

# DISTRIBUTION BACKTRACKING BUILDS A FASTER CONVERGENCE TRAJECTORY FOR DIFFUSION DISTILLATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Accelerating the sampling speed of diffusion models remains a significant challenge. Recent score distillation methods distill a heavy teacher model into a student generator to achieve one-step generation, which is optimized by calculating the difference between two score functions on the samples generated by the student model. However, there is a score mismatch issue in the early stage of the score distillation process, since existing methods mainly focus on using the endpoint of pre-trained diffusion models as teacher models, overlooking the importance of the convergence trajectory between the student generator and the teacher model. To address this issue, we extend the score distillation process by introducing the entire convergence trajectory of the teacher model and propose **Distribution Backtracking Distillation (DisBack)**. DisBack is composed of two stages: *Degradation Recording* and *Distribution Backtracking*. *Degradation Recording* is designed to obtain the convergence trajectory by recording the degradation path from the pre-trained teacher model to the untrained student generator. The degradation path implicitly represents the intermediate distributions between the teacher and the student, and its reverse can be viewed as the convergence trajectory from the student generator to the teacher model. Then *Distribution Backtracking* trains the student generator to backtrack the intermediate distributions along the path to approximate the convergence trajectory of the teacher model. Extensive experiments show that DisBack achieves faster and better convergence than the existing distillation method and achieves comparable or better generation performance, with an FID score of 1.38 on the ImageNet  $64 \times 64$  dataset. DisBack is easy to implement and can be generalized to existing distillation methods to boost performance.

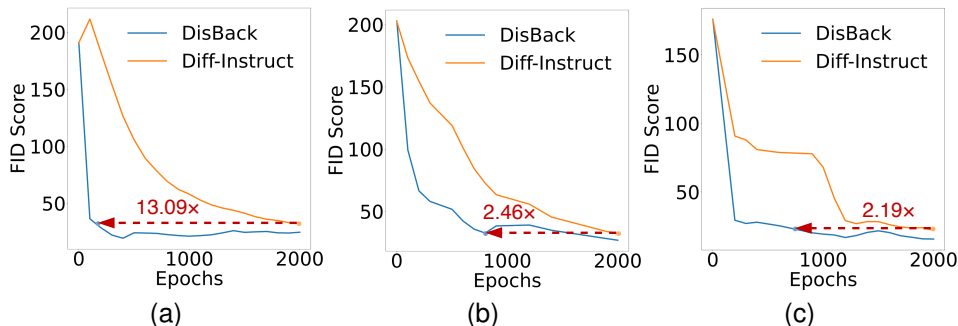


Figure 1: The comparison of the distillation process between existing SOTA score distillation method Diff-Instruct (Luo et al., 2023c) and our proposed DisBack on (a) CIFAR10, (b) FFHQ  $64 \times 64$ , and (c) ImageNet  $64 \times 64$  datasets. DisBack achieves a faster convergence speed due to the constraint of the entire convergence trajectory between the student generator and the teacher model.

# 1 INTRODUCTION

Recently, generative models have demonstrated remarkable performance across diverse domains such as images (Kou et al., 2023; Yin et al., 2024a), audio (Evans et al., 2024; Xing et al., 2024), and videos (Wang et al., 2024; Chen et al., 2024). However, existing models still grapple with the “trilemma” problem, wherein they struggle to simultaneously achieve high generation quality, fast generation speed, and high sample diversity (Xiao et al., 2021). Generative Adversarial Networks (GANs) (Goodfellow et al., 2014) can rapidly produce high-quality samples but often face mode collapse issues. Variational Autoencoders (VAEs) (Kingma & Welling, 2014) offer stable training but tend to yield lower-quality samples. Recently, Diffusion models (DMs) have emerged as a competitive contender in the generative model landscape (Fan et al., 2023; Zhou et al., 2023; Xu et al., 2024). Diffusion models can generate high-quality, diverse samples but still suffer from slow sampling speeds due to iterative network evaluations.

To accelerate the sampling speed, the score distillation method tries to distill a heavy teacher model to a student generator to reduce the sampling cost and achieve the one-step generation (Bao et al., 2023; Luo et al., 2023c; Yin et al., 2024b). The score distillation method optimizes the student generator by calculating the difference between two score functions on the samples generated by the student generator. However, as the generated distribution is far from the training distribution at the beginning, the generated sample lies outside the training data distribution. Thus, the predicted score of the generated sample from the teacher model does not match the sample’s real score in the training distribution. This mismatch issue is reflected by unreliable network predictions of the teacher model, which prevents the student model from receiving accurate guidance and leads to a decline in final generative performance. We identified that this issue arises because existing score distillation methods mainly focus on using the endpoint of the pre-trained diffusion model as the teacher model, overlooking the importance of the convergence trajectory between the student generator and the teacher model. Without the constraint of the convergence trajectory, the mismatch issue causes the student generator to deviate from a reasonable optimization path during training, leading to convergence to suboptimal solutions and a decline in final performance.

To address this problem, we extend the score distillation process by introducing the entire convergence trajectory of the teacher model and propose **Distribution Backtracking Distillation (DisBack)** for a faster and more efficient distillation. The construction of DisBack is based on the following insights. In practice, the convergence trajectory of most teacher models is inaccessible, particularly for large models like Stable Diffusion (Rombach et al., 2022). Because the trajectory of distribution changes is bidirectional, it is possible to construct a degradation path from the teacher model to the initial student generator, and the reverse of this path can be viewed as the convergence trajectory of the teacher model. Compared with fitting the teacher model directly, fitting intermediate targets along the convergence trajectory can mitigate the mismatch issue. Thus, the DisBack incorporates degradation recording and distribution backtracking stages. In the degradation recording stage, the teacher model is tuned to fit the distribution of the initial student generator and obtains a distribution degradation path. The path includes a series of in-between diffusion models to represent the intermediate distributions of the teacher model implicitly. In the distribution backtracking stage, the degradation path is reversed and viewed as the convergence trajectory. Then the student generator is trained to backtrack the intermediate distributions along the path to optimize towards the convergence trajectory of the teacher model. In practice, the degradation recording stage typically requires only a few hundred iterations. Therefore, the proposed method incurs trivial additional computational costs. Compared to the existing score distillation method, DisBack exhibits a significantly increased convergence speed (Fig. 1), and it also delivers superior generation performance (Fig. 2).

Our main contributions are summarized as follows. (1) We extend the score distillation process by introducing the constraint of the entire convergence trajectory of the teacher model and propose Distribution Backtracking Distillation (DisBack), which achieves a faster and more efficient distillation (Sec. 4). (2) Extensive experiments demonstrate that the proposed DisBack accelerates the convergence speed of the score distillation process while achieving comparable or better generation quality compared to existing methods (Sec. 5). (3) The contribution of DisBack is orthogonal to those of other distillation methods. Researchers are encouraged to incorporate our DisBack training strategy into their distillation methods.

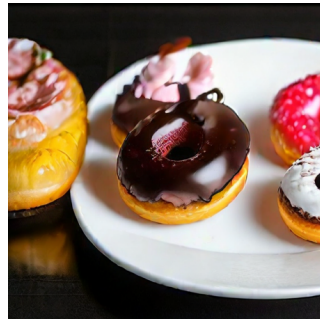
108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161



Wolf in space nebula.



An ocean made of liquid gold, set in a glass bottle, a pirate sailing on a leaf.



Donuts and assorted pastries fill this white plate.



Phoenix emerging from fire with galaxy.



Magical world, valley.



A quantic vintage Futurist, Space rocket air hostess.



A floating island level suspended in the clouds.



The joker walking through streets of New York.



A dog laying on its stomach on a skateboard.

Figure 2: Several examples of  $1024 \times 1024$  images generated by our proposed one-step DisBack model distilled from SDXL (Podell et al., 2024).

## 2 RELATED WORKS

**Efficient diffusion models.** To improve the efficiency of the diffusion model, existing methods use the knowledge distillation method to distill a large teacher diffusion model to a small and efficient student diffusion model (Yang et al., 2022). The progressive distillation model (Salimans & Ho, 2021) progressively distills the entire sampling process into a new diffusion model with half the number of steps iteratively. Building on this, the classifier-guided distillation model (Sun et al., 2023) introduces a dataset-independent classifier to focus the student model on the crucial features to enhance the distillation process. Guided-distillation model (Meng et al., 2023) proposes a classifier-free guiding framework to avoid the computational cost of additional classifiers and achieve high-quality sampling in only 2-4 steps. Recently, the Consistency Model (Song et al., 2023) uses the self-consistency of the ODE generation process to achieve one-step distillation, but this is at the expense of generation quality. To mitigate the surface of the sample quality caused by the acceleration, the Consistency Trajectory Model (Kim et al., 2024) combines the adversarial training and denoising

score matching loss to further improve the performance. Latent Adversarial Diffusion Distillation (Sauer et al., 2024) leverages generative features from pre-trained latent diffusion models to achieve high-resolution, multi-aspect ratio, few-step image generation.

**Score distillation for one-step generation.** Diff-Instruct (Luo et al., 2023c) proposes a distillation method from the pre-trained diffusion model to the one-step generator that involves optimizing the generator by the gradient of the difference between two score functions. One score function represents the pre-trained diffusion distribution, while the other represents the generated distribution. Adversarial Score Distillation (Wei et al., 2024) further employs the paradigm of WGAN and retains an optimizable discriminator to improve performance. Additionally, Swiftbrush (Nguyen & Tran, 2024) leverages score distillation to distill a Stable Diffusion v2.1 into a one-step text-to-image generation model and achieve competitive results. DMD (Yin et al., 2024b) suggests the inclusion of a regression loss between noisy images and corresponding outputs to alleviate instability in the distillation process in text-to-image generation tasks. DMD2 (Yin et al., 2024a) introduces a two-time-scale update rule and an additional GAN loss to address the issue of generation quality being limited by the teacher model in DMD, achieving superior performance. Recently, HyperSD (Ren et al., 2024) integrates score distillation with trajectory segmented consistency distillation and human feedback learning, which achieves SOTA performance from 1 to 8 inference steps.

### 3 PRELIMINARY

In this part, we briefly introduce the score distillation approach. Let  $q_0^G$  and  $q_t^G$  be the distribution of the student generator  $G_{stu}$  and its noisy distribution at timestep  $t$ . In addition,  $q_0$  and  $q_t$  are the training distribution and its noisy distribution at timestep  $t$ . By optimizing the KL divergence in Eq. (1), we can train a student generator to enable one-step generation (Wang et al., 2023).

$$\min_{\eta} D_{KL}(q_0^G(\mathbf{x}_0) \| q_0(\mathbf{x}_0)) \quad (1)$$

Here  $\mathbf{x}_0 = G_{stu}(\mathbf{z}; \eta)$  is the generated samples, and  $\eta$  is the trainable parameter of  $G_{stu}$ . However, due to the complexity of  $q_0$  and its sparsity in high-density regions, directly solving Eq.(1) is challenging (Song & Ermon, 2019). Inspired by Variational Score Distillation (VSD) (Wang et al., 2023), Eq.(1) can be extended to optimization problems at different timesteps  $t$  in Eq. (2). As  $t$  increases, the diffusion distribution becomes closer to a Gaussian distribution.

$$\min_{\eta} \mathbb{E}_{t \sim \mathcal{U}(0,1), \epsilon \sim \mathcal{N}(0,I)} D_{KL}(q_t^G(\mathbf{x}_t) \| q_t(\mathbf{x}_t)) \quad (2)$$

Here  $\mathbf{x}_t$  is the noisy data and  $p(\mathbf{x}_t | \mathbf{x}_0) \sim \mathcal{N}(\mathbf{x}_0, \sigma_t^2 I)$ . Theorem 1 proves that introducing the additional KL-divergence for  $t > 0$  does not affect the global optimum of the original optimization problem in Eq.(1).

**Theorem 1 (The global optimum of training (Wang et al., 2023))** *Given  $t > 0$ , we have,*

$$D_{KL}(q_t^G(\mathbf{x}_t) \| q_t(\mathbf{x}_t)) = 0 \Leftrightarrow D_{KL}(q_0^G(\mathbf{x}_0) \| q_0(\mathbf{x}_0)) = 0 \quad (3)$$

Therefore, by minimizing the KL divergence in Eq. (2), the student generator can be optimized through the following gradients:

$$\nabla_{\eta} D_{KL}(q_t^G(\mathbf{x}_t) \| q_t(\mathbf{x}_t)) = \mathbb{E}_{t, \epsilon} \left[ \left[ \nabla_{\mathbf{x}_t} \log q_t^G(\mathbf{x}_t) - \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t) \right] \frac{\partial \mathbf{x}_t}{\partial \eta} \right] \quad (4)$$

Here the score of perturbed training data  $\nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t)$  can be approximated by a pre-trained diffusion model  $s_{\theta}$ . The score of perturbed generated data  $\nabla_{\mathbf{x}_t} \log q_t^G(\mathbf{x}_t)$  is estimated by another diffusion model  $s_{\phi}$ , which is optimized by score matching with generated data (Song et al., 2021b):

$$\min_{\phi} \mathbb{E}_{t, \epsilon} \left\| s_{\phi}(\mathbf{x}_t, t) - \frac{\mathbf{x}_0 - \mathbf{x}_t}{\sigma_t^2} \right\|_2^2 \quad (5)$$

Thus, the gradient of student generator in Eq.(4) is estimated as

$$\nabla_{\eta} D_{KL}(q_t^G(\mathbf{x}_t) \| q_t(\mathbf{x}_t)) \approx \mathbb{E}_{t, \epsilon} [s_{\phi}(\mathbf{x}_t, t) - s_{\theta}(\mathbf{x}_t, t)] \frac{\partial \mathbf{x}_t}{\partial \eta} \quad (6)$$

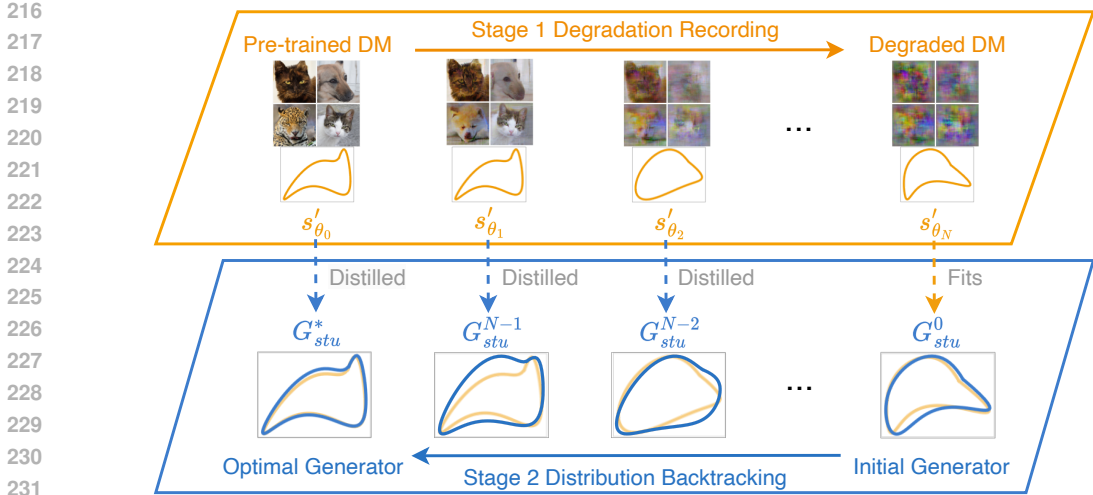


Figure 3: The overall framework of DisBack. Stage 1: An auxiliary diffusion model is initialized with the teacher model  $s_{\theta}$  and then fits the distribution of the initial student generator  $G_{stu}^0$ . The intermediate checkpoints  $\{s'_{\theta_i} \mid i = 0, \dots, N\}$  are saved to form a degradation path. The degradation path is then reversed and viewed as the convergence trajectory. Stage 2: The intermediate node  $s_{\theta_i}$  along the convergence trajectory is distilled to the student generator sequentially until the generator converges to the distribution of the teacher model.

The distribution of the student generator changes after its update. Therefore,  $s_{\phi}$  also needs to be optimized based on the newly generated images to ensure the timely approximation of the generated distribution. Thus, the student generator and  $s_{\phi}$  are optimized alternately.

In practice,  $s_{\phi}$  has three initialization strategies: (1)  $s_{\phi}$  is randomly initialized (Franceschi et al., 2023). (2)  $s_{\phi}$  is initialized as  $s_{\theta}$  or its LoRA (Hu et al., 2021; Wei et al., 2024). (3)  $s_{\phi}$  is initialized by fitting the generated samples of student generator (Luo et al., 2023c). Beyond unconditional image generation (Ye & Liu, 2023), this method has also been applied to tasks such as text-to-image and image-to-image generation across various structures (Yin et al., 2024b; Hertz et al., 2023).

## 4 METHOD

### 4.1 INSIGHT

In this section, we introduce the **Distribution Backtracking Distillation (DisBack)**. The key insight behind DisBack is the importance of the convergence trajectory. As mentioned in Sec.3, there are two score functions in score distillation, one representing the pre-trained diffusion distribution and the other representing the generated distribution. The student model is optimized using the gradient of the difference between these two score functions. Existing methods (Luo et al., 2023c; Yin et al., 2024b;a) directly use the endpoint of the pre-trained diffusion model as the teacher model, overlooking the intermediate convergence trajectory between the student generator and the teacher model. The resulting score mismatch issue between the predicted scores of the generated sample from the teacher model and the real scores causes the student model to receive inaccurate guidance. It ultimately leads to a decline in final performance. Constraining the convergence trajectory between the student generator and the teacher model during the distillation process can mitigate the mismatch issue and help the student generator approximate the convergence trajectory of teacher models to achieve faster convergence. In practice, it is infeasible to obtain the convergence trajectory of most teacher models, especially for large models such as Stable Diffusion (Rombach et al., 2022). Reversely, it is possible to obtain the degradation path from the teacher model to the initial student generator. The reverse of this degradation path can be viewed as the convergence trajectory of the teacher model. Based on the above insights, we structure the proposed DisBack in two stages including the degradation recording stage and the distribution backtracking stage (Fig. 3).

## 4.2 DEGRADATION RECORDING

This stage aims to obtain the degradation path from the teacher model to the initial student generator. The degradation path is then reversed and viewed as the convergence trajectory of the teacher model. The teacher model here is the pre-trained diffusion model  $s_\theta$  and the student generator is represented by  $G_{stu}^0$ .

Let  $s'_\theta$  be a diffusion model initialized with the teacher model  $s_\theta$ , and it is trained on generated samples to fit the initial student generator’s distribution  $q_0^G$  with Eq. (7). By saving the multiple intermediate checkpoints during the training, we can obtain a series of diffusion models  $\{s'_{\theta_i} \mid i = 0, \dots, N\}$ , where  $s'_{\theta_0} = s_\theta \approx q_0$  and  $s'_{\theta_N} \approx q_0^G$ . These diffusion models describe the scores of non-existent distributions on the path, recording how the training distribution  $q_0$  degrades to the initial generated distribution  $q_0^G$ . Algorithm 1 shows the process of obtaining the degradation path. Since distribution degradation is easily achievable, the degradation recording stage only needs trivial additional computational resources (200 iterations in most cases).

$$\min_{\theta} \mathbb{E}_{t, \epsilon} \left\| s'_\theta(\mathbf{x}_t, t) - \frac{\mathbf{x}_0 - \mathbf{x}_t}{\sigma_t^2} \right\|_2^2 \quad (7)$$

---

### Algorithm 1 Degradation Recording.

---

**Input:** Initial student generator  $G_{stu}^0$  and pre-trained diffusion model  $s_\theta$ .

**Output:** Degradation path checkpoints  $\{s'_{\theta_i} \mid i = 0, \dots, N\}$   
 $s'_\theta \leftarrow s_\theta$

**while** not converge **do**

$\mathbf{x}_0 = G_{stu}^0(\mathbf{z}; \eta)$

$\mathbf{x}_t = \mathbf{x}_0 + \sigma_t \epsilon$

Update  $\theta$  with gradient

$$\frac{\partial}{\partial \theta_i} \mathbb{E}_{t, \epsilon} \left\| s'_\theta(\mathbf{x}_t, t) - \frac{\mathbf{x}_0 - \mathbf{x}_t}{\sigma_t^2} \right\|_2^2$$

Save intermediate checkpoints  $s'_{\theta_i}$

**end while**

---



---

### Algorithm 2 Distribution Backtracking.

---

**Input:** Initial student generator  $G_{stu}^0$  and reverse path checkpoints  $\{s'_{\theta_i} \mid i = N, \dots, 0\}$

**Output:** One-step generator  $G_{stu}^*$

$s_\phi \leftarrow s'_{\theta_N}$

**for**  $i \leftarrow N - 1$  to 0 **do**

**while** not converge **do**

$\mathbf{x}_0 = G_{stu}^0(\mathbf{z}; \eta)$

$\mathbf{x}_t = \mathbf{x}_0 + \sigma_t \epsilon$

Update  $\eta$  with gradient

$$\mathbb{E}_{t, \epsilon} [s_\phi(\mathbf{x}_t, t) - s'_{\theta_i}(\mathbf{x}_t, t)] \frac{\partial \mathbf{x}_t}{\partial \eta}$$

Update  $\phi$  with gradient

$$\frac{\partial}{\partial \phi} \mathbb{E}_{t, \epsilon} \left\| s_\phi(\mathbf{x}_t, t) - \frac{\mathbf{x}_0 - \mathbf{x}_t}{\sigma_t^2} \right\|_2^2$$

**end while**

**end for**

---

## 4.3 DISTRIBUTION BACKTRACKING

Given the degradation path from the teacher model to the initial student generator, the reverse path is viewed as a representation of the convergence trajectory between the initial student generator  $G_{stu}^0$  and the teacher model  $s_\theta$ . The key to the distribution backtracking is to sequentially distill checkpoints in the convergence trajectory into the student generator. The last node  $s'_{\theta_N}$  in the path is close to the initially generated distribution  $q_0^G$ . Therefore, in the distribution backtracking stage, we use  $s'_{\theta_{N-1}}$  as the first target to distill the student generator. When near convergence, we switch the target to  $s'_{\theta_{N-2}}$ . The checkpoints  $s'_{\theta_i}$  is sequentially distilled to  $G_{stu}$  until the final target  $s'_{\theta_0}$  is reached. During the distillation, the gradient of  $G_{stu}$  is:

$$\text{Grad}(\eta) = \mathbb{E}_{t, \epsilon} [s_\phi(\mathbf{x}_t, t) - s'_{\theta_i}(\mathbf{x}_t, t)] \frac{\partial \mathbf{x}_t}{\partial \eta} \quad (8)$$

In this stage,  $G_{stu}$  and  $s_\phi$  are also optimized alternately and the optimization of  $s_\phi$  is the same as in the original score distillation (Eq. 5). Compared to existing score distillation methods, the final target of DisBack is the same while constraining the convergence trajectory to achieve more efficient distillation of the student generator. Algorithm 2 summarizes the distribution backtracking stage.

Table 1: The unconditional generation performance of DisBack. The FID ( $\downarrow$ ) scores are shown.

Model	NFE ( $\downarrow$ )	FFHQ	AFHQv2	LSUN-bedroom	LSUN-cat
DDPM (Ho et al., 2020)	1000	3.52	-	4.89	17.10
ADM (Dhariwal & Nichol, 2021)	1000	-	-	1.90	5.57
NCSN++ (Song et al., 2021b)	79	25.95	18.52	-	-
DDPM++ (Song et al., 2021b)	79	3.39	2.58	-	-
EDM (Karras et al., 2022)	79	2.39	1.96	3.57	6.69
EDM (Karras et al., 2022)	15	15.81	13.67	-	-
Diff-Instruct (Luo et al., 2023c)	1	19.93	-	-	-
PD (Salimans & Ho, 2021)	1	-	-	16.92	29.60
CT (Song et al., 2023)	1	-	-	16.00	20.70
CD (Song et al., 2023)	1	12.58	10.75	7.80	11.00
<b>DisBack</b>	<b>1</b>	<b>10.88</b>	<b>9.97</b>	<b>6.99</b>	<b>10.30</b>

## 5 EXPERIMENT

Experiments are conducted on different models across various datasets. We first compare the performance of DisBack with other multi-step diffusion models and distillation methods (Sec. 5.1). Secondly, we compare the convergence speed of DisBack with its variants without the constraint of the convergence trajectory (Sec. 5.2). Thirdly, further experiments are conducted to demonstrate DisBack’s effectiveness in mitigating the score mismatch issues (Sec. 5.3). Then, we also conduct the ablation study to show the effectiveness of introducing the convergence trajectory (Sec. 5.4). Finally, we show the results of DisBack on text-to-image generation tasks (Sec. 5.5).

### 5.1 QUANTITATIVE EVALUATION

DisBack can achieve performance comparable to or even better than the existing diffusion models or distillation methods. Experiments are conducted on different datasets. (1) The unconditional generation on FFHQ 64x64, AFHQv2 64x64, LSUN-bedroom 256x256 and LSUN cat 256x256. (2) The conditional generation on ImageNet 64x64. The performance of DisBack is shown in Tab. 1 and Tab. 2. All the DisBack models are distilled from the pre-trained EDM model (Karras et al., 2022).

For unconditional generation, the one-step generator distilled by the DisBack achieves comparable performance across different datasets compared to multi-step generation diffusion models. Specifically, it outperforms the original EDM model with 15 NFEs (10.88 of DisBack and 15.81 of EDM on FFHQ64). Compared to existing one-step generators and distillation methods, DisBack achieves optimal performance. For conditional generation, the DisBack achieves the best performance compared to the existing models. Moreover, DisBack requires no training data and additional constraints during training. In conclusion, DisBack can achieve competitive distillation performance compared to existing models.

Table 2: The conditional generation performance of DisBack on ImageNet 64x64 dataset.

Model	NFE ( $\downarrow$ )	FID ( $\downarrow$ )
DDPM (Ho et al., 2020)	1000	3.77
DDM (Zhang et al., 2024)	1000	2.11
EDM (Karras et al., 2022)	79	2.44
EDM (Karras et al., 2022)	15	10.46
Moment Matching (Salimans et al., 2024)	8	3.3
SlimFlow (Zhu et al., 2024)	1	12.34
BOOT (Gu et al., 2024)	1	12.30
DDM (Zhang et al., 2024)	1	3.47
CTM (Kim et al., 2024)	1	2.06
Sid (Zhou et al., 2024)	1	1.52
DMD2 (Yin et al., 2024a)	1	1.51
Diff-Instruct (Luo et al., 2023c)	1	5.57
PD (Salimans & Ho, 2021)	1	8.95
CT (Song et al., 2023)	1	13.00
CD (Song et al., 2023)	1	6.20
<b>DisBack</b>	<b>1</b>	<b>1.38</b>

Table 3: Ablation study on constraining the convergence trajectory to the score distillation process. The FID ( $\downarrow$ ) scores in each case are shown.

Model	FFHQ	AHFQv2	ImageNet	LSUN-bedroom	LSUN-cat
DisBack	<b>10.88</b>	<b>9.97</b>	<b>1.38</b>	<b>6.99</b>	<b>10.30</b>
w/o Convergence Trajectory	12.26	10.29	5.96	7.43	10.63

## 5.2 CONVERGENCE SPEED

We conducted a series of experiments to demonstrate the advantages of DisBack in accelerating the convergence speed of the score distillation process on unconditional CIFAR10 (Krizhevsky, 2009), FFHQ 64x64 (Karras et al., 2019), and conditional ImageNet 64x64 (Deng et al., 2009) datasets. Diff-Instruct (Luo et al., 2023c) is the existing SOTA score distillation method, which can be regarded as a variation of DisBack not introducing the convergence trajectory. We compared the FID trends of DisBack and Diff-Instruct during the distillation process in the same situation.

The results are shown in Fig. 1. As for unconditional generation, DisBack achieves a convergence speed 2.46 times faster than the variant without the constraint of convergence trajectory on the FFHQ 64x64 dataset and 13.09 times faster on the CIFAR10 dataset. For the conditional generation on the ImageNet 64x64 dataset, DisBack is 2.19 times faster than the variant without the constraint of convergence trajectory. The fast convergence speed is because constraining the convergence trajectory of the generator provides a clear optimization direction, avoiding the generator falling into suboptimal solutions and enabling faster convergence to the target distribution.

## 5.3 EXPERIMENTS ON SCORE MISMATCH ISSUE

In this part, experiments are conducted to validate the positive impact of constraining the convergence trajectory on mitigating the mismatch issues. We propose a new metric called mismatch degree to assess whether the predicted score of the teacher model matches the distribution’s real score given a data distribution. This score is inspired by the score-matching loss.

$$d_{mis} = \mathbb{E}_{\mathbf{x}_t} \|s_{\theta}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)\|_2 \quad (9)$$

Here  $\mathbf{x}_t$  is the noisy data from the assessed distribution. Besides,  $s_{\theta}(\mathbf{x}_t, t)$  represents the predicted score and  $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)$  represents the real score. Because the real scores are not available in practice, we use Stable Target Field (STF) (Xu et al., 2022) to approximate the real score  $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)$  on the assessed distribution. STF estimation leverages reference batches to reduce the variance of training objectives, which has been proven to yield accurate asymptotically unbiased estimates of the real score (Xu et al., 2022).

When the assessed distribution is close to the distribution of the teacher model  $s_{\theta}$ , the mismatch degree is small, and vice versa. When calculated directly on the training data, the resulting mismatch degree represents the ideal lower bound. Therefore, the mismatch degree can be used to assess the convergence degree of the generated distribution during the distillation process and visualize the convergence speed under the constraint of the convergence trajectory.

We conduct experiments on the FFHQ 64x64 dataset with Diff-Instruct (Luo et al., 2023c) as a baseline. We calculate the mismatch degree on the distribution of the student generator of both Diff-

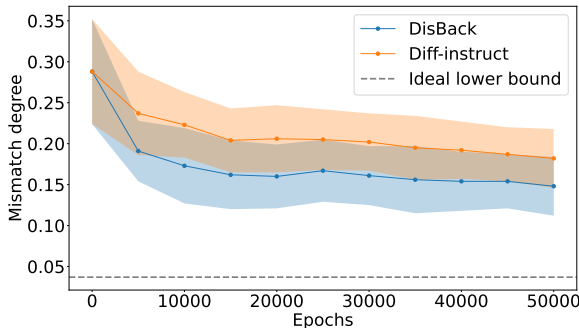


Figure 4: The mismatch degree during the distillation process of Diff-Instruct and proposed DisBack. The standard deviation is visualized. DisBack effectively mitigates the mismatch degree during the entire distillation process.



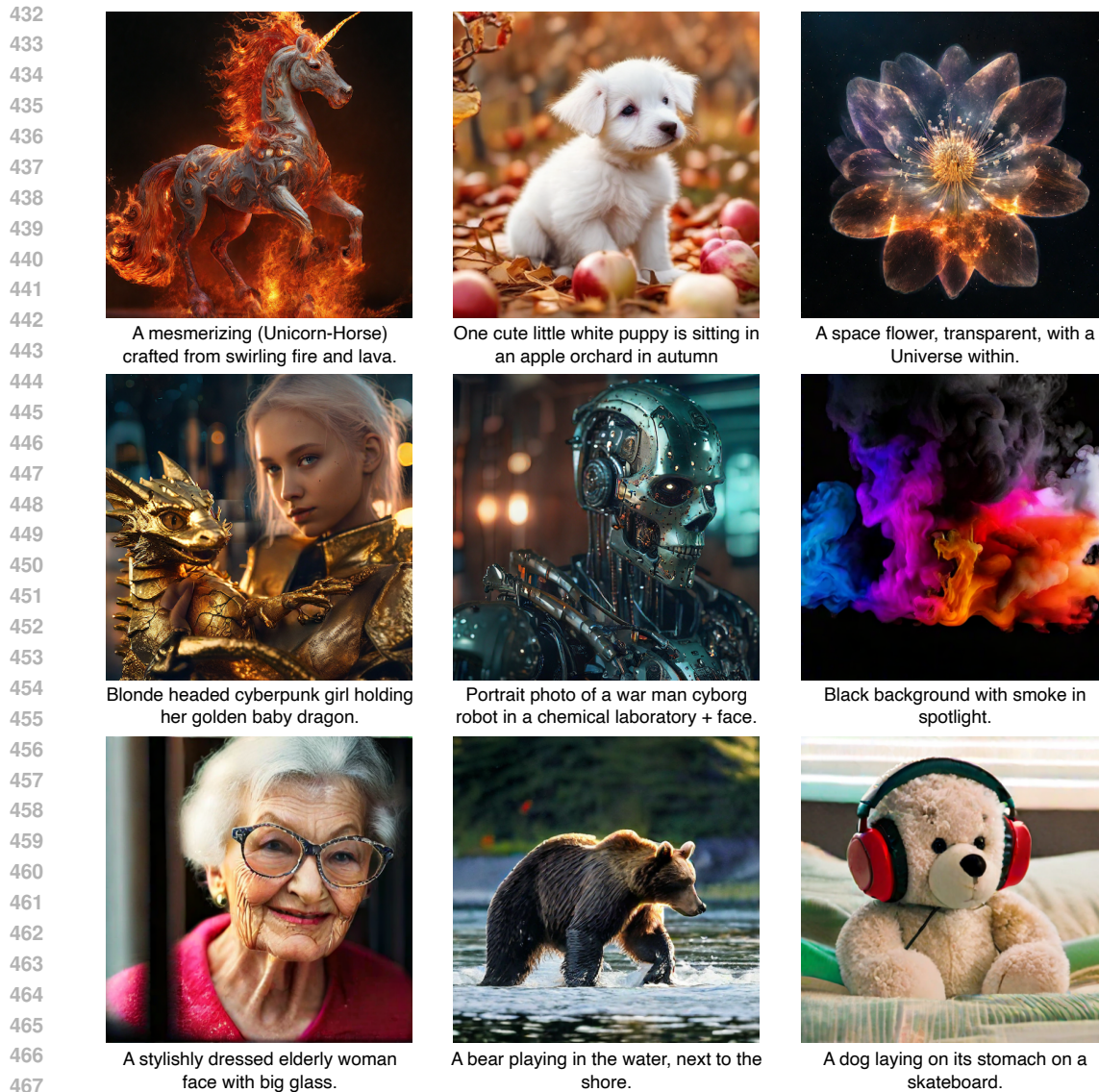


Figure 5: Generation samples by DisBack distilled from SDXL with  $1024 \times 1024$  resolution.

Instruct and the proposed DisBack. The pre-trained EDM model is chosen as the teacher model. In this scenario, the ideal lower bound of the mismatch degree is 0.037. We visualized the mismatch degree in Fig. 4. With the constraining of the convergence trajectory, the mismatch degree of the proposed DisBack is lower during the distillation process, meaning the student generator converges faster and better. Thus, by constraining the convergence trajectory, the mismatch issue can be mitigated and DisBack can achieve more efficient distillation.

#### 5.4 ABLATION STUDY

Ablation studies are conducted to compare the performance of DisBack with its variant without the constraint of the convergence trajectory. The results are shown in Tab. 3. Results show that the variant without the constraint of convergence trajectory suffers from a performance decay in different cases. This confirms the efficacy of constraining the convergence trajectory between the student generator and the teacher model can improve the final performance of the generation.

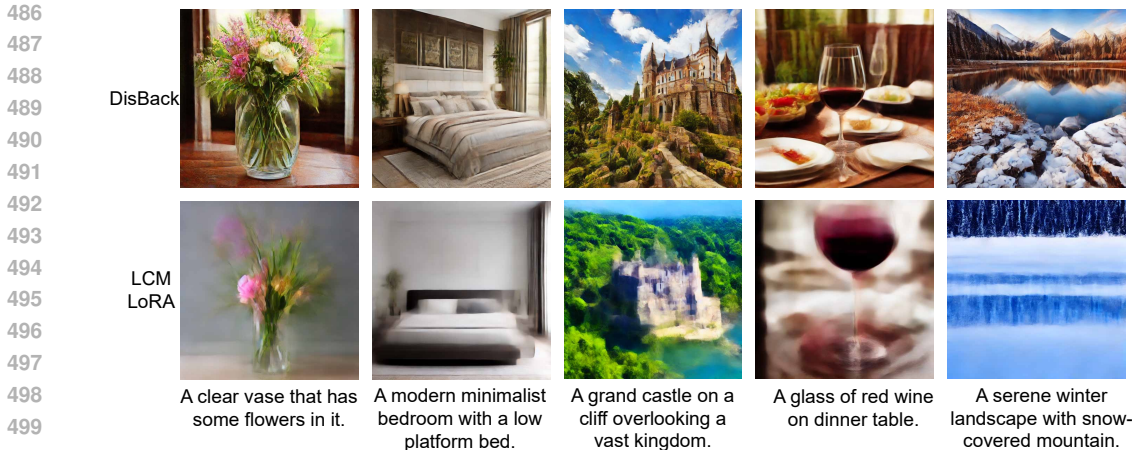


Figure 6: One step generation samples by original LCM-LoRA and its variant distilled from SD v1.5 with DisBack in 512x512. LCM-LoRA with DisBack can generate images with higher quality.

### 5.5 TEXT TO IMAGE GENERATION

Further experiments are conducted on text-to-image generation tasks. We use DisBack to distill the SDXL model (Podell et al., 2024) and evaluate the FID scores of the distilled SDXL and the original SDXL on the COCO 2014 (Radford et al., 2021). The user studies are conducted to verify the effectiveness of DisBack.

Model	FID (↓)	NFE (↓)	User Preference
SDXL	19.36	100	38.7%
DisBack	18.96	1	61.3%

We randomly select 128 prompts from the LAION-Aesthetics (Schuhmann et al., 2022) to generate images and ask volunteer participants to choose the images they think are better. Detailed information about the user study is included in Sec. B.3. The results of the FID evaluation and user study are presented in Tab. 4. DisBack achieved better results in single-step generation compared to the original SDXL with the 100-step DDIM sampler (Song et al., 2021a). The preference scores of DisBack over the original SDXL are 61.3%. Some generation samples are shown in Fig. 2 and Fig. 5.

We also conducted experiments on LCM-LoRA (Luo et al., 2023b). The LCM-LoRA distilled from SDv1.5 using DisBack has an FID score of 36.37 on one-step generation, while the FID score of the original LCM-LoRA is 78.26. Some generated samples of DisBack and original LCM-LoRA are shown in Fig. 6. The details of experiments and results are provided in Sec. A.

## 6 CONCLUSION

**Summary.** This paper proposes Distribution Backtracking Distillation (DisBack) to introduce the entire convergence trajectory of the teacher model in the score distillation. The DisBack can also be used to distill large-scale text-to-image models. DisBack performs a faster and more efficient distillation and achieves a comparable or better performance in one-step generation compared to existing multi-step generation diffusion models and one-step diffusion distillation models.

**Limitation.** The performance of DisBack is inherently limited by the teacher model. The better the original performance of the teacher model, the better the performance of DisBack will also be. Additionally, to achieve optimal performances in both accelerated distillation and generation quality, DisBack requires careful design of the distribution degradation path and the setting of various hyperparameters (such as how many epochs are used to fit each intermediate node in distribution backtracking stage). While with no meticulous design, it can also achieve better performance, further exploration is required to enable the model to reach optimal performance.

## REFERENCES

- 540  
541  
542 Zhiqiang Bao, Zihao Chen, Changdong Wang, Wei-Shi Zheng, Zhenhua Huang, and Yunwen Chen.  
543 Post-distillation via Neural Resuscitation. *IEEE Transactions on Multimedia*, pp. 3046 – 3060,  
544 2023.
- 545 Haoxin Chen, Yong Zhang, Xiaodong Cun, Menghan Xia, Xintao Wang, Chao Weng, and Ying  
546 Shan. Videocrafter2: Overcoming data limitations for high-quality video diffusion models. In  
547 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.  
548 7310–7320, 2024.
- 549  
550 Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. StarGAN v2: Diverse image synthesis  
551 for multiple domains. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
552 *Pattern Recognition*, pp. 8188–8197, 2020.
- 553 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hi-  
554 erarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,  
555 pp. 248–255. Ieee, 2009.
- 556 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat GANs on image synthesis. In  
557 *Advances in Neural Information Processing Systems*, pp. 8780–8794, 2021.
- 558  
559 Zach Evans, CJ Carr, Josiah Taylor, Scott H Hawley, and Jordi Pons. Fast timing-conditioned latent  
560 audio diffusion. *arXiv preprint arXiv:2402.04825*, 2024.
- 561  
562 Wenqi Fan, Chengyi Liu, Yunqing Liu, Jiatong Li, Hang Li, Hui Liu, Jiliang Tang, and Qing Li.  
563 Generative diffusion models on graphs: Methods and applications. In *International Joint Confer-*  
564 *ence on Artificial Intelligence*, pp. 6702–6711, 2023.
- 565 Jean-Yves Franceschi, Mike Gartrell, Ludovic Dos Santos, Thibaut Issenhuth, Emmanuel  
566 de Bézenac, Mickaël Chen, and Alain Rakotomamonjy. Unifying GANs and Score-Based Dif-  
567 fusion as Generative Particle Models. In *Advances in Neural Information Processing Systems*,  
568 volume 36, pp. 59729–59760, 2023.
- 569 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,  
570 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Infor-*  
571 *mation Processing Systems*, pp. 2672–2680, 2014.
- 572  
573 Jiatao Gu, Chen Wang, Shuangfei Zhai, Yizhe Zhang, Lingjie Liu, and Joshua M Susskind. Data-  
574 free Distillation of Diffusion Models with Bootstrapping. In *International Conference on Machine*  
575 *Learning*, 2024.
- 576 Amir Hertz, Kfir Aberman, and Daniel Cohen-Or. Delta denoising score. In *International Confer-*  
577 *ence on Computer Vision*, pp. 2328–2337, 2023.
- 578  
579 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances*  
580 *in Neural Information Processing Systems*, volume 33, pp. 6840–6851, 2020.
- 581 Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, Weizhu Chen,  
582 et al. LoRA: Low-Rank Adaptation of Large Language Models. In *International Conference on*  
583 *Learning Representations*, 2021.
- 584  
585 Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative  
586 adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
587 *Pattern Recognition*, pp. 4401–4410, 2019.
- 588 Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyz-  
589 ing and improving the image quality of StyleGAN. In *Proceedings of the IEEE/CVF Conference*  
590 *on Computer Vision and Pattern Recognition*, pp. 8110–8119, 2020.
- 591  
592 Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and  
593 Timo Aila. Alias-free generative adversarial networks. In *Advances in Neural Information Pro-*  
*cessing Systems*, pp. 852–863, 2021.

- 594 Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-  
595 based generative models. In *Advances in Neural Information Processing Systems*, pp. 26565–  
596 26577, 2022.
- 597  
598 Dongjun Kim, Chieh-Hsin Lai, Wei-Hsiang Liao, Naoki Murata, Yuhta Takida, Toshimitsu Uesaka,  
599 Yutong He, Yuki Mitsufuji, and Stefano Ermon. Consistency trajectory models: Learning proba-  
600 bility flow ode trajectory of diffusion. In *International Conference on Learning Representations*,  
601 2024.
- 602 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *International Conference*  
603 *on Learning Representations*, 2014.
- 604  
605 Ziyi Kou, Shichao Pei, Yijun Tian, and Xiangliang Zhang. Character as pixels: a controllable  
606 prompt adversarial attacking framework for black-box text guided image generation models. In  
607 *International Joint Conference on Artificial Intelligence*, pp. 4912–4920, 2023.
- 608  
609 A Krizhevsky. Learning multiple layers of features from tiny images. *University of Tront*, 2009.
- 610 Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models: Synthe-  
611 sizing high-resolution images with few-step inference. *arXiv preprint arXiv:2310.04378*, 2023a.
- 612  
613 Simian Luo, Yiqin Tan, Suraj Patil, Daniel Gu, Patrick von Platen, Apolinário Passos, Longbo  
614 Huang, Jian Li, and Hang Zhao. Lcm-lora: A universal stable-diffusion acceleration module.  
615 *arXiv preprint arXiv:2311.05556*, 2023b.
- 616 Weijian Luo, Tianyang Hu, Shifeng Zhang, Jiacheng Sun, Zhenguo Li, and Zhihua Zhang. Diff-  
617 Instruct: A Universal Approach for Transferring Knowledge From Pre-trained Diffusion Models.  
618 In *Advances in Neural Information Processing Systems*, volume 36, pp. 76525–76546, 2023c.
- 619  
620 Paul C. Matthews. *Vector Calculus*. Springer undergraduate mathematics series. Springer, London,  
621 1998. ISBN 978-3-540-76180-8. doi: doi.org/10.1007/978-1-4471-0597-8.
- 622  
623 Chenlin Meng, Robin Rombach, Ruiqi Gao, Diederik Kingma, Stefano Ermon, Jonathan Ho, and  
624 Tim Salimans. On distillation of guided diffusion models. In *Proceedings of the IEEE/CVF*  
625 *Conference on Computer Vision and Pattern Recognition*, pp. 14297–14306, 2023.
- 626  
627 Thuan Hoang Nguyen and Anh Tran. Swiftbrush: One-step text-to-image diffusion model with  
628 variational score distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
*and Pattern Recognition*, pp. 7807–7816, 2024.
- 629  
630 Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models.  
631 In *International Conference on Machine Learning*, pp. 8162–8171, 2021.
- 632  
633 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe  
634 Penna, and Robin Rombach. SDXL: Improving latent diffusion models for high-resolution image  
synthesis. In *International Conference on Learning Representations*, 2024.
- 635  
636 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
637 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
638 models from natural language supervision. In *International conference on machine learning*, pp.  
639 8748–8763. PMLR, 2021.
- 640  
641 Yuxi Ren, Xin Xia, Yanzuo Lu, Jiacheng Zhang, Jie Wu, Pan Xie, Xing Wang, and Xuefeng Xiao.  
642 Hyper-sd: Trajectory segmented consistency model for efficient image synthesis. *arXiv preprint*  
*arXiv:2404.13686*, 2024.
- 643  
644 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
645 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Con-*  
646 *ference on Computer Vision and Pattern Recognition*, pp. 10684–10695, 2022.
- 647  
Tim Salimans and Jonathan Ho. Progressive Distillation for Fast Sampling of Diffusion Models. In  
*International Conference on Learning Representations*, 2021.

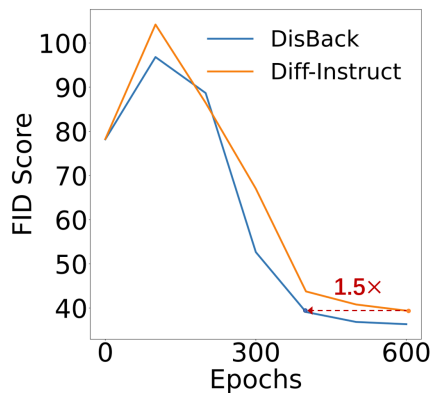
- 648 Tim Salimans, Thomas Mensink, Jonathan Heek, and Emiel Hooeboom. Multistep Distillation of  
649 Diffusion Models via Moment Matching. *arXiv preprint arXiv:2406.04103*, 2024.  
650
- 651 Axel Sauer, Frederic Boesel, Tim Dockhorn, Andreas Blattmann, Patrick Esser, and Robin Rom-  
652 bach. Fast high-resolution image synthesis with latent adversarial diffusion distillation. *arXiv*  
653 *preprint arXiv:2403.12015*, 2024.
- 654 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi  
655 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An  
656 open large-scale dataset for training next generation image-text models. *Advances in Neural*  
657 *Information Processing Systems*, 35:25278–25294, 2022.
- 658 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising Diffusion Implicit Models. In *Inter-*  
659 *national Conference on Learning Representations*, 2021a.  
660
- 661 Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution.  
662 In *Advances in Neural Information Processing Systems*, pp. 11895–11907, 2019.
- 663 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben  
664 Poole. Score-Based Generative Modeling through Stochastic Differential Equations. In *Intern-*  
665 *ational Conference on Learning Representations*, 2021b.  
666
- 667 Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency Models. In *International*  
668 *Conference on Machine Learning*, pp. 32211–32252, 2023.
- 669 Wujie Sun, Defang Chen, Can Wang, Deshi Ye, Yan Feng, and Chun Chen. Accelerating diffusion  
670 sampling with classifier-based feature distillation. In *International Conference on Multimedia*  
671 *and Expo*, pp. 810–815, 2023.
- 672 Fei-Yue Wang, Qinghai Miao, Lingxi Li, Qinghua Ni, Xuan Li, Juanjuan Li, Lili Fan, Yonglin Tian,  
673 and Qing-Long Han. When does sora show: The beginning of tao to imaginative intelligence and  
674 scenarios engineering. *IEEE/CAA Journal of Automatica Sinica*, 11(4):809–815, 2024.  
675
- 676 Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolific-  
677 Dreamer: High-Fidelity and Diverse Text-to-3D Generation with Variational Score Distillation.  
678 In *Advances in Neural Information Processing Systems*, volume 36, pp. 8406–8441, 2023.
- 679 Min Wei, Jingkai Zhou, Junyao Sun, and Xuesong Zhang. Adversarial Score Distillation: When  
680 score distillation meets GAN. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
681 *and Pattern Recognition*, 2024.
- 682 Zhisheng Xiao, Karsten Kreis, and Arash Vahdat. Tackling the Generative Learning Trilemma with  
683 Denoising Diffusion GANs. In *International Conference on Learning Representations*, 2021.  
684
- 685 Yazhou Xing, Yingqing He, Zeyue Tian, Xintao Wang, and Qifeng Chen. Seeing and hearing: Open-  
686 domain visual-audio generation with diffusion latent aligners. In *Proceedings of the IEEE/CVF*  
687 *Conference on Computer Vision and Pattern Recognition*, pp. 7151–7161, 2024.
- 688 Yanwu Xu, Yang Zhao, Zhisheng Xiao, and Tingbo Hou. Ufogen: You forward once large scale  
689 text-to-image generation via diffusion gans. In *Proceedings of the IEEE/CVF Conference on*  
690 *Computer Vision and Pattern Recognition*, pp. 8196–8206, 2024.  
691
- 692 Yilun Xu, Shangyuan Tong, and Tommi S Jaakkola. Stable Target Field for Reduced Variance Score  
693 Estimation in Diffusion Models. In *International Conference on Learning Representations*, 2022.
- 694 Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Wentao Zhang,  
695 Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of methods and  
696 applications. *ACM Computing Surveys*, pp. 1–39, 2022.
- 697 Senmao Ye and Fei Liu. Score Mismatching for Generative Modeling. *arXiv preprint*  
698 *arXiv:2309.11043*, 2023.  
699
- 700 Mingxuan Yi, Zhanxing Zhu, and Song Liu. MonoFlow: Rethinking divergence GANs via the  
701 perspective of Wasserstein gradient flows. In *International Conference on Machine Learning*, pp.  
39984–40000, 2023.

- 702 Tianwei Yin, Michaël Gharbi, Taesung Park, Richard Zhang, Eli Shechtman, Fredo Durand, and  
 703 William T Freeman. Improved Distribution Matching Distillation for Fast Image Synthesis. *arXiv*  
 704 *preprint arXiv:2405.14867*, 2024a.
- 705 Tianwei Yin, Michaël Gharbi, Richard Zhang, Eli Shechtman, Fredo Durand, William T Freeman,  
 706 and Taesung Park. One-step diffusion with distribution matching distillation. In *Proceedings of*  
 707 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6613–6623, 2024b.
- 709 Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. Lsun:  
 710 Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv*  
 711 *preprint arXiv:1506.03365*, 2015.
- 712 Dan Zhang, Jingjing Wang, and Feng Luo. Directly Denoising Diffusion Model. *arXiv preprint*  
 713 *arXiv:2405.13540*, 2024.
- 715 Dewei Zhou, Zongxin Yang, and Yi Yang. Pyramid Diffusion Models For Low-light Image En-  
 716 hancement. In *International Joint Conference on Artificial Intelligence*, pp. 1795–1803, 2023.
- 717 Mingyuan Zhou, Huangjie Zheng, Zhendong Wang, Mingzhang Yin, and Hai Huang. Score identity  
 718 distillation: Exponentially fast distillation of pretrained diffusion models for one-step generation.  
 719 In *International Conference on Machine Learning*, 2024.
- 721 Yuanzhi Zhu, Xingchao Liu, and Qiang Liu. SlimFlow: Training Smaller One-Step Diffusion Mod-  
 722 els with Rectified Flow. *arXiv preprint arXiv:2407.12718*, 2024.

## 724 A MORE EXPERIMENT RESULTS ON TEXT-TO-IMAGE GENERATION

726 We conducted the experiment on LCM-LoRA (Luo et al.,  
 727 2023b). LCM-LoRA is a Low-Rank Adaptation (LoRA)  
 728 version of the Latent Consistency Model (LCM) Luo  
 729 et al. (2023a), applicable across fine-tuned Stable Diffu-  
 730 sion models for high-quality, single-step or few-step gen-  
 731 eration. In this experiment, we use LCM-LoRA as the  
 732 student generator and Stable Diffusion v1.5 as the teacher  
 733 model. We observed that the score distillation under-  
 734 performs when LCM-LoRA serves as the teacher model.  
 735 This issue likely stems from the infeasibility of directly  
 736 converting the outputs of LCM-LoRA into scores.

737 We distill the LCM-LoRA with the proposed DisBack  
 738 and evaluate the FID scores on the COCO 2014 dataset  
 739 (Radford et al., 2021) with the resolution of 512×512.  
 740 50,000 real images and 30,000 generated images were  
 741 used to calculate FID scores. The 30,000 generated im-  
 742 ages were obtained by generating one image for each of  
 743 the 30,000 distinct prompts. In the case of one-step gen-  
 744 eration, the original LCM-LoRA has an FID score of 78.26,  
 745 while the DisBack achieves an FID of 36.37. The change  
 746 in FID scores over training steps is illustrated in Fig. 7,  
 747 showing that DisBack achieves a 1.5 times acceleration  
 748 in convergence speed and yields superior generation performance within the same training period.



749 Figure 7: The FID scores of LCM-  
 750 LoRA distilled from SD1.5 across training  
 751 steps. DisBack achieves faster conver-  
 752 gence and better performance.

## 753 B IMPLEMENTATION DETAILS

### 754 B.1 DATASET SETUP

755 We experiment on the following datasets:

The FFHQ (Flickr-Faces-HQ) dataset (Karras et al., 2019) is a high-resolution dataset of human face images used for face generation tasks. It includes high-definition face images of various ages,

756 genders, skin tones, and expressions from the Flickr platform. This dataset is commonly employed  
 757 to train large-scale generative models. In this paper, we utilize a derivative dataset of the FFHQ  
 758 called FFHQ64, which involves downsampling the images from the original FFHQ dataset to a  
 759 resolution of  $64 \times 64$ .

760 The AFHQv2 (Animal Faces-HQ) dataset (Choi et al., 2020) comprises 15,000 high-definition ani-  
 761 mal face images with a resolution of  $512 \times 512$ , including 5,000 images each for cats, dogs, and wild  
 762 animals. AFHQv2 is commonly employed in tasks such as image-to-image translation and image  
 763 generation. Similar to the FFHQ dataset, we downscale the original AFHQv2 dataset to a resolution  
 764 of  $64 \times 64$  for the experiment.

765 The ImageNet dataset (Deng et al., 2009) was established as a large-scale image dataset to facilitate  
 766 the development of computer vision technologies. This dataset comprises over 14,197,122 images  
 767 spanning more than 20,000 categories, indexed by 21,841 Synsets. In this paper, we use the Image-  
 768 Net64 dataset, a subsampled version of the ImageNet dataset. The Imagenet64 dataset consists  
 769 of a vast collection of images with a resolution of  $64 \times 64$ , containing 1,281,167 training samples,  
 770 50,000 testing samples, and 1,000 labels.

771 The LSUN (Large Scale Scene Understanding) dataset (Yu et al., 2015) is a large-scale dataset for  
 772 scene understanding in visual tasks within deep learning. Encompassing numerous indoor scene im-  
 773 ages, it spans various scenes and perspectives. The LSUN dataset comprises multiple sub-datasets,  
 774 in this study, we use the LSUN Cat and Bedroom sub-datasets with a resolution of  $256 \times 256$ .

## 776 B.2 EXPERIMENT SETUP

777 For experiments on FFHQ  $64 \times 64$ , AFHQv2  $64 \times 64$ , and ImageNet  $64 \times 64$  datasets, the pre-trained  
 778 models are provided by the official release of EDM Karras et al. (2022). We use Adam optimizers  
 779 to train the student generator  $G$  and  $s_\phi$ , with both learning rates set to  $1e^{-5}$ . The training consisted  
 780 of 500,000 epochs on four NVIDIA 3090 GPUs, and the batch size per GPU is set to 8. The training  
 781 ratio between  $s_\phi$  and  $G$  remains at 1 : 1. In the Degradation stage, we trained for 200 epochs  
 782 total, saving a checkpoint every 50 epochs, resulting in a total of 5 intermediate nodes along the  
 783 degradation path  $\{s'_{\theta_i} | i = 0, 1, 2, 3, 4\}$ . In the Distribution Backtracking stage, when  $i \geq 3$ , each  
 784 checkpoint was trained for 1,000 steps. When  $i < 3$ , each checkpoint was trained for 10,000 steps.  
 785 The remaining steps were used to distill the original teacher model  $s'_{\theta_0}$ .

786 For experiments on LSUN bedroom and LSUN cat datasets, the pre-trained EDM models are pro-  
 787 vided by the official release of Consistency Model Song et al. (2023). During the training, we set  
 788  $\sigma_{max}$  to 80 and keep it constant during the single-step generation process. We use SGD and AdamW  
 789 optimizers during training to train the generator  $G$  and  $s_\phi$ , with learning rates set to  $1e^{-3}$  and  $1e^{-4}$ ,  
 790 respectively. The training consisted of 10,000 epochs on one NVIDIA A100 GPU, and the batch  
 791 size per GPU is set to 2. The training ratio between  $s_\phi$  and  $G$  remains at 4 : 1. In the Degradation  
 792 stage, we trained for 200 epochs total and saved the checkpoint every 50 epochs, resulting in a total  
 793 of 5 intermediate nodes along the degradation path  $\{s'_{\theta_i} | i = 0, 1, 2, 3, 4\}$ . In the Distribution Back-  
 794 tracking stage, when  $i \geq 3$ , each checkpoint was trained for 500 steps. When  $i < 3$ , each checkpoint  
 795 was trained for 1000 steps. The remaining steps were used to distill the original teacher model  $s'_{\theta_0}$ .

796 When distilling the SDXL model, the teacher model and the student generator are both initialed  
 797 by the pre-trained SDXL model on the huggingface (model id is 'stabilityai/stable-diffusion-xl-  
 798 base-1.0'). We use Adam optimizers to train  $G$  and  $s_\phi$ , with learning rates set to  $1e^{-3}$  and  $1e^{-2}$ ,  
 799 respectively. The training consisted of 50,000 epochs on one NVIDIA A100 GPU, and the batch  
 800 size per GPU is set to 1. The training ratio between  $s_\phi$  and  $G$  remains at 1 : 1. The training prompts  
 801 are obtained from LAION-Aesthetics. In the Degradation stage, we trained for 1,000 epochs total  
 802 and saved the checkpoint every 100 epochs, resulting in a total of 10 intermediate nodes along the  
 803 degradation path  $\{s'_{\theta_i} | i = 0, 1, \dots, 9\}$ . In the Distribution Backtracking stage, each checkpoint was  
 804 trained for 1,000 steps. The remaining steps were used to distill the original teacher model  $s'_{\theta_0}$ .

## 806 B.3 USER STUDY SETUP

807 Firstly, we randomly selected 128 prompts from the LAION-Aesthetics (Schuhmann et al., 2022).  
 808 Then we use the original SDXL model and the distilled SDXL model to generate 128 pairs of  
 809 images. Subsequently, we randomly recruit 10 volunteers, instructing each to individually evaluate

the fidelity, detail, and vividness of these pairwise images. 10 volunteers included 6 males and 4 females, aged between 24 and 29. 5 of them have artificial intelligence or related majors and the other 5 of them have other majors. They were given unlimited time for the experiment, and all of the volunteers completed the assessment with an average time of 30 minutes. Finally, we took the average of the evaluation results of 10 volunteers as the final user study result.

## C THEORETICAL DEMONSTRATION

### C.1 KL DIVERGENCE OF DISBACK

As the KL divergence follows

$$D_{\text{KL}}(q \parallel p) = \mathbb{E}_q \left[ \log \frac{q}{p} \right] \quad (10)$$

The KL divergence of generated distribution and training distribution at timestep  $t$  can be written as

$$\begin{aligned} D_{\text{KL}}(q_t^G(\mathbf{x}_t) \parallel q_t(\mathbf{x}_t)) &= \mathbb{E}_{\mathbf{x}_t \sim q_t^G(\mathbf{x}_t)} \log \frac{q_t^G(\mathbf{x}_t)}{q_t(\mathbf{x}_t)} \\ &= \mathbb{E}_{\mathbf{x}_0 \sim G(z; \eta)} [\log q_t^G(\mathbf{x}_t) - \log q_t(\mathbf{x}_t)] \\ &= \mathbb{E}_z [\log q_t^G(\mathbf{x}_t) - \log q_t(\mathbf{x}_t)] \end{aligned} \quad (11)$$

Thus, the gradient of KL divergence can be estimated as

$$\nabla_{\eta} D_{\text{KL}}(q_t^G(\mathbf{x}_t) \parallel q_t(\mathbf{x}_t)) = \mathbb{E}_{t, \epsilon} [s_{\phi}(\mathbf{x}_t, t) - s_{\theta}(\mathbf{x}_t, t)] \frac{\delta \mathbf{x}_t}{\delta \eta} \quad (12)$$

### C.2 STABLE TARGET FIELD

Given  $\mathbf{x}_0 \sim q_0$  is the training data,  $\mathbf{x}_t \sim p(\mathbf{x}_t | \mathbf{x}_0)$  is the disturbed data, Xu *et al.* (Xu et al., 2022) presents an estimation of the score as:

$$\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) = \frac{\nabla_{\mathbf{x}_t} p_t(\mathbf{x}_t)}{p_t(\mathbf{x}_t)} = \frac{\mathbb{E}_{\mathbf{x}_0} \nabla_{\mathbf{x}_t} p(\mathbf{x}_t | \mathbf{x}_0)}{p_t(\mathbf{x}_t)} \quad (13)$$

The transition kernel  $p(\mathbf{x}_t | \mathbf{x}_0)$  follows the Gaussian distribution  $p(\mathbf{x}_t | \mathbf{x}_0) \sim \mathcal{N}(\mu_t, \sigma_t^2 I)$ . Here  $\mu_t = \mathbf{x}_0$  in Variance Exploding SDE (Song et al., 2021b) but is defined differently in other diffusion models.

$$p(\mathbf{x}_t | \mathbf{x}_0) = \frac{1}{\sqrt{(2\pi^k)\sigma_t}} \exp\left(-\frac{(\mathbf{x}_t - \mu_t)^T(\mathbf{x}_t - \mu_t)}{2\sigma_t^2}\right) \quad (14)$$

$$\begin{aligned} &\nabla_{\mathbf{x}_t} p(\mathbf{x}_t | \mathbf{x}_0) \\ &= \nabla_{\mathbf{x}_t} \left[ \frac{1}{\sqrt{(2\pi^k)\sigma_t}} \exp\left(-\frac{(\mathbf{x}_t - \mu_t)^T(\mathbf{x}_t - \mu_t)}{2\sigma_t^2}\right) \right] \\ &= p(\mathbf{x}_t | \mathbf{x}_0) \nabla_{\mathbf{x}_t} \left( -\frac{(\mathbf{x}_t - \mu_t)^T(\mathbf{x}_t - \mu_t)}{2\sigma_t^2} \right) \\ &= p(\mathbf{x}_t | \mathbf{x}_0) \frac{\mu_t - \mathbf{x}_t}{\sigma_t^2} \end{aligned} \quad (15)$$

Combine Eq. (13) to Eq. (15), we have

$$\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) = \mathbb{E}_{\mathbf{x}_0} \frac{p(\mathbf{x}_t | \mathbf{x}_0)}{p_t(\mathbf{x}_t)} \frac{\mu_t - \mathbf{x}_t}{\sigma_t^2} = \frac{1}{p_t(\mathbf{x}_t)} \mathbb{E}_{\mathbf{x}_0} p(\mathbf{x}_t | \mathbf{x}_0) \frac{\mu_t - \mathbf{x}_t}{\sigma_t^2} \quad (16)$$

Let  $B$  be a set of reference samples for Monte Carlo estimation, we have

$$p_t(\mathbf{x}_t) = \mathbb{E}_{\mathbf{x}_0} p(\mathbf{x}_t | \mathbf{x}_0) \approx \frac{1}{|B|} \sum_{\mathbf{x}_0^{(i)} \in B} p(\mathbf{x}_t | \mathbf{x}_0^{(i)}) \quad (17)$$



Combine the Eq. (16) and Eq. (17), we can get

$$\nabla_{\mathbf{x}_t} \log p_t(x_t) = \mathbb{E}_{\mathbf{x}_0} \frac{p(\mathbf{x}_t | \mathbf{x}_0) \mu_t - \mathbf{x}_t}{p_t(\mathbf{x}_t) \sigma_t^2} \approx \frac{1}{p_t(\mathbf{x}_t)} \frac{1}{|B|} \sum_{\mathbf{x}_0^{(i)} \in B} p(\mathbf{x}_t | \mathbf{x}_0^{(i)}) \frac{\mu_t - \mathbf{x}_t}{\sigma_t^2} \quad (18)$$

Here the “ $\approx$ ” represents the Monte Carlo estimate.

Depending on the network prediction, the diffusion model can be divided into different types, including  $\epsilon$  prediction (Karras et al., 2022) and  $x_0$  prediction (Song et al., 2021a; Ho et al., 2020; Nichol & Dhariwal, 2021). When the score  $\nabla_{\mathbf{x}_t} \log p_t(x_t)$  is estimated by Eq.(18), it can be converted to  $\epsilon$ ,  $x_0$  and  $v$  by a series of transformations.

$$\hat{\epsilon} \approx -\sigma_t \nabla_{\mathbf{x}_t} \log p_t(x_t) \quad (19)$$

$$\hat{\mathbf{x}}_0 \approx \nabla_{\mathbf{x}_t} \log p_t(x_t) * \sigma_t^2 + \mathbf{x}_t \quad (20)$$

## D DISCUSSION

### D.1 TRAINING EFFICIENCY OF DISBACK

While DisBack involves an iterative optimization process during training, the optimization objective of  $s_\phi(\mathbf{x}_t, t)$  aims to minimize the loss of the standard diffusion model based on Eq.(21), and the objective of student generator aims to minimize the KL divergence in Eq.(22). These two optimization processes do not entail adversarial training as in GANs. Consequently, the optimization process tends to be more stable. A recent work Monoflow (Yi et al., 2023) also discusses in GANs training a vector field is obtained to guide the optimization of the generator, but the vector field derives from the discriminator and the instability is not mitigated.

$$\min_{\phi} \mathbb{E}_{t, \epsilon} \left\| s_\phi(\mathbf{x}_t, t) - \frac{\mathbf{x}_0 - \mathbf{x}_t}{\sigma_t^2} \right\|_2^2 \quad (21)$$

$$\min_{\eta} \mathbb{E}_{t, \epsilon} D_{KL}(q_t^G(\mathbf{x}_t) \| q_t(\mathbf{x}_t)) \quad (22)$$

For the DisBack, training the student generator only requires two U-Nets to perform inference and subtraction. Training  $s_\phi$  only involves training a single U-Net, and gradients do not need to be back-propagated to  $G_{stu}$ . Therefore, these models can be naturally deployed to different devices, making computational resource requirements more distributed. This ease of distribution allows for joint training on computational devices with limited capacity. In contrast, for GANs and VAEs, which require gradient propagation between models (discriminator to generator, decoder to encoder), computational requirements are more centralized, necessitating the use of a single device or tools like DeepSpeed to manage the workload.

### D.2 VECTOR FIELD

In our research, each of the estimated score functions  $s'_{\theta_i}$ , for  $i$  ranging from 0 to  $N$ , delineates a vector field  $\mathbb{R}^{3 \times W \times H} \mapsto \mathbb{R}^{3 \times W \times H}$ . We make a strong assumption behind our proposed method that these score functions represent existing or non-existent distributions and that they altogether imply a transformation path between  $s_\theta$  and the student generator  $G_{stu}^0$ . Nevertheless, a score fundamentally constitutes a gradient field, signifying the gradient of the inherent probability density. A vector field is a gradient field when several conditions are satisfied, including path independence, continuous partial derivatives, and zero curls (Matthews, 1998). The vector field, as characterized by the score functions, may not meet these conditions, and thus there is not a potential function or a probability density function. Such deficiencies could potentially hinder the successful training of the student generator and introduce unforeseen difficulties in the distillation process. Specifically, in instances where  $s_\theta$  does not precisely represent a gradient field, a highly probable scenario considering  $s_\theta$  is a neural network, the samples generated from  $s_\theta$  could encompass failure cases. Although our empirical studies exemplify the effectiveness of the proposed DisBack, the detrimental effects of the discussed issue remain unclear. We will further explore this issue in our future work.

918  
 919  
 920  
 921  
 922  
 923  
 924  
 925  
 926  
 927  
 928  
 929  
 930  
 931  
 932  
 933  
 934  
 935  
 936  
 937  
 938  
 939  
 940  
 941  
 942  
 943  
 944  
 945  
 946  
 947  
 948  
 949  
 950  
 951  
 952  
 953  
 954  
 955  
 956  
 957  
 958  
 959  
 960  
 961  
 962  
 963  
 964  
 965  
 966  
 967  
 968  
 969  
 970  
 971

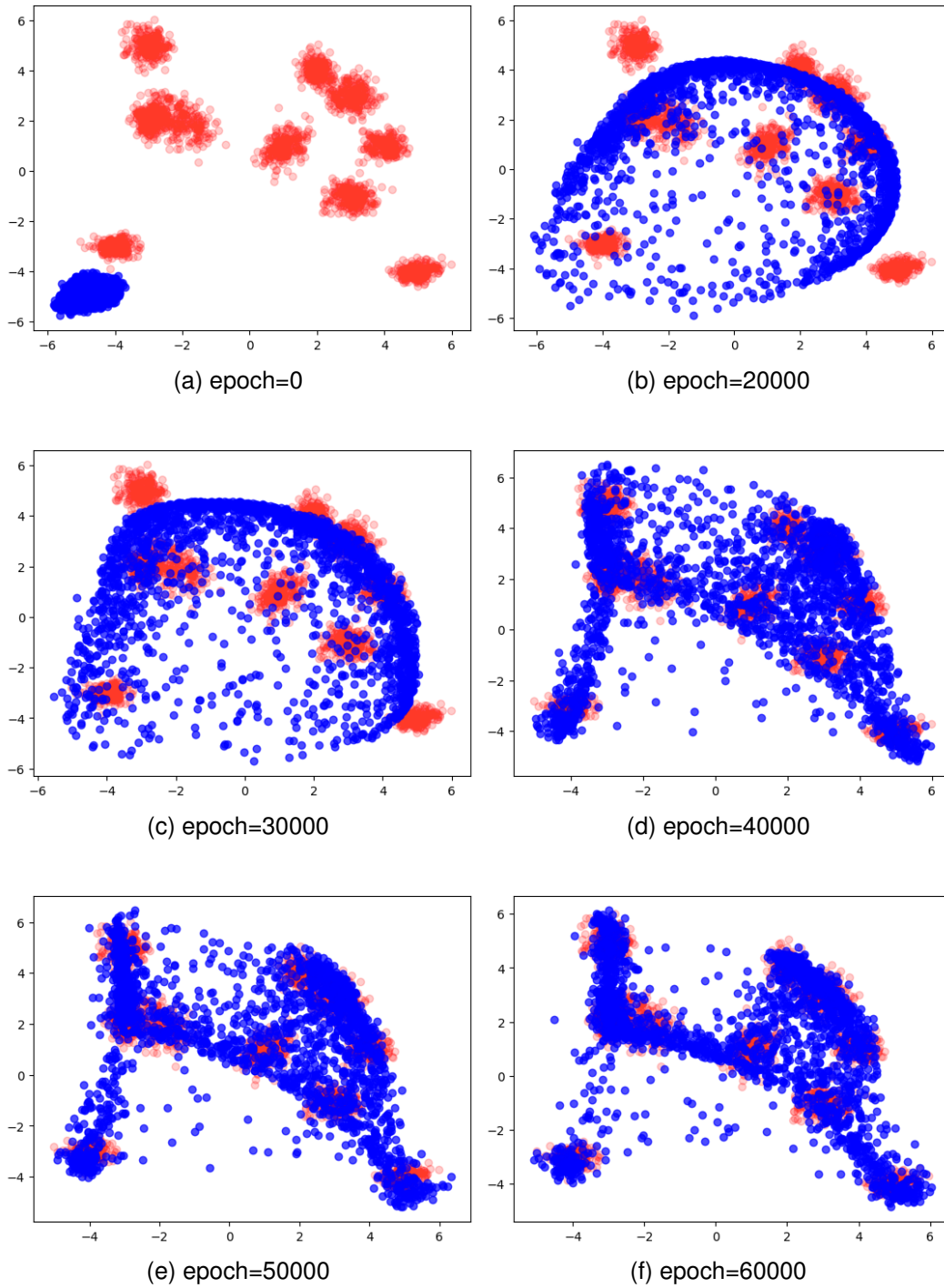


Figure 8: The distribution of student generator during the training process. Blue points visualize the generated distribution  $q_t^G$  and the red points visualize the training distribution  $q_0$ .

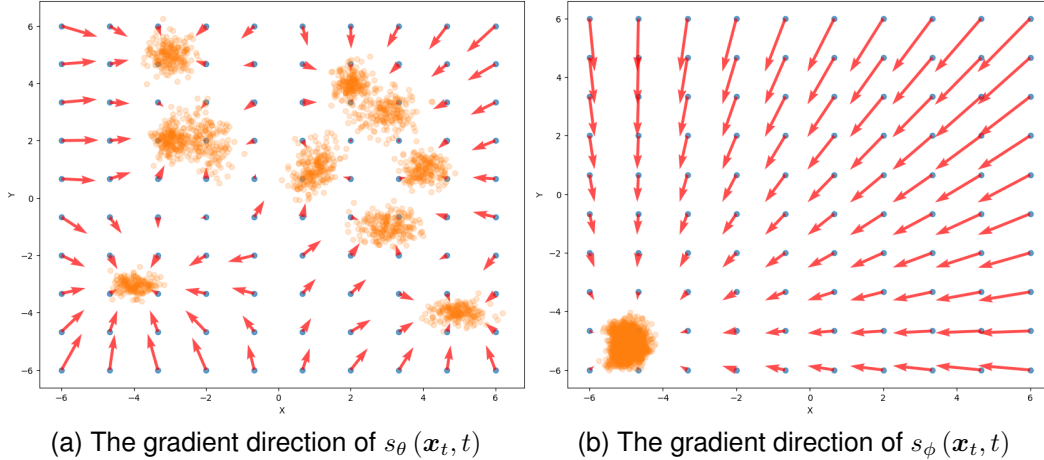


Figure 9: The gradient direction of  $s_\theta(\mathbf{x}_t, t)$  and  $s_\phi(\mathbf{x}_t, t)$  on  $\mathbf{x}_t$ . The points in (a) are sampled from the training distribution and the points in (b) are sampled from the generated distribution.

## E ADDITIONAL DETAILS IN PRE-EXPERIMENTS

### E.1 A TOY EXPERIMENT ON GAUSSIAN MIXTURE DISTRIBUTION

To validate the feasibility of the proposed DisBack, we conduct experiments on two-dimensional Gaussian mixture data. First, we randomly select 10 Gaussian distributions mixed as the training distribution  $q_0$ . Next, we construct a ResNet MLP as the two-dimensional diffusion model  $s_\theta$  and train it using the created mixture Gaussian distribution. Similarly, we construct a simple MLP as the student generator  $G_{stu}$  and train a model  $s_\phi$  with the same architecture as  $s_\theta$  using generated data. Therefore, we can use  $s_\theta$  and  $s_\phi$  to train the student generator  $G_{stu}$ . During the training process, we visualize the distribution of the student generator and training data to intuitively demonstrate the changes in the student generator distribution under the proposed training framework. The distribution of  $G_{stu}$  during the training process is shown in Figure 8. As training progresses, the generated distribution  $q^G$  initially expands outward and then gradually converges towards the training distribution. The results show that the proposed method for training the student generator is effective.

### E.2 GRADIENT ORIENTATION VERIFICATION OF DISBACK

As mentioned in Sec. 3, when updating  $G_{stu}$  using Eq.(6),  $s_\theta(\mathbf{x}_t, t)$  provides a gradient towards the training distribution, while  $s_\phi(\mathbf{x}_t, t)$  provides a gradient toward the generated distribution.

To validate the correctness of these gradient directions, we experiment on two-dimensional data. We evenly sample  $N$  data points within the range of  $(x, y) \in [-6, 6]$  as the noisy data  $\mathbf{x}_t$ . Subsequently, we depict the gradient directions of  $\mathbf{x}_t$  based on  $s_\theta$  and  $s_\phi$  respectively. As shown in Figure 9, consistent with theoretical derivation, for any given  $\mathbf{x}_t$ , the gradient direction of  $s_\theta(\mathbf{x}_t, t)$  points toward the training distribution, and the magnitude of the gradient decreases as the distance to the training distribution decreases. Similarly, for any given  $\mathbf{x}_t$ , the gradient direction of  $s_\phi(\mathbf{x}_t, t)$  points toward the generated distribution.

## F ADDITIONAL SAMPLES FROM DISBACK

We provide additional samples from DisBack on FFHQ  $64 \times 64$  (Figure 11), AFHQv2  $64 \times 64$  (Figure 12), ImageNet  $64 \times 64$  (Figure 13), LSUN Bedroom  $256 \times 256$  (Figure 14) and LSUN Cat  $256 \times 256$  (Figure 15).

## 1026 G FAILURE EXAMPLES

1027 Fig. 10 presents several failure cases of DisBack.

1028 In terms of FFHQ, AFHQv2, and ImageNet, while these images already capture the features of the  
1029 corresponding datasets, the generated results lack accurate and clear backgrounds. The potential  
1030 reasons for this include the fact that these datasets primarily focus on learning foreground content,  
1031 with low requirements for image backgrounds, making the model difficult to clear backgrounds.  
1032

1033 As for LSUN Cat and Bedroom, DisBack successfully generates details such as the cat’s fur and  
1034 the bed’s texture, but it does not generate the overall shape and the detailed structure. This may be  
1035 because the model does not capture the overall information of the data, only capturing local content.  
1036 This issue may stem from the inherent limitations of U-Net, resulting in poor generation of overall  
1037 structures in rare cases.  
1038

1039 In the future, attempts will be made to use more advanced teacher models or improve the distillation  
1040 algorithm to overcome these limitations. Moreover, we will further explore more advanced generator  
1041 architectures such as StyleGAN Karras et al. (2020; 2021) to achieve higher-quality generation.



1057

1058 Figure 10: Failure examples.

## 1061 H ETHICAL STATEMENT

### 1062 H.1 ETHICAL IMPACT

1063 The potential ethical impact of our work is about fairness. As “human face” is included as a kind  
1064 of generated image, our method can be used in face generation tasks. Human-related datasets may  
1065 have data bias related to fairness issues, such as the bias to gender or skin color. Such bias can be  
1066 captured by the generative model in the training.  
1067

### 1068 H.2 NOTIFICATION TO HUMAN SUBJECTS

1069 In our user study, we present the notification to subjects to inform the collection and use of data  
1070 before the experiments.  
1071

1072

1073

1074 Dear volunteers, we would like to thank you for supporting our study. We propose  
1075 the Distribution Backtracking Distillation, which introduces the convergence tra-  
1076 jectory into the score distillation process to achieve efficient and fast distillation  
1077 and high-quality single-step generation.

1078 All information about your participation in the study will appear in the study  
1079 record. All information will be processed and stored according to the local law  
and policy on privacy. Your name will not appear in the final report. Only an

1080 individual number assigned to you is mentioned when referring to the data you  
1081 provided.

1082 We respect your decision whether you want to be a volunteer for the study. If you  
1083 decide to participate in the study, you can sign this informed consent form.  
1084

1085 The Institutional Review Board approved the use of users' data of the main authors' affiliation.  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187



Figure 11: Additional Samples form conditional FFHQ 64x64.

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

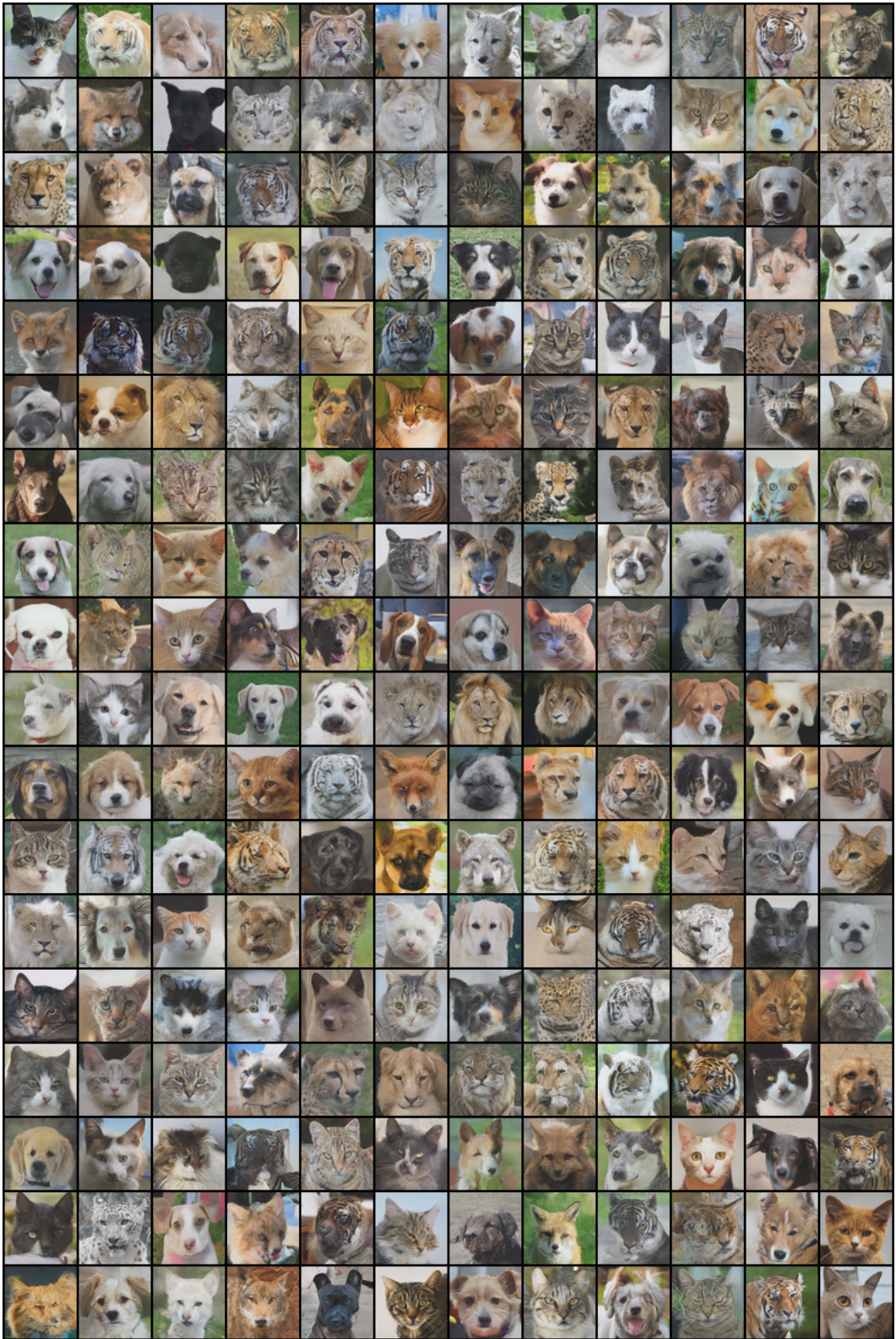


Figure 12: Additional Samples form conditional AFHQv2 64x64.

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

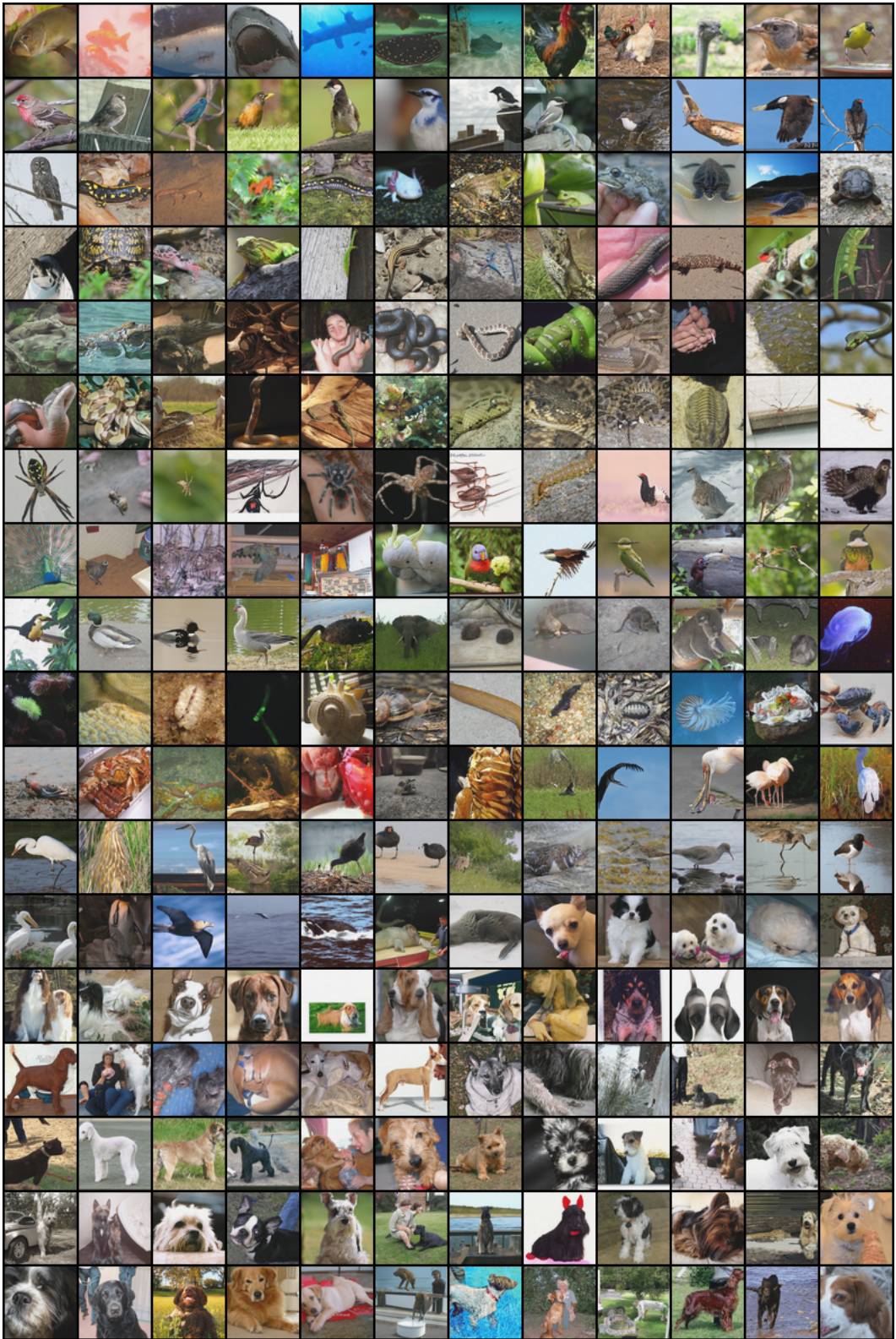


Figure 13: Additional Samples form conditional ImageNet 64x64.



1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349



Figure 14: Additional Samples form conditional LSUN bedroom.

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403



Figure 15: Additional Samples form conditional LSUN cat.