

)iff-Instruct: Knowledge From Pre-trained Diffusion Universal Approach for Tra

Background

Diffusion Models DM learns score functions by minimizing:

 $\mathcal{J}(\theta) = \int_0^{-} \mathbb{E}_{p_0(\mathbf{x}_0)p_{0t}(\mathbf{x}_t'|\mathbf{x}_0)} [\lambda(t)||\nabla_{\mathbf{x}_t'} \log p_{0t}(\mathbf{x}_t'|\mathbf{x}_0) - s_{\theta}(\mathbf{x}_t', t)||_2^2] \mathrm{d}t$

Knowledge Transferring of DM

- After training, a DM learns rich knowledge about data distributions, making them valuable assets.
- Moreover, in many applications, high-quality data (such as 3D datasets) is expensive (or even impossible) to obtain;
- V we are motivated to study transferring knowledge of as implicit generative models or generators. pre-trained diffusion models to other generative models, such
- Diff-Instruct is tailored for such a knowledge transfer.

Our Work: The Diff-Instruct

Brief Summary of Diff-Instruct

- Set up a Math foundation of knowledge transferring of DMs;
- SOTA single-step diffusion distillation model;
- Consistently improving pre-trained GAN generator models;

Mathematical Setup

- We have a pre-trained teacher diffusion model $s_{p(t)}(x_t)$ and a samples efficiently through $x_0 = g_{\theta}(z), z \sim p_z$. We assume x_0 is differentiable to parameter θ . student implicit generative model g_{θ} , which can generate data
- Let $p^{(t)}$ denote the underlying distribution of teacher DM, and $q^{(t)}$ the distribution of diffused student distribution.

The goal is to minimize some divergence $\mathcal{D}(q^{(0)}, p^{(0)})$

The Integral Kullback-Leibler divergence.

between two distributions p, q is defined as weighting function $w(t) > 0, t \in [0, T]$, the IKL divergence Definition. Given a forward diffusion process and a proper

> Repeat: 2. up G

Weijian Luo¹, Tianyang Hu², Shifeng Zhang², Jiacheng Sun², Zhenguo Li², Zh Peking University¹, Huawei Noah's Ark Lab²

TL; DR: Unified Framework for Understanding Knowledge Transferring of

$$\mathcal{D}_{I\!K\!L}^{[0,T]}(q,p) \coloneqq \int_{t=0}^{T} w(t) \mathcal{D}_{K\!L}(q^{(t)},p^{(t)}) \mathrm{d}t \coloneqq \int_{t=0}^{T} w(t) \mathbb{E}_{x_t \sim q^{(t)}} \Big[\log \frac{q^{(t)}(x_t)}{p^{(t)}(x_t)}\Big] \mathrm{d}t, \ (1)$$

where $q^{(t)}$ and $p^{(t)}$ denote the marginal densities of the forward diffusion process at time t initialized with $q^{(0)} = q$ and $p^{(0)} = p$ respectively.

Diff-Instruct

 \blacktriangleright The θ gradient of such an objective has a formula The goal is to minimize the IKL divergence between $q^{(0)}$ and $p^{(0)}$

 $\mathsf{Grad}(\theta) = \int_{t=0}^{T} w(t) \mathbb{E}_{z \sim p_{z}, x_{0} = g_{\theta}(z), \atop x_{t} \mid x_{0} \sim p_{t}(x_{t} \mid x_{0})} \left[s_{q(t)}(x_{t}, t) - s_{p(t)}(x_{t}) \right] \frac{\partial x_{t}}{\partial \theta} \mathrm{d}t.$ (2)

But $s_{q(t)}(x_t, t)$ is unknown, so we fine-tune an auxiliary diffusion model $s_{\phi}(x_t, t)$ with data consistently sampled from $q^{(0)}$ to approximate $s_{q^{(t)}}(x_t, t)$.

Diff-Instruct Algorithm:

Input:pre-trained DM $s_{q(t)}$, student g_{θ} , prior distribution p_z , DM s_{ϕ} ;

1. update ϕ using <u>SGD</u> with gradient:

$$\begin{split} \mathsf{rad}(\phi) &= \frac{\partial}{\partial \phi} \int_{t=0}^{T} w(t) \mathbb{E}_{\substack{x_0 = g(\theta), \\ x_t \mid x_0 \sim p_t(x_t \mid x_0)}} \| s_{\phi}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t \mid x_0) \|_2^2 \mathrm{d}t. \\ \mathsf{pdate} \ \theta \ \mathsf{using} \ \mathsf{SGD} \ \mathsf{with} \ \mathsf{the} \ \mathsf{gradient:} \\ \mathsf{rad}(\theta) &= \int_{t=0}^{T} w(t) \mathbb{E}_{x_0 = g(\theta), x_0 = g(\theta), x_0$$

$$\operatorname{\mathsf{Grad}}(\theta) = \int_{t=0}^{\tau} w(t) \mathbb{E}_{\substack{x_0 = g(\theta), \\ x_t \mid x_0 \sim p_t(x_t \mid x_0)}} \left[s_{\phi}(\mathbf{x}_t, t) - s_{\rho(t)}(\mathbf{x}_t) \right] \frac{\partial x_t}{\partial \theta} \mathrm{d}t.$$

Untill Converge.

 \blacktriangleright The only requests g_{θ} is that its generated sample is differentiable to θ ; The algorithm has a diffusion fine-tune step and a student model update step; $s_{p(t)}(\mathbf{x}_t, t)$ is teacher, $s_{\phi}(\mathbf{x}_t, t)$ TA, and g_{θ} student;

Connection to DreamFusion: the generator's output is a Dirac's Delta distribution with learnable

parameters, i.e. $q(\mathbf{x}_0) = \delta_{g(\theta)}(\mathbf{x}_0)^{-a}$

Then the Diff-Instruct gradient formula becomes

$$\operatorname{Grad}(\theta) = \int_{t=0}^{t} w(t) \mathbb{E}_{\substack{x_0 = g(\theta), \\ x_t | x_0 \sim p_t(x_t | x_0)}} \left[\nabla_{x_t} \log p_t(x_t | x_0) - s_{p(t)}(x_t) \right] \frac{\partial x_t}{\partial \theta} d$$
This shows that **DreamFusion is a special case of Diff-Instruct**.

.+

^aWe switch the notation from $g_{\theta}(z)$ to $g(\theta)$ since under the assumptions the generator has no randomness.



Diffusion Models.

Application . . step Diffusion Distillation



Figure 2: Generated samples from one-step generators that are distilled from pre-trained diffusion models on different datasets. *Left*: FFHQ-64 (unconditional); *Mid*: ImageNet-64 (conditional); *Right*: CIFAR-10 (unconditional).

Performance:

 Table 1: Unconditional sample quality on CIFAR

 10 through diffusion generations. * Methods that
 require synthetic data construction for distillation. Methods that require real data for distillation

CT [68] 1-ReFlow (+distill)* [42] 2-ReFlow (+distill)* [42] 3-ReFlow (+distill)* [42] PD [63] CD-L2[†] [68] CD-LPIPS[†] [68] DPM-solver-2 [43] DPM-solver-3 [43] 3-DEIS [80] UniPC [82] UniPC [82] DDIM [66] DDIM [66] LSGM [70] PFGM [76] EDM [34] TDPM [85] G Diff-Instruct Single Step Denoise KD* [44] PD **Multiple Steps** -ReFlow [42] 168] [63] [68] THOD e Diffusion GAN(T=2) [75] Diffusi (include Diffusion Distill 147 110 35 50 000 FID (↓) 23.228.70 6.18 4.85 5.21 8.34 9.36 8.91 378 4.08 5.58 5.83 2.93 4.17 5.10 4.67 8.23 5.28 6.03 7.90 **3.55** 4.12 1.97 2.10 2.35 8.65 1.13 8.49 9.08 9.01 8.79 8.69 IS (†) 9.48 9.89 **9.80** 9.05 8.85 9.75 9.68 9.46 Table 2: Class-conditional sample quality on CIFAR-10 and ImageNet 64×64 through diffusion generations. *Methods that require synthetic data construction for distillation. [†]Methods that CT [68] CD-L2[†] [68] CD-LPIPS[†] [68] require real data for distillation. EDM-Heun[34] GGDM [72] Single Step EDM [34] Diff-Instruct EDM-Heun [34] EDM-Euler [34] Multiple Steps ADM [13] EDM[34] PD[†] [63] Class-condition METHOD CD[†] [68] PD CT [68]

 Multiple Steps (include Diffusion Distillation)

 EDM [34]
 35
 1 7

 Diff-Instruct Single Steps EDM-Heun [34] EDM [34] N-DDIM [5] 63 ImageNo NFE (↓) 100 79 10 20 20 10 NNN 6464.†D FID (↓) 15.39 13.00 12.10 18.4 11.1 8.95 2.44 17.25 6.23 4.70 15.56 7 2.07 14.8 4.19

Application Improving GAN

4.24

le quality on Models that t CIFAR10 through GAN models. [†] M we implemented. METHOD quality on Models that

 $IS(\uparrow)$

 $IS(\uparrow)$

AN+DGfi GAN [29] 3AN2 [32 3AN2† 3AN2† 3AN2† + 3AN2+A]	AN+DGflow [1] GAN [29] 3AN2 [32] 3AN2 [†] 3AN2 [†] + DI 3AN2 [†] + DI 3AN2+ADA [31] 3AN2+ADA [31]	AN+DGflow [1] 9.35 GAN [29] 9.02 BAN2 [32] 8.32 BAN2 [†] DI 7.56 BAN2+ADA [31] 5.33 BAN2+ADA+Tune [311 2.92
AN [19]	AN [19]	AN [19] 12.42
N [18]	N [18]	N [18] 15.52
AN [51]	AN [51]	AN [51] 21.70