other leading open-sourced 424 models that are either based on SD1.5 diffusion models or larger models such as SDXL. For all 425 models, we compute four standard scores in Table 1: the Image Reward (Xu et al., 2023a), the Aes-426 thetic Score (Schuhmann, 2022), the PickScore(Kirstain et al., 2023), and the CLIP score(Radford 427 et al., 2021) on the same 1K prompts randomly sampled from COCO-2017-validation(Lin et al., 428 2014) set on the same computing devices. We have compared the results on 1K prompts and 30K 429 COCO prompts and found similar results. Since our training prompts do not involve COCO valida-430 tion prompts, the evaluation of COCO prompts can be viewed as an out-of-training set evaluation. 431 We also evaluate models with Human Preference Score v2.0 (HPSv2.0) (Wu et al., 2023) in Table 2. The HPSv2.0 is a standard score for evaluating models' human preferences across different styles.

432 5.2 PERFORMANCES AND FINDINGS

Quantitative Comparison: DI* Achieves SoTA Human Preference Scores. As Table 1 shows, our best 0.6B DiT-DI* model outperforms all other open-sourced models. It achieves a COCO (out-of-sample) Image Reward of 1.31, which is 50% better than the second best 2.6B SDXL-DMD2-4Step model of 1024 resolution. It also outperforms the SDXL model with a margin of 70%. It also shows an Aesthetic Score of 6.30, which is 14.3% better than the 2.6B SDXL-DMD2-4Step model. Among SD1.5-based models, our best SD1.5-DI* model outperforms other models with significant margins. This result demonstrates the compatibility of DI* across UNet and DiT architectures.

As Table 2 shows, our 0.6B DiT-DI* one-step model achieves a record-breaking HPSv2.0 score of
28.70 across open-sourced text-to-image models. It clearly outperforms the autoregressive models
such as 6B CogView2 (Ding et al., 2022), diffusion models such as 5.5B Dalle-E 2 (Ramesh et al.,
2022), 5B GLIDE(Nichol et al., 2021), 2.6B SDXL(Podell et al., 2023), and 4.3B DeepFloyd-XL
with up to 30+ generation steps. Such a lightweight, high-performance, and high-efficiency model
will bring big impacts on applications that need real-time generations.

447 Qualitative Comparisons. As Figure 1 and Figure 2 show, models trained with DI* show better layouts, richer colors, and more aesthetic face preferences. We find such an advantage is stronger 448 for scene generations. Figure 4 shows a visualization of the SD1.5-DI* model in Table 2 before 449 and after alignments with DI*. It shows that models after alignment produce images with richer 450 colors. We also evaluate the COCO-FID values of SD1.5-based one-step models with and without 451 preference alignment with DI*. The FID of the alignment model is 18.44, while the un-alignment 452 model (the SiD-LSG which we use as the initialization of the one-step generator) has an FID of 453 8.27. This phenomenon suggests that only using ImageReward(which tends to be subjective) to 454 align models with DI* can potentially harm objective metrics such as FID and CLIPScore. A naive 455 solution to hack CLIPScore is to simply include the CLIP score and PickScore as two other sources 456 of reward when using DI*, which will definitely improve these scores. However, since we want to 457 verify the generalization ability of one reward to others using DI*, we only use the ImageReward 458 and observed other scores in this paper.

459 Score-based Divergence is Better than KL Divergence. In Figure 3 we compare the use of 460 score-based divergence of DI* and the KL divergence of DI++(Luo, 2024) for RLHF regularization 461 in (3.5). We fix the $\alpha_{rew} = 1000$ and $\alpha_{cfg} = 1.5$ for both DI++ and DI* and use the same Image 462 Reward as an explicit reward model for training the SD1.5-based one-step generators. As we can 463 see in Figure 3, for each iteration, DI* has a better score than DI++. Besides, DI* achieves the best 464 final results. The reason for the worse performance of KL divergence might be its definition, which 465 involves the ratio of two distributions, which may lead to unstable numerical performances when 466 two distributions have misaligned density supports.

467 6 CONCLUSION AND LIMITATIONS

In this paper, we present Diff-Instruct*, a novel approach for aligning human-preferred one-step text-to-image generators. By formulating the training objective as a maximization of expected human reward functions with a score-based divergence regularization, we have developed practical losses and easy-to-implement algorithms. Our results show that DiT-based one-step text-to-image generators trained with DI* achieve new state-of-the-art human preference performances.

474 Nonetheless, DI* has its limitations. We empirically find three typical mistakes the generator often 475 makes. (1) The Generator Sometimes Generates Bad Human Faces and Hands. (2) The aligned 476 model still can not count correctly. Figure 5 shows some occasional bad generation cases. We be-477 lieve that consistently improving both the generator architecture and the reward models will lead to 478 better models. Besides, we are also interested in following directions that call for further research. 479 First, in this paper, we only study using only one reward model to train the generator. However, 480 training generators with multiple rewards is still unexplored. Second, the DI* needs a pre-trained reward model. However, methods like DPO (Rafailov et al., 2024; Wallace et al., 2024) have put a 481 new setup that trains generative models directly using human feedback data. Whether it is possible 482 to develop DPO-like algorithms based on DI* is a promising research direction. Third, it is also in-483 teresting to explore training multistep models instead of one-step generators for better performance. 484 We hope these future directions could contribute more to the community. 485

Table 2: **HPSv2.0** (upper table) and **HPSv2.1** (under table). We compare open-sourced models regardless of their base model and architecture. † indicates our implementation. ‡ indicates the same 4-step model of DMD2 but with different inference steps.

_						
	Model	ANIMATION	CONCEPT-ART	PAINTING	Рното	AVERAGE
	GLIDE (NICHOL ET AL., 2021)	23.34	23.08	23.27	24.50	23.55
	LAFITE (ZHOU ET AL., 2022)	24.63	24.38	24.43	25.81	24.81
	VQ-DIFFUSION (GU ET AL., 2022)	24.97	24.70	25.01	25.71	25.10
	FUSEDREAM (LIU ET AL., 2021)	25.26	25.15	25.13	25.57	25.28
	LATENT DIFFUSION (ROMBACH ET AL., 2022)	25.73	25.15	25.25	26.97	25.78
	COGVIEW2 (DING ET AL., 2022)	26.50	26.59	26.33	26.44	26.47
	DALL-E MINI	26.10	25.56	25.56	26.12	25.83
	VERSATILE DIFFUSION (XU ET AL., 2023B)	26.59	26.28	26.43	27.05	26.59
	VQGAN + CLIP (ESSER ET AL., 2021) DALL E 2 (PAMESH ET AL., 2022)	20.44	20.53	20.47	20.12	26.39
	STARLE 2 (RAMESH ET AL., 2022) STARLE DIFFUSION V1.4 (ROMBACH ET AL., 2022)	27.34	20.54	20.08	27.24	20.95
	STABLE DIFFUSION V1.4 (ROMBACH ET AL., 2022) STABLE DIFFUSION V2.0 (ROMBACH ET AL., 2022)	27.20	26.89	26.86	27.46	20.55
	EPIC DIFFUSION	27.57	26.96	27.03	27.49	27.26
	DEEPFLOYD-XL	27.64	26.83	26.86	27.75	27.27
	OPENJOURNEY	27.85	27.18	27.25	27.53	27.45
	MAJICMIX REALISTIC	27.88	27.19	27.22	27.64	27.48
	CHILLOUTMIX	27.92	27.29	27.32	27.61	27.54
	Deliberate	28.13	27.46	27.45	27.62	27.67
	REALISTIC VISION	28.22	27.53	27.56	27.75	27.77
	SDXL-BASE(PODELL ET AL., 2023)	28.42	27.63	27.60	27.29	27.73
	SDXL-REFINER(PODELL ET AL., 2023)	28.45	27.66	27.67	27.46	27.80
	DREAMLIKE PHOTOREAL 2.0	28.24	27.60	27.59	27.99	27.86
	SD15-DPO-15STEP(WALLACE ET AL., 2024)	27.11	26.75	26.70	27.30	26.97
	SD15-DPO-25STEP(WALLACE ET AL., 2024)	27.54	26.97	26.99	27.49	27.25
	SD15-LCM-1STEP(LUO ET AL., 2023A)	23.35	23.41	23.53	23.81	23.52
	SD15-LCM-4Step(Luo et al., 2023a)	26.42	25.79	25.95	26.91	26.27
	SD15-TCD-1STEP(ZHENG ET AL., 2024)	23.37	23.16	23.26	23.88	23.42
	SD15-TCD-4Step(Zheng et al., 2024)	26.67	26.25	26.26	27.19	26.59
	SD15-Hyper-1Step(Ren et al., 2024)	27.76	27.36	27.41	27.63	27.54
	SD15-Hyper-4Step(Ren et al., 2024)	28.04	27.39	27.42	27.89	27.69
	SD15-INSTAFLOW-1STEP(LIU ET AL., 2023)	26.07	25.80	25.89	26.32	26.02
	SD15-PEREFLOW-ISTEP(YAN ET AL., 2024)	25.70	25.45	25.57	25.96	25.67
	SD15-BOOI-ISTEP(GUETAL., 2023) SD21 SWIETPRIJSH 1STEP(NCUVEN & TRAN. 2022)	25.29	24.40	24.01	25.10	24.80
	SD21- $SWIFIBRO3R-ISTEP(NOUTEN & TRAN, 2023)SD21-TURBO_1STEP(SAUEP ET AL 2023P)$	20.31	26.86	20.37	26.89	20.70
	$SDXL-DMD2^{\ddagger}-1STEP-1024(YIN ET AL. 2024)$	27.40 27.67	20.00 27.02	27.40	26.03 26.94	27.11
	SDXL-DMD2 [‡] -4STEP-1024(YIN ET AL., 2024)	28.97	27.99	27.90	28.28	28.29
	SDXL-DMD2 [‡] -1Step-512(Yin et al., 2024)	27.70	27.07	27.02	26.94	27.18
	SDXL-DMD2 [‡] -4STEP-512(YIN ET AL., 2024)	27.22	26.65	26.62	26.57	26.76
	SD15-DMD2-1STEP-512(YIN ET AL., 2024)	26.31	25.75	25.78	26.59	26.11
	SD15-15Step(Rombach et al., 2022)	26.76	26.37	26.41	27.12	26.66
	SD15-25STEP(ROMBACH ET AL., 2022)	27.04	26.57	26.61	27.30	26.88
	SDXL-BASE-15STEP(PODELL ET AL., 2023)	28.25	27.27	27.43	27.43	27.60
	SD15-SIDLSG-1STEP(REPORT)(ZHOU ET AL., 2024A)	27.39	26.65	26.58	27.30	26.98
	SD15-SIDLSG-1STEP(ZHOU ET AL., 2024A)	26.37	25.85	25.88	26.73	26.20
	PIXELART- α -25STEP-512(CHEN ET AL., 2023)	28.77	27.92	27.96	28.37	28.25
	PIXELART- α -15STEP-512(CHEN ET AL., 2023)	28.68	27.85	27.87	28.29	28.17
	SD15-DI++-1STEP(LUO, 2024)'	28.42	27.84	28.01	28.19	28.12
	SDAL-DMD2-13TEP-1024 (TIN ET AL., 2024)($\alpha_c = 6$) SD15-DI*-1STEP ($\alpha_c = 1000, \alpha_c = 1.5$)(OURS)	28.40	27.52	27.52	21.10	27.81
	SDXL-DI*-1STEP.1024 ($\alpha = 1000, \alpha_c = 1.0$)(00RS)	28.50 28.74	28.05	28.17	28.51	28.21
	DIT-DI*-1STEP ($\alpha_r = 10, \alpha_r = 4.5$)(OURS)	28.78	28.31	28.48	28.37	28.48
	DIT-DI*-1STEP ($\alpha_r = 1, \alpha_c = 4.5$)(OURS)	29.13	28.51	28.51	28.63	28.70
		00.10	22.21	00 =0	04.75	
	SD15-15STEP(KOMBACH ET AL., 2022)	23.43	22.91	22.76	24.17	23.32
	SUAL-BASE-ISSTEP(PODELL ET AL., 2023)	29.71	27.69	27.71	25.46	27.64
	FIALAKI- α -1381EP(UHEN ET AL., 2023) SD15 SIDI SG 1STEP(ZHOU ET AL. 2024A)	31.31 22 54	29.80 21 52	29.00 21.27	28.49 23.00	29.81 22.12
	SD13-SIDLSG-15TEP(ZHOU ET AL., $2024A$) SDXL-DMD2-1STEP(YIN ET AL. 2024)	22.04 29.72	21.00 27.06	21.37 27.64	⊿3.09 26.55	22.13 27.07
	$\mathbf{RCM} = \mathbf{IR} + I$	29.65	31.15	32.00	31.03	30.95
	SD15-DI*-1STEP(OURS)	29.00	28.88	29.17	27.68	28.68
	SDXL-DI*-1STEP-1024(OURS)	30.74	30.03	30.05	28.00	29.71
	DIT-DI*-1STEP(OURS)	33.05	32.37	32.10	30.07	31.90

540 **BROADER IMPACT STATEMENT** 541

542 This work is motivated by our aim to increase the positive impact of one-step text-to-image genera-543 tive models toward satisfying human preferences. By default, one-step generators are either trained 544 over large-scale image-caption pair datasets or distilled from pre-trained diffusion models, which convey only subjective knowledge without human instructions.

546 Our results indicate that the proposed approach is promising for making one-step generative models 547 more aesthetic, and more preferred by human users. In the longer term, alignment failures could 548 lead to more severe consequences, particularly if these models are deployed in safety-critical sit-549 uations. For instance, if alignment failures occur, the one-step text-to-image model may generate 550 toxic images with misleading information, and horrible images that can potentially be scary to users. 551 We strongly recommend using our human preference alignment techniques together with AI safety 552 checkers for text-to-image generation to prevent undesirable negative impacts.

553 554

555 556

558 559

560

563

564

565

576

577

582

Reproducibility Statement

We provide extensive details of experimental settings and hyperparameters to reproduce our experimental results. We plan to release our code to ensure transparency and reproducibility of the results.

- REFERENCES
- 561 Christopher M Bishop. Pattern recognition and machine learning. Springer google schola, 2:1122– 562 1128, 2006.
 - Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion models with reinforcement learning. arXiv preprint arXiv:2305.13301, 2023.
- Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe 566 Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video 567 generation models as world simulators. 2024. URL https://openai.com/research/ 568 video-generation-models-as-world-simulators. 569
- 570 Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, 571 James T. Kwok, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart- α : Fast training of diffu-572 sion transformer for photorealistic text-to-image synthesis. ArXiv, abs/2310.00426, 2023. URL 573 https://api.semanticscholar.org/CorpusID:263334265.
- 574 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep 575 reinforcement learning from human preferences. Advances in neural information processing systems, 30, 2017.
- Kevin Clark, Paul Vicol, Kevin Swersky, and David J Fleet. Directly fine-tuning diffusion models 578 on differentiable rewards. arXiv preprint arXiv:2309.17400, 2023. 579
- 580 Guillaume Couairon, Jakob Verbeek, Holger Schwenk, and Matthieu Cord. Diffedit: Diffusion-581 based semantic image editing with mask guidance. ArXiv, abs/2210.11427, 2022.
- Xiaoliang Dai, Ji Hou, Chih-Yao Ma, Sam Tsai, Jialiang Wang, Rui Wang, Peizhao Zhang, Simon 583 Vandenhende, Xiaofang Wang, Abhimanyu Dubey, et al. Emu: Enhancing image generation 584 models using photogenic needles in a haystack. arXiv preprint arXiv:2309.15807, 2023. 585
- 586 Wei Deng, Weijian Luo, Yixin Tan, Marin Biloš, Yu Chen, Yuriy Nevmyvaka, and Ricky TQ Chen. 587 Variational schr\" odinger diffusion models. arXiv preprint arXiv:2405.04795, 2024.
- 588 Ming Ding, Wendi Zheng, Wenyi Hong, and Jie Tang. Cogview2: Faster and better text-to-image 589 generation via hierarchical transformers. Advances in Neural Information Processing Systems, 590 35:16890–16902, 2022.
- Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image 592 synthesis. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 12873-12883, 2021.

- Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first International Conference on Machine Learning*, 2024.
- Ying Fan, Olivia Watkins, Yuqing Du, Hao Liu, Moonkyung Ryu, Craig Boutilier, Pieter Abbeel, Mohammad Ghavamzadeh, Kangwook Lee, and Kimin Lee. Reinforcement learning for finetuning text-to-image diffusion models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Yasong Feng, Weijian Luo, Yimin Huang, and Tianyu Wang. A lipschitz bandits approach for continuous hyperparameter optimization. *arXiv preprint arXiv:2302.01539*, 2023.
- ⁶⁰⁵
 ⁶⁰⁶
 ⁶⁰⁶
 ⁶⁰⁶
 ⁶⁰⁷
 ⁶⁰⁷
 ⁶⁰⁸
 ⁶⁰⁸
 ⁶⁰⁷
 ⁶⁰⁸
 ⁶⁰⁹
 ⁶⁰⁹
- Zhengyang Geng, Ashwini Pokle, William Luo, Justin Lin, and J Zico Kolter. Consistency models
 made easy. *arXiv preprint arXiv:2406.14548*, 2024.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural infor- mation processing systems*, pp. 2672–2680, 2014.
- Jiatao Gu, Shuangfei Zhai, Yizhe Zhang, Lingjie Liu, and Josh Susskind. Boot: Data-free distillation of denoising diffusion models with bootstrapping. *arXiv preprint arXiv:2306.05544*, 2023.
- Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and
 Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10696–10706, 2022.
- Jonathan Heek, Emiel Hoogeboom, and Tim Salimans. Multistep consistency models. arXiv preprint arXiv:2403.06807, 2024.
- Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. arXiv preprint arXiv:2207.12598, 2022.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. Advances in
 Neural Information Processing Systems, 33:6840–6851, 2020.
- Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *arXiv preprint arXiv:2204.03458*, 2022.
- Minguk Kang, Jun-Yan Zhu, Richard Zhang, Jaesik Park, Eli Shechtman, Sylvain Paris, and Taesung
 Park. Scaling up gans for text-to-image synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023a.
- Minguk Kang, Jun-Yan Zhu, Richard Zhang, Jaesik Park, Eli Shechtman, Sylvain Paris, and Taesung Park. Scaling up gans for text-to-image synthesis. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pp. 10124–10134, 2023b.
- Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyz ing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 8110–8119, 2020.
- Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion based generative models. In *Proc. NeurIPS*, 2022.
- Dongjun Kim, Chieh-Hsin Lai, Wei-Hsiang Liao, Naoki Murata, Yuhta Takida, Toshimitsu Uesaka, Yutong He, Yuki Mitsufuji, and Stefano Ermon. Consistency trajectory models: Learning probability flow ode trajectory of diffusion. *arXiv preprint arXiv:2310.02279*, 2023.
- Heeseung Kim, Sungwon Kim, and Sungroh Yoon. Guided-tts: A diffusion model for text-to-speech
 via classifier guidance. In *International Conference on Machine Learning*, pp. 11119–11133.
 PMLR, 2022.

658

659

660

685

687 688

- 648
 649
 650
 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete
 Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceed- ings of the IEEE/CVF International Conference on Computer Vision*, pp. 4015–4026, 2023.
- Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy.
 Pick-a-pic: An open dataset of user preferences for text-to-image generation. *arXiv preprint arXiv:2305.01569*, 2023.
 - Kimin Lee, Hao Liu, Moonkyung Ryu, Olivia Watkins, Yuqing Du, Craig Boutilier, Pieter Abbeel, Mohammad Ghavamzadeh, and Shixiang Shane Gu. Aligning text-to-image models using human feedback. arXiv preprint arXiv:2302.12192, 2023.
- Jiachen Li, Weixi Feng, Wenhu Chen, and William Yang Wang. Reward guided latent consistency
 distillation. *arXiv preprint arXiv:2403.11027*, 2024.
- Shanchuan Lin, Anran Wang, and Xiao Yang. Sdxl-lightning: Progressive adversarial diffusion distillation. *arXiv preprint arXiv:2402.13929*, 2024.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays,
 Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft coco:
 Common objects in context, 2014. URL http://arxiv.org/abs/1405.0312. cite
 arxiv:1405.0312Comment: 1) updated annotation pipeline description and figures; 2) added new
 section describing datasets splits; 3) updated author list.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. Advances in neural information processing systems, 36, 2024.
- Kingchao Liu, Chengyue Gong, Lemeng Wu, Shujian Zhang, Hao Su, and Qiang Liu. Fuse dream: Training-free text-to-image generation with improved clip+ gan space optimization. *arXiv* preprint arXiv:2112.01573, 2021.
- Kingchao Liu, Xiwen Zhang, Jianzhu Ma, Jian Peng, et al. Instaflow: One step is enough for
 high-quality diffusion-based text-to-image generation. In *The Twelfth International Conference on Learning Representations*, 2023.
- Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models: Synthesizing high-resolution images with few-step inference. *arXiv preprint arXiv:2310.04378*, 2023a.
- Weijian Luo. A comprehensive survey on knowledge distillation of diffusion models. *arXiv preprint arXiv:2304.04262*, 2023.
- Weijian Luo. Diff-instruct++: Training one-step text-to-image generator model to align with human preferences. *arXiv preprint arXiv:2410.18881*, 2024.
 - Weijian Luo and Zhihua Zhang. Data prediction denoising models: The pupil outdoes the master, 2024. URL https://openreview.net/forum?id=wYmcfur889.
- Weijian Luo, Hao Jiang, Tianyang Hu, Jiacheng Sun, Zhenguo Li, and Zhihua Zhang. Training energy-based models with diffusion contrastive divergences. *arXiv preprint arXiv:2307.01668*, 2023b.
- Weijian Luo, Tianyang Hu, Shifeng Zhang, Jiacheng Sun, Zhenguo Li, and Zhihua Zhang. Diff instruct: A universal approach for transferring knowledge from pre-trained diffusion models.
 Advances in Neural Information Processing Systems, 36, 2024a.
- Weijian Luo, Tianyang Hu, Shifeng Zhang, Jiacheng Sun, Zhenguo Li, and Zhihua Zhang. Diff instruct: A universal approach for transferring knowledge from pre-trained diffusion models.
 Advances in Neural Information Processing Systems, 36, 2024b.
- 701 Weijian Luo, Zemin Huang, Zhengyang Geng, J Zico Kolter, and Guo-Jun Qi. One-step diffusion distillation through score implicit matching. *arXiv preprint arXiv:2410.16794*, 2024c.

- Weijian Luo, Boya Zhang, and Zhihua Zhang. Entropy-based training methods for scalable neural implicit samplers. *Advances in Neural Information Processing Systems*, 36, 2024d.
- Chenlin Meng, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021.
- Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications* of the ACM, 65(1):99–106, 2021.
- Thuan Hoang Nguyen and Anh Tran. Swiftbrush: One-step text-to-image diffusion model with variational score distillation. *arXiv preprint arXiv:2312.05239*, 2023.
- Alex Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. *arXiv preprint arXiv:2102.09672*, 2021.
- Alex 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. *arXiv preprint arXiv:2112.10741*, 2021.
- Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves,
 Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative model for
 raw audio. *arXiv preprint arXiv:1609.03499*, 2016.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744, 2022.
- William Peebles and Saining Xie. Scalable diffusion models with transformers. arXiv preprint arXiv:2212.09748, 2022.
- Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe
 Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023.
- Ashwini Pokle, Zhengyang Geng, and J Zico Kolter. Deep equilibrium approaches to diffusion models. *Advances in Neural Information Processing Systems*, 35:37975–37990, 2022.
- Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022.
- Mihir Prabhudesai, Anirudh Goyal, Deepak Pathak, and Katerina Fragkiadaki. Aligning text-toimage diffusion models with reward backpropagation. *arXiv preprint arXiv:2310.03739*, 2023.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea
 Finn. Direct preference optimization: Your language model is secretly a reward model. Advances
 in Neural Information Processing Systems, 36, 2024.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi
 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text
 transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical textconditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- Yuxi Ren, Xin Xia, Yanzuo Lu, Jiacheng Zhang, Jie Wu, Pan Xie, Xing Wang, and Xuefeng Xiao.
 Hyper-sd: Trajectory segmented consistency model for efficient image synthesis. *arXiv preprint arXiv:2404.13686*, 2024.

756 757 758 750	Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High- resolution image synthesis with latent diffusion models. In <i>Proceedings of the IEEE/CVF Con-</i> <i>ference on Computer Vision and Pattern Recognition</i> , pp. 10684–10695, 2022.
759 760 761	Tim Salimans, Thomas Mensink, Jonathan Heek, and Emiel Hoogeboom. Multistep distillation of diffusion models via moment matching. <i>arXiv preprint arXiv:2406.04103</i> , 2024.
762 763 764 765	Axel Sauer, Tero Karras, Samuli Laine, Andreas Geiger, and Timo Aila. Stylegan-t: Unlocking the power of gans for fast large-scale text-to-image synthesis. In <i>International conference on machine learning</i> , pp. 30105–30118. PMLR, 2023a.
766 767	Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin Rombach. Adversarial diffusion dis- tillation. <i>arXiv preprint arXiv:2311.17042</i> , 2023b.
768 769 770	ChristophSchuhmann.Laion-aesthetics.https://laion.ai/blog/laion-aesthetics/, 2022.Accessed: 2023 - 11- 10.
771	John Schulman. Trust region policy optimization. arXiv preprint arXiv:1502.05477, 2015.
772 773 774	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. <i>arXiv preprint arXiv:1707.06347</i> , 2017.
775 776 777 778	Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In <i>International Conference on Machine Learning</i> , pp. 2256–2265. PMLR, 2015.
779 780	Yang Song and Prafulla Dhariwal. Improved techniques for training consistency models. In <i>The Twelfth International Conference on Learning Representations</i> .
781 782 783 784	Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In <i>International Conference on Learning Representations</i> , 2020.
785 786	Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. <i>arXiv preprint arXiv:2303.01469</i> , 2023.
787 788 789	Yuda Song, Zehao Sun, and Xuanwu Yin. Sdxs: Real-time one-step latent diffusion models with image conditions. <i>arXiv preprint arXiv:2403.16627</i> , 2024.
790 791 792	Yusuke Tashiro, Jiaming Song, Yang Song, and Stefano Ermon. Csdi: Conditional score-based diffusion models for probabilistic time series imputation. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> , 2021.
793 794 795 796	Yi Tay, Mostafa Dehghani, Jinfeng Rao, William Fedus, Samira Abnar, Hyung Won Chung, Sharan Narang, Dani Yogatama, Ashish Vaswani, and Donald Metzler. Scale efficiently: Insights from pre-training and fine-tuning transformers. <i>arXiv preprint arXiv:2109.10686</i> , 2021.
797 798 799	Pascal Vincent. A Connection Between Score Matching and Denoising Autoencoders. <i>Neural Computation</i> , 23(7):1661–1674, 2011.
800 801 802 803	Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou, Aaron Lou, Senthil Purushwalkam, Stefano Ermon, Caiming Xiong, Shafiq Joty, and Nikhil Naik. Diffusion model alignment using direct preference optimization. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 8228–8238, 2024.
804 805 806 807	Yifei Wang, Weimin Bai, Weijian Luo, Wenzheng Chen, and He Sun. Integrating amortized infer- ence with diffusion models for learning clean distribution from corrupted images. <i>arXiv preprint</i> <i>arXiv:2407.11162</i> , 2024.
808 809	Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolific- dreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. <i>arXiv</i> preprint arXiv:2305.16213, 2023.

827

837

838

839

840

847

- Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li.
 Human preference score v2: A solid benchmark for evaluating human preferences of text-toimage synthesis. *arXiv preprint arXiv:2306.09341*, 2023.
- Sirui Xie, Zhisheng Xiao, Diederik P Kingma, Tingbo Hou, Ying Nian Wu, Kevin Patrick Murphy, Tim Salimans, Ben Poole, and Ruiqi Gao. Em distillation for one-step diffusion models. *arXiv* preprint arXiv:2405.16852, 2024.
- Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao
 Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation, 2023a.
- Xingqian Xu, Zhangyang Wang, Gong Zhang, Kai Wang, and Humphrey Shi. Versatile diffusion:
 Text, images and variations all in one diffusion model. In *Proceedings of the IEEE/CVF Interna- tional Conference on Computer Vision*, pp. 7754–7765, 2023b.
- 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. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https://openreview.
 net/forum?id=f6a9XVFYIO.
- Hanshu Yan, Xingchao Liu, Jiachun Pan, Jun Hao Liew, Qiang Liu, and Jiashi Feng. Perflow:
 Piecewise rectified flow as universal plug-and-play accelerator. *arXiv preprint arXiv:2405.07510*, 2024.
- Kai Yang, Jian Tao, Jiafei Lyu, Chunjiang Ge, Jiaxin Chen, Weihan Shen, Xiaolong Zhu, and Xiu Li.
 Using human feedback to fine-tune diffusion models without any reward model. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 8941–8951, 2024.
- Tianwei Yin, Michaël Gharbi, Richard Zhang, Eli Shechtman, Fredo Durand, William T Freeman, and Taesung Park. One-step diffusion with distribution matching distillation. *arXiv preprint arXiv:2311.18828*, 2023.
 - Tianwei Yin, Michaël Gharbi, Taesung Park, Richard Zhang, Eli Shechtman, Fredo Durand, and William T Freeman. Improved distribution matching distillation for fast image synthesis. *arXiv* preprint arXiv:2405.14867, 2024.
- Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan,
 Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2(3):5, 2022.
- Boya Zhang, Weijian Luo, and Zhihua Zhang. Enhancing adversarial robustness via score-based optimization. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023a.
 URL https://openreview.net/forum?id=MOAHXRzHhm.
 - Boya Zhang, Weijian Luo, and Zhihua Zhang. Purify++: Improving diffusion-purification with advanced diffusion models and control of randomness. *arXiv preprint arXiv:2310.18762*, 2023b.
- Bowen Zheng and Tianming Yang. Diffusion models are innate one-step generators. *arXiv preprint arXiv:2405.20750*, 2024.
- Jianbin Zheng, Minghui Hu, Zhongyi Fan, Chaoyue Wang, Changxing Ding, Dacheng Tao, and Tat-Jen Cham. Trajectory consistency distillation. *arXiv preprint arXiv:2402.19159*, 2024.
- Mingyuan Zhou, Zhendong Wang, Huangjie Zheng, and Hai Huang. Long and short guidance in score identity distillation for one-step text-to-image generation. *arXiv preprint arXiv:2406.01561*, 2024a.
- Mingyuan Zhou, Huangjie Zheng, Zhendong Wang, Mingzhang Yin, and Hai Huang. Score identity distillation: Exponentially fast distillation of pretrained diffusion models for one-step generation. *arXiv preprint arXiv:2404.04057*, 2024b.
- Yufan Zhou, Ruiyi Zhang, Changyou Chen, Chunyuan Li, Chris Tensmeyer, Tong Yu, Jiuxiang
 Gu, Jinhui Xu, and Tong Sun. Towards language-free training for text-to-image generation. In
 Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 17907–17917, 2022.

THEORY А

A.1 **PROOF OF THEOREM 3.1**

Proof. Recall that $p_{\theta}(\cdot)$ is induced by the generator $g_{\theta}(\cdot)$, therefore the sample is obtained by $x_0 =$ $g_{\theta}(\boldsymbol{z}|\boldsymbol{c}), \boldsymbol{z} \sim p_{z}$. The term \boldsymbol{x} contains parameter through $\boldsymbol{x}_{0} = g_{\theta}(\boldsymbol{z}|\boldsymbol{c}), \boldsymbol{z} \sim p_{z}$. To demonstrate the parameter dependence, we use the notation $p_{\theta}(\cdot)$. Note that $p_{ref}(\cdot)$ is the reference distribution. The alignment objective writes

$$\mathcal{L}_{Orig}(\theta) = \mathbb{E}_{\substack{\boldsymbol{x} \sim p_{z}, \\ \boldsymbol{x}_{0} = g_{\theta}(\boldsymbol{x}|\boldsymbol{c})}} \left[-\alpha r(\boldsymbol{x}_{0}, \boldsymbol{c}) \right] + \mathbf{D}^{[0,T]}(p_{\theta}, p_{ref})$$
(A.1)

The first loss term of (A.1) $\alpha r(\boldsymbol{x}_0, \boldsymbol{c})$ is easy to compute by directly pushing the generated sample x_0 and the text prompt c into the reward model $r(\cdot, \cdot)$. However, the second loss term (A.2) is intractable because we do not explicitly know the relation between θ and $p_{\theta,t}(\cdot)$.

$$\mathbf{D}^{[0,T]}(p_{\theta}, p_{ref}) \coloneqq \int_{t=0}^{T} w(t) \mathbb{E}_{\boldsymbol{x}_{t} \sim \pi_{t}} \left\{ \mathbf{d}(\boldsymbol{s}_{p_{\theta,t}}(\boldsymbol{x}_{t}) - \boldsymbol{s}_{q_{t}}(\boldsymbol{x}_{t})) \right\} \mathrm{d}t,$$
(A.2)

We turn to derive the equivalent loss for $\mathbf{D}^{[0,T]}(p_{\theta}, p_{ref})$. First we take the θ gradient of (A.2), show

$$\frac{\partial}{\partial \theta} \mathbf{D}^{[0,T]}(p_{\theta}, p_{ref}) = \frac{\partial}{\partial \theta} \int_{t=0}^{T} w(t) \mathbb{E}_{\boldsymbol{x}_{t} \sim \pi_{t}} \left\{ \mathbf{d}(\boldsymbol{s}_{p_{\theta,t}}(\boldsymbol{x}_{t}) - \boldsymbol{s}_{q_{t}}(\boldsymbol{x}_{t})) \right\} dt$$
(A.3)

$$= \mathbb{E}_{t, \boldsymbol{x}_t \sim \pi_t} w(t) \left\{ \mathbf{d}'(\boldsymbol{y}_t) \right\}^T \frac{\partial}{\partial \theta} \boldsymbol{s}_{p_{\theta, t}}(\boldsymbol{x}_t)$$
(A.4)

Notice that $p_{\theta,t}(\cdot)$ is induced by first generating samples with one-step generator then adding noise with diffusion process (2.1), we do not know the term $\frac{\partial}{\partial \theta} s_{p_{\theta,t}}(x_t)$. Therefore the gradient formula (A.4) is intractable. However, we will show that a tractable loss function can recover the intractable gradient (A.4), and therefore can be used for minimizing (A.2). Our proof is inspired by the theory from Vincent (2011), Zhou et al. (2024b) and Luo et al. (2024c).

We first present a so-called Score-projection identity (Theorem A.1), which has been studied in Zhou et al. (2024b) and Vincent (2011):

Theorem A.1. Let $u(\cdot)$ be a θ -free vector-valued function under mild conditions, the identity holds:

$$\mathbb{E}_{\substack{\boldsymbol{x}_{0} \sim p_{\theta,0}, \\ \boldsymbol{x}_{t} \mid \boldsymbol{x}_{0} \sim q_{t}(\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0})}} \boldsymbol{u}(\boldsymbol{x}_{t})^{T} \left\{ \boldsymbol{s}_{p_{\theta,t}}(\boldsymbol{x}_{t}) - \nabla_{\boldsymbol{x}_{t}} \log q_{t}(\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0}) \right\} = 0, \quad \forall \theta.$$
(A.5)

We give a short proof of Theorem A.1 as a clarification. Readers can also refer to Vincent (2011) or Zhou et al. (2024b) as a reference.

Recall the relation between $s_{p_{\theta,t}}(x_t)$ and $\nabla_{x_t} \log q_t(x_t|x_0)$, we know

$$oldsymbol{s}_{p_{oldsymbol{ heta},t}}(oldsymbol{x}_t) =
abla_{oldsymbol{x}_t}\log\int p_{oldsymbol{ heta},0}(oldsymbol{x}_0)q_t(oldsymbol{x}_t|oldsymbol{x}_0)\mathrm{d}oldsymbol{x}_0$$

913
914
$$\int p_{\theta,t}(\boldsymbol{x}_0) \nabla_{\boldsymbol{x}_t} q_t(\boldsymbol{x}_t | \boldsymbol{x}_0) d\boldsymbol{x}_0$$

$$=\frac{\int p_{\theta,t}(\boldsymbol{x}_0) \, \mathbf{v}_{\boldsymbol{x}_t} q_t(\boldsymbol{x}_t | \boldsymbol{x}_0)}{\mathbf{v}_{\boldsymbol{x}_t} q_t(\boldsymbol{x}_t | \boldsymbol{x}_0)}$$

$$p_{\theta}$$

915
916
917
$$= \int \frac{p_{\theta,t}(\boldsymbol{x}_t)}{\nabla_{\boldsymbol{x}_t} \log q_t(\boldsymbol{x}_t | \boldsymbol{x}_0) p_{\theta,t}(\boldsymbol{x}_0) q_t(\boldsymbol{x}_t | \boldsymbol{x}_0)}{p_{\theta,t}(\boldsymbol{x}_t)} d\boldsymbol{x}_0$$

We have $\mathbb{E}_{\substack{\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0} \sim p_{\theta,0}, \\ \boldsymbol{x}_{t} \mid \boldsymbol{x}_{0} \sim q_{t}(\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0})}} \boldsymbol{u}(\boldsymbol{x}_{t})^{T} \boldsymbol{s}_{p_{\theta,t}}(\boldsymbol{x}_{t}) = \mathbb{E}_{\boldsymbol{x}_{t} \sim p_{\theta,t}} \boldsymbol{u}(\boldsymbol{x}_{t})^{T} \boldsymbol{s}_{p_{\theta,t}}(\boldsymbol{x}_{t})$ $=\int p_{\theta,t}(\boldsymbol{x}_t)\boldsymbol{u}(\boldsymbol{x}_t)^T\boldsymbol{s}_{p_{\theta,t}}(\boldsymbol{x}_t)\mathrm{d}\boldsymbol{x}_t$ $= \int p_{\theta,t}(\boldsymbol{x}_t) \boldsymbol{u}(\boldsymbol{x}_t)^T \int \frac{\nabla_{\boldsymbol{x}_t} \log q_t(\boldsymbol{x}_t | \boldsymbol{x}_0) p_{\theta,t}(\boldsymbol{x}_0) q_t(\boldsymbol{x}_t | \boldsymbol{x}_0)}{p_{\theta,t}(\boldsymbol{x}_t)} \mathrm{d}\boldsymbol{x}_0 \mathrm{d}\boldsymbol{x}_t$ $= \int \boldsymbol{u}(\boldsymbol{x}_t)^T \int \nabla_{\boldsymbol{x}_t} \log q_t(\boldsymbol{x}_t | \boldsymbol{x}_0) p_{\theta,t}(\boldsymbol{x}_0) q_t(\boldsymbol{x}_t | \boldsymbol{x}_0) \mathrm{d} \boldsymbol{x}_0 \mathrm{d} \boldsymbol{x}_t$ $=\int\int egin{aligned} & =\int\int egin{aligned} & \mathbf{u}(oldsymbol{x}_t)^T
abla_{oldsymbol{x}_t}\log q_t(oldsymbol{x}_t|oldsymbol{x}_0)p_{ heta,t}(oldsymbol{x}_0)q_t(oldsymbol{x}_t|oldsymbol{x}_0)\mathrm{d}oldsymbol{x}_0\mathrm{d}oldsymbol{x}_t \end{aligned}$ $= \mathbb{E}_{\boldsymbol{x}_t \mid \boldsymbol{x}_0 \sim p_{\boldsymbol{\theta}, 0}, \atop \boldsymbol{x}_t \mid \boldsymbol{x}_0 \sim q_t(\boldsymbol{x}_t \mid \boldsymbol{x}_0)} \boldsymbol{u}(\boldsymbol{x}_t)^T \nabla_{\boldsymbol{x}_t} \log q_t(\boldsymbol{x}_t \mid \boldsymbol{x}_0)$

 If we take the θ gradient on both sides of (A.5), we have

$$0 = \mathbb{E}_{\substack{\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0} \sim p_{\theta,0}, \\ \boldsymbol{x}_{t} \mid \boldsymbol{x}_{0} \sim q_{t}(\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0})}} \left\{ \frac{\partial}{\partial \boldsymbol{x}_{t}} \left[\boldsymbol{u}(\boldsymbol{x}_{t})^{T} \left\{ \boldsymbol{s}_{p_{\theta,t}}(\boldsymbol{x}_{t}) - \nabla_{\boldsymbol{x}_{t}} \log q_{t}(\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0}) \right\} \right] \frac{\partial \boldsymbol{x}_{t}}{\partial \theta} - \boldsymbol{u}(\boldsymbol{x}_{t})^{T} \frac{\partial}{\partial \boldsymbol{x}_{0}} \left[\nabla_{\boldsymbol{x}_{t}} \log q_{t}(\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0}) \right] \frac{\partial \boldsymbol{x}_{0}}{\partial \theta} \right\} + \mathbb{E}_{\boldsymbol{x}_{t} \sim p_{\theta,t}} \boldsymbol{u}(\boldsymbol{x}_{t})^{T} \frac{\partial}{\partial \theta} \left\{ \boldsymbol{s}_{p_{\theta,t}}(\boldsymbol{x}_{t}) \right\}$$
(A.6)

So we have an identity

$$\mathbb{E}_{\boldsymbol{x}_t \sim p_{\theta,t}} \boldsymbol{u}(\boldsymbol{x}_t)^T \frac{\partial}{\partial \theta} \left\{ \boldsymbol{s}_{p_{\theta,t}}(\boldsymbol{x}_t) \right\} = -\frac{\partial}{\partial \theta} \mathbb{E}_{\boldsymbol{x}_t \mid \boldsymbol{x}_0 \sim p_{\theta,0}, \atop \boldsymbol{x}_t \mid \boldsymbol{x}_0 \sim q_t(\boldsymbol{x}_t \mid \boldsymbol{x}_0)} \left\{ \boldsymbol{u}(\boldsymbol{x}_t) \left\{ \boldsymbol{s}_{p_{\mathrm{sg}[\theta],t}}(\boldsymbol{x}_t) - \nabla_{\boldsymbol{x}_t} \log q_t(\boldsymbol{x}_t \mid \boldsymbol{x}_0) \right\} \right\}$$

Notice that the left-hand side of equation (A.7) can be interpreted as the gradient of the loss function when the parameter dependency of the sampling distribution is cut off, i.e.

$$\mathbb{E}_{\boldsymbol{x}_t \sim p_{\theta,t}} \boldsymbol{u}(\boldsymbol{x}_t)^T \frac{\partial}{\partial \theta} \left\{ \boldsymbol{s}_{p_{\theta,t}}(\boldsymbol{x}_t) \right\} = \frac{\partial}{\partial \theta} \mathbb{E}_{\boldsymbol{x}_t \sim p_{\mathrm{sg}[\theta],t}} \left\{ \boldsymbol{u}(\boldsymbol{x}_t)^T \boldsymbol{s}_{p_{\theta,t}}(\boldsymbol{x}_t) \right\}$$
(A.7)

Therefore we have the final equation

$$\frac{\partial}{\partial \theta} \mathbb{E}_{\boldsymbol{x}_{t} \sim p_{\mathrm{sg}[\theta],t}} \left\{ \boldsymbol{u}(\boldsymbol{x}_{t})^{T} \boldsymbol{s}_{p_{\theta,t}}(\boldsymbol{x}_{t}) \right\} = -\frac{\partial}{\partial \theta} \mathbb{E}_{\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0} \sim q_{t}(\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0})} \left\{ \boldsymbol{u}(\boldsymbol{x}_{t}) \left\{ \boldsymbol{s}_{p_{\mathrm{sg}[\theta],t}}(\boldsymbol{x}_{t}) - \nabla_{\boldsymbol{x}_{t}} \log q_{t}(\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0}) \right\} \right\}$$
(A.8)

which holds for arbitrary function $u(\cdot)$ and parameter θ . If we set

$$egin{aligned} oldsymbol{u}(oldsymbol{x}_t) &= \mathbf{d}'(oldsymbol{y}_t) \ oldsymbol{y}_t &= oldsymbol{s}_{p_{ ext{sg}[oldsymbol{ heta}],t}}(oldsymbol{x}_t) - oldsymbol{s}_{q_t}(oldsymbol{x}_t) \end{aligned}$$

m

Then we formally have

$$\frac{\partial}{\partial \theta} \mathbb{E}_{\boldsymbol{x}_{t} \sim p_{\mathrm{sg}[\theta], t}} \left\{ \mathbf{d}'(\boldsymbol{y}_{t}) \right\}^{T} \left\{ \boldsymbol{s}_{p_{\theta, t}}(\boldsymbol{x}_{t}) \right\}$$

$$= \frac{\partial}{\partial \theta} \mathbb{E}_{\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0} \sim q_{t}(\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0})}^{\boldsymbol{x}_{0} \sim p_{\theta, 0}, 0} \left\{ -\mathbf{d}'(\boldsymbol{y}_{t}) \right\}^{T} \left\{ \boldsymbol{s}_{p_{\theta, t}}(\boldsymbol{x}_{t}) - \nabla_{\boldsymbol{x}_{t}} \log q_{t}(\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0}) \right\}$$
(A.9)

This means that we can use the θ gradient of a tractable loss:

$$\mathbb{E}_{\substack{t, \boldsymbol{x}_{0} \sim p_{\theta,0}, \\ \boldsymbol{x}_{t} \mid \boldsymbol{x}_{0} \sim q_{t}(\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0})}} w(t) \left\{ -\mathbf{d}'(\boldsymbol{y}_{t}) \right\}^{T} \left\{ \boldsymbol{s}_{p_{\theta,t}}(\boldsymbol{x}_{t}) - \nabla_{\boldsymbol{x}_{t}} \log q_{t}(\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0}) \right\}$$
(A.10)

to replace the wanted θ gradient (A.3), which can minimize the regularization loss (A.2).

Combining $r(\boldsymbol{x}_0, \boldsymbol{c})$ and (A.10), we have the practical loss

$$\mathcal{L}_{DI*}(\theta) = \mathbb{E}_{\substack{\boldsymbol{z} \sim p_{z}, \\ \boldsymbol{x}_{0} = g_{\theta}(\boldsymbol{z})}} \left[-\alpha r(\boldsymbol{x}_{0}, \boldsymbol{c}) \right]$$
(A.11)

969
970
971

$$+ \mathbb{E}_{\substack{t, \boldsymbol{x}_t \mid \boldsymbol{x}_0 \\ \sim q_t(\boldsymbol{x}_t \mid \boldsymbol{x}_0)}} w(t) \left\{ -\mathbf{d}'(\boldsymbol{y}_t) \right\}^T \left\{ \boldsymbol{s}_{p_{\mathrm{sg}[\theta], t}}(\boldsymbol{x}_t) - \nabla_{\boldsymbol{x}_t} \log q_t(\boldsymbol{x}_t \mid \boldsymbol{x}_0) \right\} \mathrm{d}t \right]$$
971

972
 973
 974
 Remark A.2. In practice, most commonly used forward diffusion processes can be expressed as a form of scale and noise addition:

$$\boldsymbol{x}_t = \alpha(t)\boldsymbol{x}_0 + \beta(t)\boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\boldsymbol{\epsilon}; \mathbf{0}, \mathbf{I}).$$
 (A.12)

So the term x_t in equation (A.11) can be instantiated as $z \sim p_z$, $\epsilon \sim \mathcal{N}(\epsilon; \mathbf{0}, \mathbf{I})$, $x_t = \alpha(t)x_0 + \beta(t)\epsilon$.

A.2 PROOF OF THEOREM 3.2

Proof. Recall the definition of the classifier-free reward (3.9). The negative reward writes

$$r(\boldsymbol{x}_0, \boldsymbol{c}) = -\mathbb{E}_{t, \boldsymbol{x}_t \sim p_{\theta, t}} w(t) \log rac{p_{ref}(\boldsymbol{x}_t | t, \boldsymbol{c})}{p_{ref}(\boldsymbol{x}_t | t)}$$

This reward will put a higher reward on those samples that have higher class-conditional probability than unconditional probability, therefore encouraging class-conditional sampling. It is clear that

$$\frac{\partial}{\partial \theta} \left\{ -r(\boldsymbol{x}_{0}, \boldsymbol{c}) \right\} = -\mathbb{E}_{t, \boldsymbol{x}_{t} \sim p_{\theta, t}} w(t) \left\{ \nabla_{\boldsymbol{x}_{t}} \log p_{ref}(\boldsymbol{x}_{t} | t, \boldsymbol{c}) - \nabla_{\boldsymbol{x}_{t}} \log p_{ref}(\boldsymbol{x}_{t} | t) \right\} \frac{\partial \boldsymbol{x}_{t}}{\partial \theta} \\ = -\mathbb{E}_{t, \boldsymbol{x}_{t} \sim p_{\theta, t}} w(t) \left\{ \boldsymbol{s}_{ref}(\mathrm{sg}[\boldsymbol{x}_{t}] | t, \boldsymbol{c}) - \boldsymbol{s}_{ref}(\mathrm{sg}[\boldsymbol{x}_{t}] | t, \boldsymbol{\varnothing}) \right\} \frac{\partial \boldsymbol{x}_{t}}{\partial \theta}$$
(A.13)

Therefore, we can see that the equivalent loss

$$\mathcal{L}_{cfg}(\theta) = \mathbb{E}_{\substack{t, \boldsymbol{z} \sim p_{z}, \boldsymbol{x}_{0} = g_{\theta}(\boldsymbol{z}|\boldsymbol{c})\\\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0} \sim p(\boldsymbol{x}_{t}\mid \boldsymbol{x}_{0})}} w(t) \left\{ \boldsymbol{s}_{ref}(\mathrm{sg}[\boldsymbol{x}_{t}]|t, \boldsymbol{c}) - \boldsymbol{s}_{ref}(\mathrm{sg}[\boldsymbol{x}_{t}]|t, \boldsymbol{\varnothing}) \right\}^{T} \boldsymbol{x}_{t}$$
(A.14)

recovers the gradient formula (A.13).

B IMPORTANT MATERIALS FOR MAIN CONTENT

999 1000 1001

1004

1008

1009

1010

1011

1012 1013

1014

1015

1016

1017

1020

1021

1023 1024

1025

997 998

975

978 979

980

985

986 987

989 990 991

B.1 PROMPTS FOR FIGURE 1 AND FIGURE 2

- Prompts for Figure 1 (from upper left to bottom right):
 - A girl examining an ammonite fossil;
 - A squirrel driving a toy car;
 - A portrait of a statue of the Egyptian god Anubis wearing aviator goggles, white t-shirt and leather jacket. The city of Los Angeles is in the background;
 - A still image of a humanoid cat posing with a hat and jacket in a bar;
 - A photograph of the inside of a subway train. There are red pandas sitting on the seats. One of them is reading a newspaper. The window shows the jungle in the background;
 - A capybara made of voxels sitting in a field;
 - A teddy bear on a skateboard in times square;
 - A sloth in a go kart on a race track. The sloth is holding a banana in one hand. There is a banana peel on the track in the background;
 - A close-up photo of a wombat wearing a red backpack and raising both arms in the air;
 - A small cactus with a happy face in the Sahara desert;
 - Baker proudly displays her white dog cake in her kitchen;
 - A bowl with rice, broccoli and a purple relish;
 - An inlet filled with boats of all kinds;
 - A black cat sitting on top of the hood of a car;
 - A woman wearing a cowboy hat face to face with a horse.
 - Prompts for Figure 2. The prompts are listed from the up rows to the bottom row:

• Pirate ship sailing into a bioluminescence sea with a galaxy in the sky, epic, 4k, ultra;

- Digital 2D, Miyazaki's style, ultimate detailed, tiny finnest details, futuristic, sci-fi, magical dreamy landscape scenery, small cute girl living alone with plushified friendly big tanuki in the gigantism of wilderness, intricate round futuristic simple multilayered architecture, habitation cabin in the trees, dramatic soft lightning, rule of thirds, cinematic;
- 1030 1031 1032

1026

1027

1028

1029

1033

• saharian landscape at sunset, 4k ultra realism, BY Anton Gorlin, trending on artstation, sharp focus, studio photo, intricate details, highly detailed, by greg rutkowski.

1034 B.2 MEANINGS OF HYPER-PARAMETERS.

1036 Meanings of Hyper-parameters. As in Algorithm 1, the overall algorithms consist of two al-1037 ternative updating steps. The first step is to update ψ of the assistant diffusion model by fine-1038 tuning it with student-generated data. Therefore the assistant diffusion $s_{\psi}(x_t|t,c)$ can approximate 1039 the score function of student generator distribution. This step means that the assistant diffusion 1040 needs to communicate with the student to know the student's status. The second step updates the 1041 generator by minimizing the tractable loss (3.7) using SGD-based optimization algorithms such as 1042 Adam (Kingma & Ba, 2014). This step means that the teacher and the assistant diffusion discuss 1043 and incorporate the student's interests to instruct the student generator.

1044 As we can see in Algorithm 1 (as well as Algorithm 2). Each hyperparameter has its intuitive 1045 meaning. The reward scale parameter α_{rew} controls the strength of human preference alignment. 1046 The larger the α_{rew} is, the stronger the generator is aligned with human preferences. However, the 1047 drawback for a too large α_{rew} might be the loss of diversity and reality. Besides, we empirically 1048 find that larger α_{rew} leads to richer generation details and better generation layouts. But a very large 1049 α_{rew} results in unrealistic and painting-like images.

1050 The CFG reward scale controls the strength of using CFG rewards when training. We empirically 1051 find that the best CFG scale for Diff-Instruct* is the same as the best CFG scale for sampling from 1052 the reference diffusion model. However, α_{cfg} may conflicts with α_{rew} . In the Stable Diffusion 1053 1.5 experiment, we find that using a large CFG reward scale leads to worse human preferences. 1054 Therefore, the proper combination of $(\alpha_{rew}, \alpha_{cfg})$ asks for careful tuning.

1055 The diffusion model weighting $\lambda(t)$ and the generator loss weighting w(t) controls the strengths put 1056 on each time level of updating assistant diffusion and the student generator. We empirically find 1057 that it is decent to set $\lambda(t)$ to be the same as the default training weighting function for the reference 1058 diffusion. And it is decent to set the w(t) = 1 for all time-levels in practice. In the following section, 1059 we give more discussions on Diff-Instruct*.

1060

1061 B.3 MORE DISCUSSIONS ON DIFF-INSTRUCT*

Flexible Choices of Divergences. Clearly, various choices of distance function $\mathbf{d}(.)$ result in different training algorithms. In this part, we discuss two instances. The first choice distance function is a simple squared distance, i.e. $\mathbf{d}(y_t) = \|y_t\|_2^2$. The corresponding derivative term writes $\mathbf{d}'(y_t) = 2y_t$. In fact, such a distance function recovers the practical diffusion distillation loss studied in Zhou et al. (2024b;a). The second distance is the pseudo-Huber distance, which shows more robust performances than the simple squared distance. The pseudo-Huber distance is defined with $\mathbf{d}(y) \coloneqq \sqrt{\|y_t\|_2^2 + c^2} - c$, where c is a pre-defined positive constant. The corresponding regularization loss (3.6) writes

1070 1071 1072

$$\mathbf{D}^{[0,T]}(p_{\theta}, p_{ref}) = -\left\{\frac{\boldsymbol{y}_t}{\sqrt{\|\boldsymbol{y}_t\|_2^2 + c^2}}\right\}^T \left\{\boldsymbol{s}_{\psi}(\boldsymbol{x}_t, t) - \nabla_{\boldsymbol{x}_t} \log q_t(\boldsymbol{x}_t | \boldsymbol{x}_0)\right\}.$$
(B.1)

1073 Here $\boldsymbol{y}_t\coloneqq \boldsymbol{s}_{p_{\mathrm{sg}[\theta],t}}(\boldsymbol{x}_t) - \boldsymbol{s}_{q_t}(\boldsymbol{x}_t).$

DI* Does Not Need Image Data When Training. One appealing advantage of DI* is the image data-free property, which means that DI* requires neither the image datasets nor synthetic images that are generated by reference diffusion models. This advantage distinguishes DI* from previous fine-tuning methods such as generative adversarial training (Goodfellow et al., 2014) which require training additional neural classifiers over image data, as well as those fine-tuning methods over large-scale synthetic or curated datasets.

Table 3: HPSv2.0 score of SDXL-based 1-step model alignment using Diff-Instruct*. We mark the
 increment over the initial model with green number to demonstrate its trend. We can see that the
 increment converges to +0.42.

1084	K IMAGES	0	102	205	307	410	512	614	717	819	922
1085	CONCEPT-ART [↑]	27.52	27.64	27.64	27.70	27.78	27.90	27.98	28.02	28.04	28.05
	Рното†	27.75	27.85	27.83	27.88	27.98	28.05	28.10	28.09	28.12	28.09
1086	ANIMATION [↑]	28.45	28.52	28.58	28.56	28.68	28.76	28.79	28.78	28.74	28.73
1007	PAINTING ↑	27.52	27.60	27.68	27.73	27.85	27.95	28.04	28.09	28.14	28.11
1087	AVG HPSv2.0↑	27.81	27.90(+0.09)	27.93(+0.11)	27.97(+0.17)	28.07(+0.26)	28.17(+0.36)	28.23(+0.42)	28.24(+0.43)	28.26(+0.05)	28.24(+0.43)

The Choice of Generator is Flexible across Broader Applications. Another interesting prop-1090 erty of DI* for alignment is its wide flexibility in the choice of generator models. We can see 1091 that, the theory of DI* only requires the generator to be able to generate output images (or data of 1092 other modalities) that are differentiable with the generator's parameters. This makes DI* a universal 1093 training method for two reasons. 1) the choice of generator architecture is flexible. The network 1094 architectures for diffusion models require the input and output to have the same dimensions. How-1095 ever, the DI* does not assign such a restriction to generator network choices. Therefore, pre-trained 1096 GAN generators, such as StyleGAN-T (Sauer et al., 2023a) and GigaGAN (Kang et al., 2023b) are also compatible with DI*. Besides, we have also shown in Section 5.2 that DI* is compatible with both UNet-based and DiT-based generator architectures. 2) student networks in broader applications 1099 may also satisfy the requirements of DI*. For instance, the neural radiance field (Mildenhall et al., 1100 2021) model used in text-to-3D generation using text-to-2D diffusion models can also be viewed as a generator. Therefore DI* can be used for such scenarios to incorporate human preference in the 1101 training process. Readers can read Poole et al. (2022), Wang et al. (2023) for more introductions. 1102

1103

1105

1083

1088 1089

1104 B.4 EXPERIMENT DETAILS FOR PRE-TRAINING AND ALIGNMENT

The human preference score (HPSv2.0) trend of DI*-SDXL-1step model. During the training 1106 of the DI*-SDXL-1step model, we monitor the change of the HPSv2.0 score as an out-of-sample 1107 validation metric. We initialize the 1-step model with DMD2-SDXL-1step model (Yin et al., 2024), 1108 which is a pretty solid SDXL-based one-step diffusion distillation model. We set the learning rate 1109 of both the one-step generator and the online diffusion model to be 1e-5 and set an exponential 1110 moving average decay rate of 0.9 for faster convergence. Following the same setting as SD1.5 and 1111 PixelArt- α experiment, we use the ImageReward as the explicit reward, while using a CFG scale 1112 of 8.0 for implicit CFG reward. Table 3 records the HPSv2.0 trend of the alignment process using 1113 Diff-Instruct*.

1114

1115 Detailed Experiment Settings for SD1.5 Experiment. For experiments of SD1.5, we use the 1116 open-sourced SD1.5 of a resolution of 512×512 as our reference diffusion in Algorithm 1. We implement our experiments based on SiD-LSG (Zhou et al., 2024a)¹, which provides a high-quality 1117 codebase for diffusion model training and distillation. We construct the one-step generator with the 1118 same architecture as the reference SD1.5 model, following the same configuration of SiD-LSG. We 1119 use the prompts of the LAION-AESTHETIC dataset with an aesthetic score larger than 6.25, which 1120 resulting a total of 3M text prompts. We only prepare the text prompts since DI* does not need 1121 image datasets. We use the off-the-shelf Image Reward² as our explicit reward model. To better 1122 explore the advantages of our score-based divergence over traditional KL divergence, we refer to a 1123 recent work (Luo, 2024) that uses KL divergence for training and conducts a detailed comparison 1124 between score-based divergences that DI* uses and KL divergences in previous works.

1125 1126

Detailed Experiment Settings for PixelArt- α **Experiment.** Different from Stabld Diffusion models, the PixArt- α model is a high-quality open-sourced text-to-image diffusion model. It uses a diffusion transformer (Peebles & Xie, 2022) to learn marginal score functions in a latent space encoded by a down-sampled variational auto-encoder (VAE) (Rombach et al., 2022). For the text conditioning mechanism, the PixelArt- α model uses a T5-XXL text encoder(Raffel et al., 2020; Tay et al., 2021), which makes the model able to understand long prompts without an obvious length

1132 1133

¹https://github.com/mingyuanzhou/SiD-LSG

²https://github.com/THUDM/ImageReward

1134 restriction. We use the DiT architecture and the 0.6B PixArt- α of a resolution of 512×512 as 1135 our reference diffusion to demonstrate the compatibility of DI* for different kinds of neural net-1136 work architectures. We use the prompts from the SAM-LLaVA-Caption-10M dataset as our prompt 1137 dataset. The SAM-LLaVA-Caption-10M dataset contains the images collected by Kirillov et al. 1138 (2023), together with text descriptions that are captioned by LLaVA model (Liu et al., 2024). The SAM-LLaVA-Caption-10M dataset is used for training the PixelArt- α model. Since the PixelArt- α 1139 diffusion model uses a T5-XXL, which is memory and computationally expensive. To speed up 1140 the alignment training, we pre-encoded the text prompts using the T5-XXL text encoder, saved the 1141 encoded embedding vectors, and built the data loaders in-house. 1142

1143 We follow the setting of Diff-Instruct (Luo et al., 2024a) to use the same neural network architecture 1144 as the reference diffusion model for the one-step generator. The PixelArt- α model is trained using so-called VP diffusion(Song et al., 2020), which first scales the data in the latent space, then adds 1145 noise to the scaled latent data. We reformulate the VP diffusion as the form of so-called data-1146 prediction proposed in EDM paper (Karras et al., 2022) by re-scaling the noisy data with the inverse 1147 of the scale that has been applied to data with VP diffusion. Under the data-prediction formulation, 1148 we select a fixed noise σ_{init} level to be $\sigma_{init} = 2.5$ following the Diff-Instruct and SiD (Zhou et al., 1149 2024b). For generation, we first generate a Gaussian vector $z \sim p_z = \mathcal{N}(\mathbf{0}, \sigma_{init}^2 \mathbf{I})$. Then we 1150 input z into the generator to generate the latent. The latent vector can then be decoded by the VAE 1151 decoder to turn into an image if needed. 1152

We put the details of how to construct the one-step generator in the following paragraphs. We 1153 also initialize the assistant diffusion model with the same weight as the reference diffusion. We 1154 use the Image Reward as the human preference reward and use the Diff-Instruct* algorithm 1 (or 1155 equivalently the algorithm 2) to train the generator. We also used the Adam optimizer with the 1156 parameter $(\beta_1, \beta_2) = (0.0, 0.999)$ for both the generator and the assistant diffusion with a batch size 1157 of 128, implemented with BF16 numerical format and the accumulate-gradient training technique. 1158 We use a fixed exponential moving average decay (EMA) rate of 0.95 for all training trials. After 1159 the training, the generator aligned with both strong CFG and reward model shows significantly 1160 improved aesthetic appearance, better generation layout, and richer image details. Figure 2 shows a demonstration of the generated images using our aligned one-step generator with a CFG scale α_{cfg} 1161 1162 of 4.5 and a reward scale α_{rew} of 10.0.

1163

Construction of the one-step generator. We follow the experiment setting of Diff-Instruct (Luo et al., 2024b), generalizing its CIFAR10 experiment to text-to-image generation. Notice that the Diff-Instruct uses the EDM model (Karras et al., 2022) to formulate the diffusion model, as well as the one-step generator. We start with a brief introduction to the EDM model.

1168 The EDM model depends on the diffusion process

$$\mathrm{d}\boldsymbol{x}_t = t\mathrm{d}\boldsymbol{w}_t, t \in [0, T]. \tag{B.2}$$

1170 Samples from the forward process (B.2) can be generated by adding random noise to the output of 1171 the generator function, i.e., $x_t = x_0 + t\epsilon$ where $\epsilon \sim \mathcal{N}(0, I)$ is a Gaussian vector. The EDM 1172 model also reformulates the diffusion model's score matching objective as a denoising regression 1173 objective, which writes,

1169

$$\mathcal{L}(\psi) = \int_{t=0}^{T} \lambda(t) \mathbb{E}_{\boldsymbol{x}_{0} \sim p_{0}, \boldsymbol{x}_{t} \mid \boldsymbol{x}_{0} \sim p_{t}(\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0})} \|\boldsymbol{d}_{\psi}(\boldsymbol{x}_{t}, t) - \boldsymbol{x}_{0}\|_{2}^{2} \mathrm{d}t.$$
(B.3)

Where $d_{\psi}(\cdot)$ is a denoiser network that tries to predict the clean sample by taking noisy samples as inputs. Minimizing the loss (B.3) leads to a trained denoiser, which has a simple relation to the marginal score functions as:

1179 1180

1187

$$s_{\psi}(\boldsymbol{x}_t, t) = \frac{\boldsymbol{d}_{\psi}(\boldsymbol{x}_t, t) - \boldsymbol{x}_t}{t^2} \tag{B}$$

.4)

¹¹⁸¹ Under such a formulation, we actually have pre-trained denoiser models for experiments. Therefore, we use the EDM notations in later parts.

Let $d_{\theta}(\cdot)$ be pretrained EDM denoiser models. Owing to the denoiser formulation of the EDM model, we construct the generator to have the same architecture as the pre-trained EDM denoiser with a pre-selected index t^* , which writes

$$\boldsymbol{x}_0 = g_{\theta}(\boldsymbol{z}) \coloneqq \boldsymbol{d}(\boldsymbol{z}, t^*), \ \boldsymbol{z} \sim \mathcal{N}(\boldsymbol{0}, (t^*)^2 \mathbf{I}).$$
 (B.5)

We initialize the generator with the same parameter as the teacher EDM denoiser model.

x

Time index distribution. When training both the EDM diffusion model and the generator, we need to randomly select a time t in order to approximate the integral of the loss function (B.3). The EDM model has a default choice of t distribution as log-normal when training the diffusion (denoiser) model, i.e.

$$t \sim p_{EDM}(t): \quad t = \exp(s) \tag{B.6}$$

$$s \sim \mathcal{N}(P_{mean}, P_{std}^2), \ P_{mean} = -2.0, P_{std} = 2.0.$$
 (B.7)

And a weighting function

1192 1193

1194

1197

1198

1205 1206 1207

1208

1214

$$\lambda_{EDM}(t) = \frac{(t^2 + \sigma_{data}^2)}{(t \times \sigma_{data})^2}.$$
(B.8)

In our algorithm, we follow the same setting as the EDM model when updating the online diffusion (denoiser) model.

1202 Weighting function. For the assistant diffusion updates in both pre-training and alignment, we 1203 use the same $\lambda_{EDM}(t)$ (B.8) weighting function as EDM when updating the denoiser model. When 1204 updating the generator, we use a specially designed weighting function, which writes:

$$w_{Gen}(t) = \frac{1}{\|\boldsymbol{d}_{\psi}(\mathrm{sg}[\boldsymbol{x}_t], t) - \boldsymbol{d}_{q_t}(\mathrm{sg}[\boldsymbol{x}_t], t)\|_2}$$
(B.9)

$$\boldsymbol{x}_t = \boldsymbol{x}_0 + t\boldsymbol{\epsilon}, \ \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$
 (B.10)

The notation $sg[\cdot]$ means stop-gradient of the parameter. Such a weighting function helps to stabilize the training.

¹²¹¹ In the both Stable Diffusion 1.5 and the PixArt- α experiments, we rewrite the PixelArt- α model in EDM formulation:

$$D_{\theta}(\mathbf{x};\sigma) = \mathbf{x} - \sigma F_{\theta} \tag{B.11}$$

1215 Here, following the iDDPM+DDIM preconditioning in EDM, PixelArt- α is denoted by F_{θ} , x is the 1216 image data plus noise with a standard deviation of σ , for the remaining parameters such as C_1 and C_2 , we kept them unchanged to match those defined in EDM. Unlike the original model, we only 1217 retained the image channels for the output of this model. Since we employed the preconditioning 1218 of iDDPM+DDIM in the EDM, each σ value is rounded to the nearest 1000 bins after being passed 1219 into the model. For the actual values used in PixelArt- α , beta_start is set to 0.0001, and beta_end 1220 is set to 0.02. Therefore, according to the formulation of EDM, the range of our noise distribution 1221 is [0.01, 156.6155], which will be used to truncate our sampled σ . For our one-step generator, it is 1222 formulated as: 1223

$$g_{\theta}(\mathbf{x};\sigma_{\text{init}}) = \mathbf{x} - \sigma_{\text{init}}F_{\theta} \tag{B.12}$$

Here following Diff-Instruct to use $\sigma_{init} = 2.5$ and $\mathbf{x} \sim \mathcal{N}(0, \sigma_{init}\mathbf{I})$, we observed in practice that larger values of σ_{init} lead to faster convergence of the model, but the difference in convergence speed is negligible for the complete model training process and has minimal impact on the final results.

1228 **Detailed Quantitative Evaluations Metrics** To quantitatively evaluate the performances of the 1229 generators trained with different alignment settings, we compare our generators elaborately with 1230 other open-sourced models that are either based on SD1.5 diffusion models or larger models 1231 such as SDXL. For all models, we compute four standard scores: the Image Reward (Xu et al., 1232 2023a), the Aesthetic Score (Schuhmann, 2022), the PickScore(Kirstain et al., 2023), and the CLIP 1233 score(Radford et al., 2021). Since most existing literature tests the human preference scores with different prompts which are possibly not available, to make a fair comparison, in our experiment, we 1234 fix 1k prompts from the COCO-2017 (Lin et al., 2014) validation dataset and intensively evaluate a 1235 wide range of existing open-sourced models as SiD-LSG (Zhou et al., 2024a), SDXL-Turbo (Sauer 1236 et al., 2023b), Latent Consistency Model (LCM) (Luo et al., 2023a), Hyper-SD (Ren et al., 2024), 1237 SDXL-Lightning (Lin et al., 2024), Trajectory Consistency Model (TCD) (Zheng et al., 2024), PeReflow (Yan et al., 2024), InstaFLow(Liu et al., 2023), Diffusion DPO(Wallace et al., 2024), 1239 etc. All models are tested with the same prompts and the same computing devices. 1240

1241 Besides the COCO prompts, we also evaluate generators with open-sourced models with Human Preference Score v2.0 (HPSv2.0) (Wu et al., 2023) over their benchmark prompts. The HPS is a

1277 1278 1279

1281

1282

1283

1284

1285

1286

widely used standard benchmark that evaluates models' capability of generating images of 4 styles: Animation, concept art, Painting, and Photo. The score reflects the prompt following and the human preference strength of text-to-image models. We use the HPSv2's ³ default protocols for evaluations. Since our generators based on PixArt- α are trained with SAM-recaptioned datasets, we also evaluate the performances on 30K prompts from the SAM-recap dataset.

Sample Prompts from COCO validation dataests. In this paragraph, we propose some sample 1248 prompts from the COCO validation dataset that we use for calculating scores like Image Reward 1249 1250 and the Aesthetic Scores. 1251 1 This wire metal rack holds several pairs of shoes and sandals 1252 2 A motorcycle parked in a parking space next to another motorcycle. 1253 3 A picture of a dog laying on the ground. 4 A loft bed with a dresser underneath it. 1254 5 Two giraffes in a room with people looking at them. 1255 6 A woman stands in the dining area at the table. 1256 7 Birds perch on a bunch of twigs in the winter. 1257 8 A small kitchen with low a ceiling 1258 9 A group of baseball players is crowded at the mound. 1259 10 This table is filled with a variety of different dishes. 11 A toy dinosaur standing on a sink next to a running faucet. 1260 12 a man standing holding a game controller and two people sitting 1261 13 There is a small bus with several people standing next to it. 1262 14 A bottle on wine next to a glass of wine.

1263 15 A big burly grizzly bear is show with grass in the background.

1264 16 A man standing in front of a microwave next to pots and pans.

17 Three men in military suits are sitting on a bench, 18 Two people standing in a kitchen looking around.

1266 19 A group of men playing a game of baseball on top of a baseball field.

1267 20 A traffic light over a street surrounded by tall buildings.

1268 21 The snowboarder has jumped high into the air from a snow ramp.

1269 22 A smart phone with an image of a person on it's screen. 23 A man talking on his phone in the public.

1270 24 A cheesy pizza sitting on top of a table.

1271 25 A dog sitting on the inside of a white boat.

1272 26 A guy jumping with a tennis racket in his hand.

1273 27 a close up of a child next to a cake with balloons

1274 ²⁸ A man holding a camera up over his left shoulder.

29 A plane flies over water with two islands nearby.

12/5 30 A young boy getting ready to catch a baseball in a grass field.

Listing 1: Example Prompts from COCO validation dataset.

1279 B.5 MORE DISCUSSIONS ON FINDINGS OF QUALITATIVE EVALUATIONS

There are some other interesting findings when qualitatively evaluate different models.

• First, we find that the images generated by the aligned model show a better composition when organizing the contents presented in the image. For instance, the main objects of the generated image are smaller and show a more natural layout than other models, with the objects and the background iterating aesthetically. This in turn reveals human preference: human beings would prefer that the object of an image does not take up all spaces of an image;

Second, we find that the aligned model has richer details than the unaligned model. The stronger we align the model, the richer details the model will generate. Sometimes these rich details come as a hint to the readers about the input prompts. Sometimes they just come to improve the aesthetic performance. We think this phenomenon may be caused by the fact that human prefers images with rich details. Another finding is that as the reward scale for alignment becomes stronger, the generated image from the alignment model becomes more colorful and more similar to paintings. Sometimes this leads to a loss of reality to some

³https://github.com/tgxs002/HPSv2



degree. Therefore, we think that users should choose different aligned one-step models with a trade-off between aesthetic performance and image reality according to the use case.

Figure 5: Bad generation cases by aligned DiT-DI* one-step generator model (4.5 CFG + 10.0 reward).

Table 4: Hyperparameters used for Diff-Instruct* on SD1.5 and PixArt- α experiments.

Hyperparameter	SD1.5 Ex	periment	PixArt- α Experiment			
	$DM \boldsymbol{s}_\psi$	Generator g_{θ}	$DM~\boldsymbol{s}_\psi$	Generator g_{θ}		
Learning rate	1e-5	1e-5	2e-6	2e-6		
Batch size	512	512	256	256		
$\sigma(t^*)$	2.5	2.5	2.5	2.5		
Adam β_0	0.0	0.0	0.0	0.0		
Adam β_1	0.999	0.999	0.999	0.999		
EMA decay rate	0.9	0.9	0.95	0.95		
Time Distribution	$p_{EDM}(t)$ (B.6)	$p_{EDM}(t)$ (B.6)	$p_{EDM}(t)$ (B.6)	$p_{EDM}(t)$ (B.6)		
Weighting	$\lambda_{EDM}(t)$ (B.8)	1	$\lambda_{EDM}(t)$ (B.8)	1		
Number of GPUs	8×H800-80G	8×H800-80G	8×H100-80G	8×H800-80G		

With the optimal setting and EDM formulation, we can rewrite our algorithm in an EDM style in Algorithm 3.

		Algorithm 3: Diff-Instruct* Pseudo Code under EDM formulation.
		Input: prompt dataset C, generator $q_{\theta}(x_0 z, c)$, prior distribution p_z , reward model $r(x, c)$,
		reward scale α_{rew} , CFG scale α_{cfg} , reference EDM denoiser model $d_{ref}(x_t c,c)$,
		assistant EDM denoiser $d_{\psi}(x_t t, c)$, forward diffusion $p(x_t x_0)$ (2.1), assistant EDM
		denoiser updates rounds K_{TA} , time distribution $\pi(t)$, diffusion model weighting $\lambda(t)$,
		generator IKL loss weighting $w(t)$.
		while not converge do
		fix θ , update ψ for K_{TA} rounds by
		1. sample prompt $c \sim C$; sample time $t \sim \pi(t)$; sample $z \sim p_z(z)$;
		2. generate fake data: $\boldsymbol{x}_0 = \mathrm{sg}[g_\theta(\boldsymbol{z}, \boldsymbol{c})]$; sample noisy data: $\boldsymbol{x}_t \sim p_t(\boldsymbol{x}_t \boldsymbol{x}_0)$;
		3. update ψ by minimizing loss: $\mathcal{L}(\psi) = \lambda(t) \ \boldsymbol{d}_{\psi}(\boldsymbol{x}_t t, \boldsymbol{c}) - \boldsymbol{x}_0 \ _2^2$;
		fix ψ , update θ using StaD:
		1. sample prompt $\boldsymbol{c} \sim \mathcal{C}$; sample time $t \sim \pi(t)$; sample $\boldsymbol{z} \sim p_z(\boldsymbol{z})$;
		2. generate fake data: $\mathbf{x}_0 = a_0(\mathbf{z}, \mathbf{c})$: sample noisy data: $\mathbf{x}_1 \sim a_0(\mathbf{x}_1 \mathbf{x}_0)$:
		2. generate rate data: $w_0 = g_0(x, c)$, sample hold data: $w_t = p_t(w_t w_0)$,
		5. explicit reward: $\mathcal{L}_{rew}(\theta) = -\alpha_{rew} r(\boldsymbol{x}_0, \boldsymbol{c});$
		4. CFG reward: $\mathcal{L}_{cfg}(\theta) = \alpha_{cfg} \cdot w(t) \{ \boldsymbol{d}_{ref}(\mathrm{sg}[\boldsymbol{x}_t] t, \boldsymbol{c}) - \boldsymbol{d}_{ref}(\mathrm{sg}[\boldsymbol{x}_t] t, \boldsymbol{\varnothing}) \}^T \boldsymbol{x}_t;$
		5. score-regularization:
		$\int (\theta) = u(t) \left[\frac{d'(d_{1}(m \mid t, a))}{d_{1}(m \mid t, a)} - \frac{d_{1}(m \mid t, a)}{d_{1}(m \mid t, a)} \right]^{T} \left[\frac{d_{1}(m \mid t, a)}{d_{1}(m \mid t, a)} - \frac{d_{1}(m \mid t, a)}{d_{1}(m \mid t, a)} \right]^{T}$
		$\mathcal{L}_{reg}(v) = -w(v) \{ \mathbf{u} \left(\mathbf{u}_{\psi}(\mathbf{x}_t v, \mathbf{c}) - \mathbf{u}_{ref}(\mathbf{x}_t v, \mathbf{c}) \right) \} \{ \mathbf{u}_{\psi}(\mathbf{x}_t v, \mathbf{c}) - \mathbf{x}_0 \};$
		6. update θ by minimizing DI* loss: $\mathcal{L}_{DI*}(\theta) = \mathcal{L}_{rew}(\theta) + \mathcal{L}_{cfg}(\theta) + \mathcal{L}_{reg}(\theta)$;
		end
		return $ heta, \psi$.
		B.6 PYTORCH STYLE PSEUDO-CODE OF SCORE IMPLICIT MATCHING
		In this section, we give a PyTorch style pseudo-code for algorithm 3.
1	1	In this section, we give a PyTorch style pseudo-code for algorithm 3.
1	1	In this section, we give a PyTorch style pseudo-code for algorithm 3. import torch import torch.nn as nn
1	1 2 3	In this section, we give a PyTorch style pseudo-code for algorithm 3. import torch import torch.nn as nn import torch.optim as optim
1	1 2 3 4	In this section, we give a PyTorch style pseudo-code for algorithm 3. import torch import torch.nn as nn import torch.optim as optim import copy
	1 2 3 4 5	In this section, we give a PyTorch style pseudo-code for algorithm 3. import torch import torch.nn as nn import torch.optim as optim import copy use cfg = True
	1 2 3 4 5 6 7	<pre>In this section, we give a PyTorch style pseudo-code for algorithm 3. import torch import torch.nn as nn import torch.optim as optim import copy use_cfg = True use_reward = True</pre>
	1 2 3 4 5 6 7 8	<pre>In this section, we give a PyTorch style pseudo-code for algorithm 3. import torch import torch.nn as nn import torch.optim as optim import copy use_cfg = True use_reward = True</pre>
	1 2 3 4 5 6 7 8 9	<pre>In this section, we give a PyTorch style pseudo-code for algorithm 3. import torch import torch.nn as nn import torch.optim as optim import copy use_cfg = True use_reward = True # Initialize generator G</pre>
	1 2 3 4 5 6 7 8 9 0	<pre>In this section, we give a PyTorch style pseudo-code for algorithm 3. import torch import torch.nn as nn import torch.optim as optim import copy use_cfg = True use_reward = True # Initialize generator G G = Generator()</pre>
	1 2 3 4 5 6 7 8 9 0 1 2	<pre>In this section, we give a PyTorch style pseudo-code for algorithm 3. import torch import torch.nn as nn import torch.optim as optim import copy use_cfg = True use_reward = True # Initialize generator G G = Generator() ## load teacher DM</pre>
	1 2 3 4 5 6 7 8 9 0 1 2 3	<pre>In this section, we give a PyTorch style pseudo-code for algorithm 3. import torch import torch.nn as nn import torch.optim as optim import copy use_cfg = True use_reward = True # Initialize generator G G = Generator() ## load teacher DM Drf = DiffusionModel().load('/path_to_ckpt').eval().requires grad (False)</pre>
1 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2	1 2 3 4 5 6 7 8 9 0 1 2 3 4	<pre>In this section, we give a PyTorch style pseudo-code for algorithm 3. import torch import torch.nn as nn import torch.optim as optim import copy use_cfg = True use_reward = True # Initialize generator G G = Generator() ## load teacher DM Drf = DiffusionModel().load('/path_to_ckpt').eval().requires_grad_(False) Dta = copy.deepcopy(Drf) ## initialize online DM with teacher DM</pre>
1 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2	1 2 3 4 5 6 7 8 9 0 1 2 3 4 5	<pre>In this section, we give a PyTorch style pseudo-code for algorithm 3. import torch import torch.nn as nn import torch.optim as optim import copy use_cfg = True use_reward = True # Initialize generator G G = Generator() ## load teacher DM Drf = DiffusionModel().load('/path_to_ckpt').eval().requires_grad_(False) Dta = copy.deepcopy(Drf) ## initialize online DM with teacher DM r = RewardModel() if use_reward else None</pre>
	1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7	<pre>In this section, we give a PyTorch style pseudo-code for algorithm 3. import torch import torch.nn as nn import torch.optim as optim import copy use_cfg = True use_reward = True # Initialize generator G G = Generator() ## load teacher DM Drf = DiffusionModel().load('/path_to_ckpt').eval().requires_grad_(False) Dta = copy.deepcopy(Drf) ## initialize online DM with teacher DM r = RewardModel() if use_reward else None # Dofine optimizers</pre>
	1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8	<pre>In this section, we give a PyTorch style pseudo-code for algorithm 3. import torch import torch.nn as nn import torch.optim as optim import copy use_cfg = True use_reward = True # Initialize generator G G = Generator() ## load teacher DM Drf = DiffusionModel().load('/path_to_ckpt').eval().requires_grad_(False) Dta = copy.deepcopy(Drf) ## initialize online DM with teacher DM r = RewardModel() if use_reward else None # Define optimizers opt G = optim.Adam(G.parameters(), lr=0.001, betas=(0.0, 0.999))</pre>
	1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9	<pre>In this section, we give a PyTorch style pseudo-code for algorithm 3. import torch import torch.nn as nn import torch.optim as optim import copy use_cfg = True use_reward = True # Initialize generator G G = Generator() ## load teacher DM Drf = DiffusionModel().load('/path_to_ckpt').eval().requires_grad_(False) Dta = copy.deepcopy(Drf) ## initialize online DM with teacher DM r = RewardModel() if use_reward else None # Define optimizers opt_G = optim.Adam(G.parameters(), lr=0.001, betas=(0.0, 0.999)) opt_Sta = optim.Adam(Dta.parameters(), lr=0.001, betas=(0.0, 0.999))</pre>
1 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2	1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0	<pre>In this section, we give a PyTorch style pseudo-code for algorithm 3. import torch import torch.nn as nn import torch.optim as optim import copy use_cfg = True use_reward = True # Initialize generator G G = Generator() ## load teacher DM Drf = DiffusionModel().load('/path_to_ckpt').eval().requires_grad_(False) Dta = copy.deepcopy(Drf) ## initialize online DM with teacher DM r = RewardModel() if use_reward else None # Define optimizers opt_G = optim.Adam(G.parameters(), lr=0.001, betas=(0.0, 0.999)) opt_Sta = optim.Adam(Dta.parameters(), lr=0.001, betas=(0.0, 0.999))</pre>
1 2 2 2 2 2 2 2 2 2 2 2 2 2 2 10 11 12 13 14 15 10 17 18 19 20 21	1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1	<pre>In this section, we give a PyTorch style pseudo-code for algorithm 3. import torch import torch.nn as nn import torch.optim as optim import copy use_cfg = True use_reward = True # Initialize generator G G = Generator() ## load teacher DM Drf = DiffusionModel().load('/path_to_ckpt').eval().requires_grad_(False) Dta = copy.deepcopy(Drf) ## initialize online DM with teacher DM r = RewardModel() if use_reward else None # Define optimizers opt_G = optim.Adam(G.parameters(), lr=0.001, betas=(0.0, 0.999)) opt_Sta = optim.Adam(Dta.parameters(), lr=0.001, betas=(0.0, 0.999))</pre>
1	1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3	<pre>In this section, we give a PyTorch style pseudo-code for algorithm 3. import torch import torch.nn as nn import torch.optim as optim import copy use_cfg = True use_reward = True # Initialize generator G G = Generator() ## load teacher DM Drf = DiffusionModel().load('/path_to_ckpt').eval().requires_grad_(False) Dta = copy.deepcopy(Drf) ## initialize online DM with teacher DM r = RewardModel() if use_reward else None # Define optimizers opt_G = optim.Adam(G.parameters(), lr=0.001, betas=(0.0, 0.999)) opt_Sta = optim.Adam(Dta.parameters(), lr=0.001, betas=(0.0, 0.999)) # Training loop while True: # "" color by the pro- # ""</pre>
$\begin{array}{c} 1 \\ 2 \\ 2 \\ 4 \\ 6 \\ 6 \\ 7 \\ 7 \\ 8 \\ 8 \\ 9 \\ 9 \\ 9 \\ 10 \\ 11 \\ 12 \\ 13 \\ 14 \\ 15 \\ 16 \\ 10 \\ 12 \\ 12 \\ 20 \\ 21 \\ 22 \\ 22 \\ 22$	1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4	<pre>In this section, we give a PyTorch style pseudo-code for algorithm 3. import torch import torch.nn as nn import torch.optim as optim import copy use_cfg = True use_reward = True # Initialize generator G G = Generator() ## load teacher DM Drf = DiffusionModel().load('/path_to_ckpt').eval().requires_grad_(False) Dta = copy.deepcopy(Drf) ## initialize online DM with teacher DM r = RewardModel() if use_reward else None # Define optimizers opt_G = optim.Adam(G.parameters(), lr=0.001, betas=(0.0, 0.999)) opt_Sta = optim.Adam(Dta.parameters(), lr=0.001, betas=(0.0, 0.999)) # Training loop while True: ## update Dta Dta train() requires grad (True)</pre>
$\begin{array}{c} 1 \\ 2 \\ 3 \\ 4 \\ 6 \\ 0 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1$	1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 8 9 0 1 2 3 4 5 8 9 0 1 2 3 4 5 8 9 0 1 2 3 4 5 8 9 0 1 2 3 8 9 0 1 2 3 4 5 8 9 0 1 2 3 4 5 8 9 0 1 2 3 3 4 5 7 8 9 0 1 2 3 3 4 5 7 8 9 0 1 2 3 4 5 7 8 9 0 1 2 3 4 5 5 7 8 9 0 1 2 3 3 5 7 8 9 0 1 2 3 3 3 5 8 9 0 1 2 3 7 8 9 0 1 2 3 4 5 7 8 9 0 1 2 3 3 4 5 7 8 9 0 1 2 3 8 9 0 1 2 3 1 2 3 4 5 7 8 9 0 1 2 3 3 7 8 9 0 1 2 3 1 2 1 2	<pre>In this section, we give a PyTorch style pseudo-code for algorithm 3. import torch import torch.nn as nn import torch.optim as optim import copy use_cfg = True use_reward = True # Initialize generator G G = Generator() ## load teacher DM Drf = DiffusionModel().load('/path_to_ckpt').eval().requires_grad_(False) Dta = copy.deepcopy(Drf) ## initialize online DM with teacher DM r = RewardModel() if use_reward else None # Define optimizers opt_G = optim.Adam(G.parameters(), lr=0.001, betas=(0.0, 0.999)) opt_Sta = optim.Adam(Dta.parameters(), lr=0.001, betas=(0.0, 0.999)) # Training loop while True: ## update Dta Dta.train().requires_grad_(True) G.eval().requires grad_(False)</pre>
1	1 2 3 4 4 5 6 7 8 9 0 0 1 2 3 3 4 5 6 7 8 9 0 0 1 2 3 3 4 5 6 7 8 9 0 0 1 2 3 4 4 5 6 7 7 8 9 0 0 1 1 2 3 4 4 5 5 6 7 7 8 9 0 0 1 1 2 3 3 4 5 5 6 7 7 8 9 9 0 0 1 1 2 3 3 4 5 5 6 7 7 8 9 9 0 0 1 1 2 3 3 4 5 5 6 7 7 8 9 9 0 0 1 1 2 3 3 4 5 5 6 7 7 8 9 9 0 0 1 1 2 3 3 4 4 5 5 6 7 7 8 9 9 0 0 1 1 2 3 3 4 4 5 5 6 7 7 8 9 9 0 0 1 1 2 3 3 4 4 5 5 7 8 9 9 0 0 1 1 2 3 3 4 4 5 5 7 7 8 9 9 0 0 1 1 2 3 3 4 4 5 5 7 7 8 9 9 0 0 1 1 2 3 3 4 4 5 5 7 7 7 8 9 9 0 0 1 1 2 3 3 4 4 5 5 9 9 0 0 1 1 2 3 3 4 4 5 5 7 7 7 8 9 9 0 0 1 1 2 3 3 4 4 5 5 5 6 7 7 8 9 9 0 0 1 1 2 3 3 4 4 5 5 7 7 8 9 9 0 0 1 1 2 3 3 4 4 5 5 7 7 7 7 8 9 9 0 0 1 1 2 3 3 4 4 5 5 7 7 7 7 8 9 9 0 0 1 1 2 3 3 7 7 8 9 9 0 0 1 1 2 3 3 3 7 8 9 9 9 0 1 1 2 3 3 1 1 1 2 3 3 1 1 1 2 3 3 1 1 1 2 3 3 1 1 1 2 3 3 1 1 1 2 3 1 1 1 2 1 2	<pre>In this section, we give a PyTorch style pseudo-code for algorithm 3. import torch import torch.nn as nn import torch.optim as optim import copy use_cfg = True use_reward = True # Initialize generator G G = Generator() ## load teacher DM Drf = DiffusionModel().load('/path_to_ckpt').eval().requires_grad_(False) Dta = copy.deepcopy(Drf) ## initialize online DM with teacher DM r = RewardModel() if use_reward else None # Define optimizers opt_G = optim.Adam(G.parameters(), lr=0.001, betas=(0.0, 0.999)) opt_Sta = optim.Adam(Dta.parameters(), lr=0.001, betas=(0.0, 0.999)) # Training loop while True: ## update Dta Dta.train().requires_grad_(True) G.eval().requires_grad_(False)</pre>
1	1 2 3 4 4 5 6 7 8 9 0 0 1 2 3 3 4 5 6 6 7 8 9 0 0 1 2 3 4 5 5 6 7 7 8 9 0 0 1 2 3 4 5 5 6 7 7 8 9 0 0 1 2 3 4 5 5 6 7 7 8 9 0 0 1 1 2 3 4 5 5 6 6 7 7 8 9 0 0 1 1 2 3 4 5 5 6 6 7 7 8 9 9 0 0 1 1 2 3 4 5 5 6 6 7 7 8 9 9 0 0 1 1 2 3 4 5 5 6 6 7 7 8 9 9 0 0 1 1 2 3 4 5 5 6 6 7 7 8 9 9 0 0 1 1 2 3 4 4 5 5 6 6 7 7 8 8 9 9 0 0 1 1 2 3 4 4 5 5 6 7 7 8 8 9 0 0 1 1 2 3 4 4 5 5 8 9 0 0 1 1 2 3 4 4 5 5 7 8 9 0 0 1 1 2 3 4 4 5 5 8 9 0 0 1 1 2 3 5 6 7 7 8 9 9 0 0 1 1 2 3 7 7 7 8 8 9 0 0 1 1 2 3 7 8 7 7 8 9 9 0 0 1 1 2 3 7 8 9 9 0 0 1 1 2 3 7 7 7 8 8 9 9 0 0 1 1 2 3 7 7 8 8 9 9 0 0 1 1 2 3 3 4 5 7 7 8 8 9 9 0 0 1 1 2 3 7 8 8 8 9 9 0 0 1 1 2 3 7 8 8 9 9 9 0 0 1 1 2 3 3 8 8 9 9 9 0 1 1 1 2 3 3 8 8 9 9 9 1 1 1 2 3 3 8 9 9 1 1 1 2 3 1 1 1 2 3 3 1 1 1 2 3 1 1 1 1	<pre>In this section, we give a PyTorch style pseudo-code for algorithm 3. import torch import torch.nn as nn import torch.optim as optim import copy use_cfg = True use_reward = True # Initialize generator G G = Generator() ## load teacher DM Drf = DiffusionModel().load('/path_to_ckpt').eval().requires_grad_(False) Dta = copy.deepcopy(Drf) ## initialize online DM with teacher DM r = RewardModel() if use_reward else None # Define optimizers opt_G = optim.Adam(Dta.parameters(), lr=0.001, betas=(0.0, 0.999)) opt_Sta = optim.Adam(Dta.parameters(), lr=0.001, betas=(0.0, 0.999)) # Training loop while True: ## update Dta Dta.train().requires_grad_(True) G.eval().requires_grad_(False) ## update assistant diffusion</pre>
1	1 2 3 4 5 6 7 8 9 0 1 1 2 3 4 5 6 7 8 9 0 1 1 2 3 4 5 6 7 8 9 0 1 1 2 3 4 5 6 7 8 9 0 1 1 2 3 4 5 6 7 7 8 9 0 0 1 1 2 3 4 5 6 7 7 8 9 0 0 1 1 2 3 4 5 6 7 7 8 9 0 0 1 1 2 3 4 5 6 7 7 8 9 0 0 1 1 2 3 4 5 6 7 7 8 9 0 0 1 1 2 3 4 5 6 7 7 8 9 0 0 1 1 2 3 4 5 6 7 7 8 9 0 0 1 1 2 3 4 5 6 7 7 8 9 0 0 1 1 2 3 4 5 6 7 7 8 9 0 0 1 1 2 3 4 5 7 7 8 9 0 0 1 1 2 3 4 4 5 7 8 9 0 0 1 1 2 3 3 4 5 5 6 7 7 8 9 0 0 1 1 2 3 3 4 5 7 7 8 9 0 0 1 1 2 3 3 4 5 7 7 8 9 0 0 1 1 2 3 7 7 8 9 0 0 1 1 2 3 3 4 5 7 8 9 0 0 1 1 2 3 3 4 5 5 6 7 7 8 9 0 0 1 1 2 3 7 8 9 0 0 1 1 2 3 3 4 5 5 6 7 7 8 9 0 0 1 1 2 3 4 5 5 7 8 9 0 0 1 1 2 3 4 5 7 8 9 0 0 1 1 2 3 4 5 7 8 9 0 0 1 1 1 2 3 3 4 5 5 7 8 9 0 0 1 1 2 3 4 5 7 8 9 0 1 1 2 3 4 5 7 8 9 9 0 1 1 1 2 3 1 1 1 2 3 1 1 1 2 3 1 1 1 1 2 3 1 1 1 1	<pre>In this section, we give a PyTorch style pseudo-code for algorithm 3. import torch import torch.nn as nn import torch.optim as optim import copy use_cfg = True use_reward = True # Initialize generator G G = Generator() ## load teacher DM Drf = DiffusionModel().load('/path_to_ckpt').eval().requires_grad_(False) Dta = copy.deepcopy(Drf) ## initialize online DM with teacher DM r = RewardModel() if use_reward else None # Define optimizers opt_G = optim.Adam(G.parameters(), lr=0.001, betas=(0.0, 0.999)) opt_Sta = optim.Adam(Dta.parameters(), lr=0.001, betas=(0.0, 0.999)) # Training loop while True: ## update Dta Dta.train().requires_grad_(True) G.eval().requires_grad_(False) ## update assistant diffusion prompt = batch['prompt']</pre>
$\begin{array}{c} 1\\ 2\\ 3\\ 4\\ 6\\ 6\\ 7\\ 8\\ 8\\ 9\\ 9\\ 10\\ 11\\ 11\\ 13\\ 14\\ 15\\ 10\\ 11\\ 13\\ 14\\ 15\\ 10\\ 21\\ 22\\ 22\\ 22\\ 22\\ 22\\ 22\\ 22\\ 22\\ 22$	1 2 3 4 4 5 6 6 7 8 9 0 0 1 2 3 4 4 5 6 6 7 8 9 0 0 1 2 3 4 4 5 6 7 7 8 9 0 0 1 2 3 4 4 5 6 7 7 8 9 9 0 1 1 2 3 4 4 5 7 8 9 9 0 0 1 1 2 3 4 4 5 7 8 9 9 0 0 1 1 2 3 4 4 5 8 9 9 0 0 1 1 2 3 4 4 5 8 9 9 0 0 1 1 2 3 4 4 5 8 9 9 0 0 1 1 2 3 4 4 5 8 9 0 0 1 1 2 3 4 4 5 8 9 9 0 0 1 1 2 3 4 4 5 5 6 6 7 7 8 9 9 0 0 1 1 2 3 4 4 5 5 6 6 7 7 8 8 9 9 0 0 1 1 2 3 4 4 5 5 6 6 7 7 8 8 9 0 0 1 1 2 3 4 4 5 5 7 8 8 9 0 0 1 1 2 3 4 4 5 5 8 9 0 0 1 1 2 3 3 4 4 5 5 8 9 0 0 1 1 2 3 3 4 5 7 8 9 0 0 1 1 2 3 3 4 5 7 8 9 9 0 0 1 1 2 3 3 4 5 7 8 9 9 0 0 1 1 2 3 3 4 5 7 8 9 0 0 1 1 2 3 3 4 5 7 8 9 0 0 1 1 2 3 3 4 4 5 5 8 9 0 0 1 1 2 3 3 4 5 7 8 9 0 0 1 1 2 3 3 4 4 5 5 8 9 0 0 1 1 2 3 3 4 4 5 5 8 9 0 1 1 2 3 3 1 2 3 3 1 1 2 3 3 3 4 5 5 5 7 8 8 9 0 0 1 1 2 3 5 8 1 2 8 9 9 0 1 1 2 3 1 2 8 1 2 8 1 2 1 2 3 8 1 2 8 1 8 1 8 9 1 1 1 2 8 1 8 1 8 1 8 1 1 1 1 2 8 1 8 1	<pre>In this section, we give a PyTorch style pseudo-code for algorithm 3. import torch import torch.optim as optim import torch.optim as optim import copy use_cfg = True use_reward = True # Initialize generator G G = Generator() ## load teacher DM Drf = DiffusionModel().load('/path_to_ckpt').eval().requires_grad_(False) Dta = copy.deepcopy(Drf) ## initialize online DM with teacher DM r = RewardModel() if use_reward else None # Define optimizers opt_G = optim.Adam(G.parameters(), lr=0.001, betas=(0.0, 0.999)) opt_Sta = optim.Adam(Dta.parameters(), lr=0.001, betas=(0.0, 0.999)) # Training loop while True: ## update Dta Dta.train().requires_grad_(True) G.eval().requires_grad_(False) ## update assistant diffusion prompt = batch['prompt'] z = torch.randn((l024, 4, 64, 64), device=G.device) with torceh no gma().</pre>

```
1458
    31
                fake_x0 = G(z, prompt)
1459 32
1460 33
           sigma = torch.exp(2.0*torch.randn([1,1,1,1], device=fake_x0.device) -
            2.0)
1461
1462 <sup>34</sup>
           noise = torch.randn_like(fake_x0)
           fake_xt = fake_x0 + sigma*noise
     35
1463
           pred_x0 = Dta(fake_xt, sigma, prompt)
     36
1464 37
1465 38
           weight = compute_diffusion_weight(sigma)
1466 39
1467 <sup>40</sup>
           batch_loss = weight * (pred_x0 - fake_x0) **2
1468 <sup>41</sup>
           batch_loss = batch_loss.sum([1,2,3]).mean()
     42
1469 43
           optimizer_Dta.zero_grad()
1470 44
           batch_loss.backward()
           optimizer_Dta.step()
1471 45
1472 <sup>46</sup>
1473 47
    48
           ## update G
1474 49
           Dta.eval().requires_grad_(False)
1475 50
           G.train().requires_grad_(True)
1476 51
1477 <sup>52</sup>
          prompt = batch['prompt']
           z = torch.randn((1024, 4, 64, 64), device=G.device)
    53
1478
    54
           fake_x0 = G(z, prompt)
1479 55
1480 56
           sigma = torch.exp(2.0*torch.randn([1,1,1,1], device=fake_x0.device) -
1481
            2.0)
           noise = torch.randn_like(fake_x0)
1482 57
           fake_xt = fake_x0 + sigma*noise
    58
1483
     59
1484 60
           with torch.no_grad():
1485 61
               if use_cfg:
                    cfg_vector = (Drf(fake_xt, sigma, prompt) - Drf(fake_xt,
1486 62
           sigma, None)
1487
1488<sup>63</sup>
                else:
    64
                    cfg_vector = None
1489 65
1490 66
                pred_x0_rf = Drf(fake_xt, sigma, prompt)
                pred_x0_ta = Dta(fake_xt, sigma, prompt)
1491 67
1492 <sup>68</sup>
           denoise_diff = pred_x0_ta - pred_x0_rf
    69
1493
           adp_wgt = torch.sqrt(denoise_diff.square().sum([1,2,3], keepdims=True
     70
1494
           ) + phuber_c**2)
1495 71
           weight = compute_G_weight(sigma, denoise_diff)
1496 72
1497 <sup>73</sup>
            # compute score regularization loss
           batch_loss = weight * denoise_diff * (fake_D_yn - D_yn)/adp_wgt
    74
1498
1499 76
            # compute explicit reward loss if needed
1500 77
           if use_reward:
                reward_loss = -reward_scale * r(fake_x0, prompt)
1501 78
                batch_loss += reward_loss
1502 <sup>79</sup>
1503 <sup>80</sup>
    81
            # compute cfg reward loss if needed
1504 82
           if use_cfg:
1505 83
                cfg_reward_loss = cfg_scale * cfg_vector*fake_x0
1506 84
                batch_loss += cfg_reward_loss
1507 <sup>85</sup>
           batch_loss = batch_loss.sum([1,2,3]).mean()
    86
1508
     87
1509 <sub>88</sub>
           optimizer_G.zero_grad()
1510 89
           batch_loss.backward()
1511
```

1512	entimizer (eter()
1513	optimizer_G.step()
1514	Listing 2: Pytorch Style Pseudo-code of Diff-Instruct*
1515	
1516	
1517	
1518	
1519	
1520	
1521	
1522	
1523	
1524	
1525	
1526	
1527	
1528	
1529	
1530	
1531	
1532	
1533	
1534	
1535	
1536	
1537	
1538	
1539	
1540	
1541	
1542	
1543	
1544	
1545	
1546	
1547	
1548	
1549	
1550	
1551	
1552	
1553	
1554	
1555	
1556	
1557	
1558	
1559	
1560	
1561	
1562	
1563	
1564	
1565	