

FAST DIRECT: QUERY-EFFICIENT ONLINE BLACK-BOX GUIDANCE FOR DIFFUSION-MODEL TARGET GENERATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Guided diffusion-model generation is a promising direction for customizing the generation process of a pre-trained diffusion-model to address the specific downstream tasks. Existing guided diffusion models either rely on training of the guidance model with pre-collected datasets or require the objective functions to be differentiable. However, for most real-world tasks, the offline datasets are often unavailable, and their objective functions are often not differentiable, such as image generation with human preferences, molecular generation for drug discovery, and material design. Thus, we need an **online** algorithm capable of collecting data during runtime and supporting a **black-box** objective function. Moreover, the **query efficiency** of the algorithm is also critical because the objective evaluation of the query is often expensive in the real-world scenarios. In this work, we propose a novel and simple algorithm, **Fast Direct**, for query-efficient online black-box target generation. Our Fast Direct builds a pseudo-target on the data manifold to update the noise sequence of the diffusion model with a universal direction, which is promising to perform query-efficient guided generation. Extensive experiments on twelve high-resolution (1024×1024) image target generation tasks and six 3D-molecule target generation tasks show $6\times$ up to $10\times$ query efficiency improvement and $11\times$ up to $44\times$ query efficiency improvement, respectively.

1 INTRODUCTION

Diffusion models have become the state-of-the-art generative model for image synthesis (Ho et al., 2020; Nichol & Dhariwal, 2021; Dhariwal & Nichol, 2021) and video synthesis (Ho et al., 2022), etc. Its remarkable success is due to its powerful capability in modeling the complex multi-mode high-dimensional data distributions.

One promising direction for utilizing the generative power of diffusion models is through target generation. This allows users to customize the generation process to meet specific downstream objectives, effectively extending the models’ capabilities beyond synthesis problems to engineering optimization and science discovery problems, such as image generation with human preferences, drug discovery (Corso et al., 2022; Guan et al., 2023) and material design (Vlassis & Sun, 2023; Giannone et al., 2023).

The pre-trained diffusion models often struggle to generate desired samples for these applications, especially when the target data lies out of the training data distribution. Therefore, the target generation often involves model fine-tuning or guidance techniques. Krishnamoorthy et al. (2023) proposes to train the diffusion model with re-weighted training loss, while Clark et al. (2023); Black et al. (2023); Fan et al. (2024); Yang et al. (2024) advocates for fine-tuning the parameters of the pre-trained model. As oppose to training-time approaches, Bansal et al. (2023) introduces an inference-time approach, which replaces the classifier in classifier guidance (Dhariwal & Nichol, 2021) with a differentiable objective function to achieve the downstream target. However, the requirement of the differentiable objective limits the practical usage of it for real applications with black-box feedback.

Most real-world applications of target generation involve evaluating **black-box objectives** in an **online** manner. For example, image generation with human preferences requires human users to evaluate the generated images; drug design requires real-world experiments to evaluate the generated

054 molecules. These applications typically require several iterative query-feedback cycles with a black-
055 box objective to achieve satisfactory target generation. As objective evaluations are often expensive,
056 it becomes crucial to develop **query-efficient** algorithms to minimize the cost of these evaluations.

057 Although the **online black-box** target generation tasks are important, the existing works are not suit-
058 able to address this task. Bansal et al. (2023); Krishnamoorthy et al. (2023); Lu et al. (2023) require
059 training an offline guidance model with pre-collected data, while Bansal et al. (2023); Clark et al.
060 (2023); Prabhudesai et al. (2023); He et al. (2023) require the objective function to be differentiable.
061 Most recently, Black et al. (2023); Fan et al. (2024); Yang et al. (2024) can be employed for online
062 black-box target generation, but they require online updates of the huge number of the parameters,
063 which is both time-consuming and not query-efficient.

064 In Section 3.1, we first propose a novel guided noise sequence optimization (GNSO) technique
065 to guide the diffusion model sampling process towards a given target. GNSO updates the noise
066 sequence in a *universal direction* on the data manifold. Our GNSO is **efficient**: it enables fast
067 adaptation (with around only 50 steps) from any initial generation towards the given input target.
068 Moreover, GNSO is **robust**: empirically, even the input target image is noisy, it can still generate a
069 clear image semantically similar to the target image. The GNSO itself cannot directly handle black-
070 box target generation tasks, but it provides a backbone for further designing black-box algorithms.

071 In Section 3.2, based on our GNSO technique, we further propose a novel algorithm, **Fast Direct**, to
072 tackle the **online black-box** target generation tasks in a **query-efficient** manner. Fast Direct builds
073 a pseudo-target on the data manifold as the input for our GNSO, to guides the diffusion model at
074 inference time. Notably, our Fast Direct is easy to implement and support any stochastic diffusion
075 scheduler. Moreover, our Fast Direct provides a very flexible framework for extension. Any designs
076 of update methods for the pseudo-target can serve as a plug-in for our Fast Direct.

077 In Section 4, we evaluate our Fast Direct algorithm on the real-world applications: high-resolution
078 (1024×1024) image target generation tasks and 3D-molecule target generation tasks. For image
079 tasks, we employ the black-box API of the modern Large Language Model (LLM), Gemini 1.5,
080 as the black-box objective, which is much more practical compared with the synthetic toy score
081 function used in the literature. Extensive experiments on twelve image target generation tasks and
082 six 3D-molecule target generation tasks show significant query efficiency improvement: **6**-times
083 up to **10**-times query efficiency improvement on image tasks and **11**-times up to **44**-times query
084 efficiency improvement on 3D-molecule tasks, compared with baselines. **Additionally, we evalu-
085 ate Fast Direct on compressibility, incompressibility, and atheistic quality tasks, demonstrating its
086 generalization ability on the unseen prompts.** Our contributions are summarized as follows:

- 087 • We propose a novel guided noise sequence optimization (GNSO) technique to guide the diffu-
088 sion model sampling process toward a given target in an efficient and robust manner.
- 089 • Based on GNSO, we further propose a novel algorithm, Fast Direct, to address online black-box
090 target generation tasks in a query-efficient manner.
- 091 • Fast Direct achieves significant query efficiency improvement in the real-world applications:
092 high-resolution image target generation tasks and 3D-molecule target generation tasks. **Addi-
093 tionally, we demonstrate its generalization capability for unseen prompts.**

095 2 RELATED WORKS

097 2.1 GUIDED GENERATION

099 Diffusion guided generation (also called inference-time guidance) refers to the technique to guide
100 the sampling trajectory of the pre-trained diffusion to generate target data (Croitoru et al., 2023;
101 Chen et al., 2024). The key advantage of this approach is that it does not require updating the model
102 parameters, which is computationally expensive, especially for large models.

103 Dhariwal & Nichol (2021) proposed classifier guidance. However, it requires a guidance model that
104 trained on noisy images with different noise scales, which generally not readily available and often
105 requires training from scratch for each domain. Chung et al. (2022); Bansal et al. (2023); He et al.
106 (2023) extend the classifier guidance by using a differentiable objective function that is defined only
107 for clean data. Instead of using noisy images, the predicted clean images at each sampling steps are
used as the input for the guidance objective function. In this way, the guidance process can operate

108 on the clean image space. While the predicted clean image is naturally imperfect, empirically it still
 109 provide informative feedback to guide image generation (Bansal et al., 2023).

110 However, this approximation can harm image quality. Bansal et al. (2023) proposed universal guid-
 111 ance that comprises of backward guidance followed by a self-recurrence step to preserve image qual-
 112 ity. On the other hand, He et al. (2023) leverages the differentiability of a well-trained auto-encoder
 113 to project the image to the data manifold and thus preserve image quality during the guidance pro-
 114 cess. Instead of guiding the sampling trajectory, Karunratanakul et al. (2024); Eyring et al. (2024)
 115 treat the diffusion process as a black-box and only optimize the initial (prior) noise.

116 While the aforementioned approaches assume a differentiable objective function, Lu et al. (2023)
 117 tackles black-box objective f by learning a differentiable proxy neural network h to match their
 118 gradients, i.e., $\nabla f \approx \nabla h$. Li et al. (2024) eliminates the need for a differentiable proxy model
 119 by employing importance sampling weighted by the objective values during the sampling process.
 120 Most recently, DNO (Tang et al., 2024) proposes optimize the diffusion noise sequence by using
 121 ZO-SGD (Nesterov & Spokoiny, 2017) to tackle black-box objective function. However, it runs at
 122 the instance level; namely, each run only produces one image.

124 2.2 DIFFUSION-MODEL FINE-TUNING

125 Diffusion-model fine-tuning refer to the technique that updating the pre-trained diffusion-model
 126 parameters to improve its performance on a specific use case. In this section, we review the fine-
 127 tuning methods that support online learning of the black-box objective function.

128 Black et al. (2023); Fan et al. (2024) formulate the diffusion fine-tuning problem as a reinforce-
 129 ment learning (RL) problem within Markov Decision Processes (MDPs), and proposes an iterative
 130 algorithm to fine-tune diffusion model by using Proximal Policy Optimization (PPO) (Schulman
 131 et al., 2017; Uehara et al., 2024) loss function. Fan et al. (2024) integrates the PPO with a KL
 132 regularization to prevent the fine-tuned model deviated too much from the pre-trained model.

133 On the other hand, Yang et al. (2024) does not requires an absolute objective values, instead it
 134 uses the relative reward on pair of samples by integrating DPO (Direct Preference Optimization)
 135 (Rafailov et al., 2024), a technique for fine-tuning large language models, into diffusion model.

136 Fine-tuning diffusion model requires large amount of GPU memory. Existing works mitigates the
 137 memory consumption by using LoRA (Low-Rank Adaptation) (Hu et al., 2021) technique, and only
 138 fine-tune parameters of the attention blocks in UNet. We categorize the related works based on
 139 whether they support the online and black-box objective tasks in Appendix F Table 2.

142 3 METHODS

143 In this section, we first present a novel inference-time guidance generation method by guided noise
 144 sequence optimization. Based on this method, we further present our query-efficient online black-
 145 box guidance algorithm, Fast Direct, to address online black-box guidance tasks.

148 3.1 NOISE SEQUENCE OPTIMIZATION WITH TARGET GUIDANCE

149 Take i.i.d. Gaussian samples $\{\epsilon_0, \dots, \epsilon_K\} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, for $k \in \{1, \dots, K\}$, the inference process
 150 of the diffusion model can be formulated as follows:

$$151 \mathbf{x}_k = \mathcal{S}_\theta(\mathbf{x}_{k-1}, \epsilon_k), \quad (1)$$

152 where $\mathcal{S}_\theta(\cdot, \cdot)$ denotes the diffusion model sampler that depends on the concrete SDE solver used,
 153 $\mathbf{x}_0 = \epsilon_0$ is a Gaussian sample, and \mathbf{x}_K denotes the generated data.

154 *Question 1: Given an input target \mathbf{x}^* , can we guide a pre-trained diffusion model at inference time
 155 to generate the target data in an efficient and robust manner?*

156 To answer this question, we need to develop an algorithm to address two challenges simultaneously:
 157 (1) Guidance Efficiency: the algorithm can use only a few iteration steps of guidance update to
 158 generate target data \mathbf{x}^* . (2) Robust to noisy target: given a noisy target (e.g., a noise perturbed
 159 image), the algorithm can generate meaningful target data (e.g., a clear image).

Algorithm 1: Guided Noise Sequence Optimization

Input: Step size α , target data \mathbf{x}^* , number of diffusion sampling steps K , pre-trained diffusion sampler \mathcal{S}_θ , number of iterations T .

Output: Generated target \mathbf{x}_K .

```

1 Take i.i.d. noise  $\{\epsilon_0, \dots, \epsilon_K\} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2 Compute the norm of noise  $\epsilon_k^{\text{norm}} = \|\epsilon_k\|$  for  $k \in \{0, \dots, K\}$ .
3 for  $t \leftarrow 1$  to  $T$  do
4    $\mathbf{x}_0 \leftarrow \epsilon_0$ 
5   for  $k \leftarrow 1$  to  $K$  do
6      $\mathbf{x}_k \leftarrow \mathcal{S}_\theta(\mathbf{x}_{k-1}, \epsilon_k)$ 
7   for  $k \leftarrow 0$  to  $K$  do
8      $\hat{\epsilon}_k \leftarrow \epsilon_k + \alpha(\mathbf{x}^* - \mathbf{x}_K)$ 
9      $\epsilon_k \leftarrow \frac{\hat{\epsilon}_k}{\|\hat{\epsilon}_k\|} \epsilon_k^{\text{norm}}$ 

```

Line 4 - Line 6: Call diffusion-model sampler to generate data \mathbf{x}_K .
 Line 7 - Line 9: Update the noise ϵ_k by moving along direction $\mathbf{x}^* - \mathbf{x}_K$.

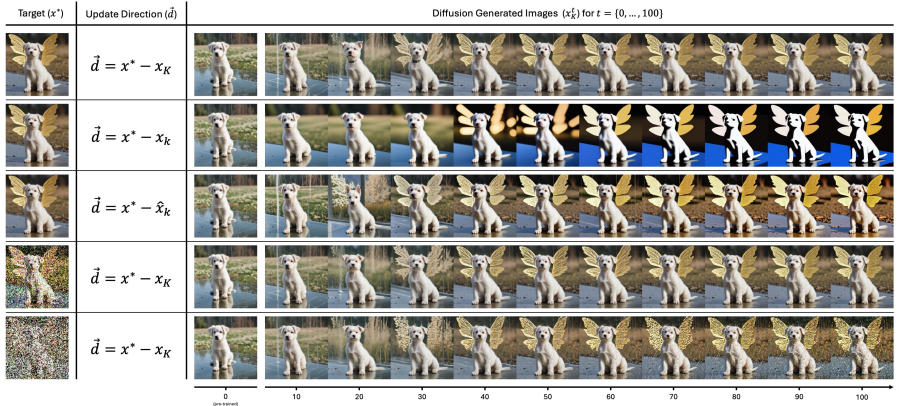


Figure 1: Demonstration of guided generation for a given target by Algorithm 1. Column 2 (Update Direction) indicates the update term of Algorithm 1 Line 8. Row 1 to 3 analyze how different update directions can affect the generated images. Row 4 and 5 show that by using the update direction of $\mathbf{x}^* - \mathbf{x}_K$, diffusion model is able to generate visually satisfied images even when the target image is noisy. The noisy target image (\mathbf{x}^*) of the row 4 is obtained by clean image added with noise $\mathcal{N}(\mathbf{0}, \mathbf{I})$, and row 5 is added with noise $\mathcal{N}(\mathbf{0}, 9 \times \mathbf{I})$.

We propose a novel guided noise sequence optimization method to address the above two challenges, which is presented in Algorithm 1. Our method optimizes the noise sequence by updating noise along direction on data manifold. Specifically, Algorithm 1 consists of two procedures: data generation procedure and noise update procedure. In data generation procedure, Line 4 - Line 6 in Algorithm 1 employ current noise sequence $\{\epsilon_0, \dots, \epsilon_K\}$ to generate a data \mathbf{x}_K . In noise update procedure, Line 7 - Line 9 in Algorithm 1, we update each noise ϵ_k by moving along a universal direction $\mathbf{x}^* - \mathbf{x}_K$. Then, Algorithm 1 continues to employ the updated noise to generate a new data \mathbf{x}_K . By looping the data generation procedure and noise update procedure, Algorithm 1 can generate a clear and meaningful target.

The direction $\mathbf{x}^* - \mathbf{x}_K$ point from the current generation data \mathbf{x}_K to the target \mathbf{x}^* . Intuitively, moving in this direction will shift the diffusion sampling trajectory toward the target. We highlight the advantage of updating noise with the *universal direction* $\mathbf{x}^* - \mathbf{x}_K$ compared with direction $\mathbf{x}^* - \mathbf{x}_k$ as the following remark.

Remark: In our Algorithm 1, we guided the noise update by moving along a universal direction $\mathbf{x}^* - \mathbf{x}_K$. This direction lies in the data manifold, which alleviates the quality degeneration of the generated data. We employ the direction $\mathbf{x}^* - \mathbf{x}_K$ instead of the direction $\mathbf{x}^* - \mathbf{x}_k$ because the internal \mathbf{x}_k during inference sampling process may be far away from the data manifold (e.g., the \mathbf{x}_k can be a very noisy and distorted image), which may degenerate the quality.

We empirically evaluate Algorithm 1 with different update directions: $\vec{d} = \mathbf{x}^* - \mathbf{x}_K$ and $\vec{d} = \mathbf{x}^* - \mathbf{x}_k$. Additionally, we consider the approximation $\mathbf{x}_K \approx \hat{\mathbf{x}}_k^1$. Lastly, we set the noisy target to test the robustness. We use the stable-diffusion model as the pre-trained model. The stepsize α is set to $\alpha = 2 \times 10^{-3}$. We present the generated images at every 10-iteration for each cases in Figure 1.

From Figure 1, we can observe that Algorithm 1 (with update direction $\mathbf{x}^* - \mathbf{x}_K$) can successfully guide the diffusion process towards the target image, while $\mathbf{x}^* - \mathbf{x}_k$ leads to a degenerated image. The approximation $\mathbf{x}^* - \hat{\mathbf{x}}_K$ can achieve comparable target image but not as accurate as $\mathbf{x}^* - \mathbf{x}_K$.

Interestingly, even the input target \mathbf{x}^* is perturbed by noise and lies out of the data manifold, Algorithm 1 will generate a pseudo-data on data manifold that is similar to the noisy target \mathbf{x}^* but with much higher quality compared with the input noisy target. This property is important for designing black-box guidance methods because the update direction is an approximation or estimation of the true gradient of the black-box objective.

Moreover, we present the updating direction $\mathbf{x}^* - \mathbf{x}_{K'}$ for $K' \in \{K, K/2, K/4, K/8\}$ in Appendix C. We observe that the generated image quality decreases as the $\mathbf{x}_{K'}$ becomes more noisy.

3.2 QUERY-EFFICIENT ONLINE BLACK-BOX GUIDANCE

In section 3.1, we address the diffusion guidance sampling with a given optimal target input. However, in reality, we often only have access to the black-box objective function $f(\mathbf{x})$, where the optimal target data $\mathbf{x}^* = \arg \min f(\mathbf{x})$ is unknown. In this section, we further investigate the diffusion guidance sampling with only black-box objective feedback. For this task, it is natural to ask the following question.

Question 2: Can we guide a pre-trained diffusion model at inference time to generate target data with only black-box objective feedback in a query-efficient online manner?

To answer this question, we need to address two challenges: (1) Black-box challenge and (2) Online guidance challenge. The black-box challenge means that we cannot access the gradient of the objective. The online guidance challenge means we don't have a prior dataset to train a surrogate (classifier) for guidance; we can only access the black-box objective through query feedback online. Usually, the query evaluation is expensive. Thus, the query efficiency is critical.

We start addressing the above two challenges based on our target guidance generation Algorithm 1. Although Algorithm 1 itself cannot handle black-box guidance tasks because it needs input a target \mathbf{x}^* , it provides a basis for us to design query efficient online black-box guidance methods. To be specific, Algorithm 1 enables us to generate a target even updating with a noisy direction $\mathbf{x}^* - \mathbf{x}_K$. This property is important for black-box guidance tasks because we can use a noisy estimation of the gradient of the black-box objective to guide the generation.

The key idea is to set a pseudo target $\hat{\mathbf{x}}^*$ to guide the generation process based on our Algorithm 1. We present our Fast Direct method as in Algorithm 2. In Algorithm 2, we call Algorithm 1 inside the for-loop w.r.t. the number of batch queries t . The only difference compared with Alg. 1 is that we set a pseudo target $\hat{\mathbf{x}}^*$ in Line 10 of Algorithm 2 instead of a given fixed target \mathbf{x}^* .

The choice of models for updating the pseudo target $\hat{\mathbf{x}}^*$ in Algorithm 2 is flexible, which supports various black-box target generation method designs based on our Fast Direct algorithm framework. In this work, we set the pseudo target $\hat{\mathbf{x}}^*$ through nonparametric methods without additional training. Specifically, we employ two methods for setting the pseudo target $\hat{\mathbf{x}}^*$, namely, Gaussian process (GP) update and historical optimal update.

Set pseudo target $\hat{\mathbf{x}}^*$ by GP update: For the GP update case, we set the pseudo target $\hat{\mathbf{x}}^*$ as the one gradient step update using the gradient of the posterior mean prediction of the GP surrogate model as follows:

$$\hat{\mathbf{x}}^* = \mathbf{x}_K - \nabla \hat{f}(\mathbf{x}_K; \mathbf{X}^n) \quad (2)$$

where \mathbf{x}_K denotes the generated data in Algorithm 2 and $\hat{f}(\mathbf{x}_K; \mathbf{X}^n)$ denotes the posterior prediction of the GP surrogate model (Seeger, 2004) evaluated at \mathbf{x}_K , which has closed-form as below:

$$\hat{f}(\mathbf{x}_K; \mathbf{X}^n) = \mathbf{k}(\mathbf{x}_K, \mathbf{X}^n)^\top (\mathcal{K}(\mathbf{X}^n, \mathbf{X}^n) + \lambda \mathbf{I})^{-1} \mathbf{y} \quad (3)$$

¹The $\hat{\mathbf{x}}_k$ is the "predicted clean image" at step k , see the Eq.12 of DDIM (Song et al., 2020). Note that it is undefined for initial step, so we simply set $\hat{\mathbf{x}}_0 := \hat{\mathbf{x}}_1$.

Algorithm 2: Fast Direct

Input: Max number of batch queries N , batch size B , step-size α , number of diffusion sampling steps K , pre-trained diffusion sampler \mathcal{S}_θ .

Output: Set of optimized samples $\{x_K^1, \dots, x_K^B\}$, pseudo target model that learns from dataset \mathcal{D} .

```

1 Initialize  $\mathcal{D} \leftarrow \{\}$ 
2 for  $i \leftarrow 1$  to  $N$  do
3   do parallel for  $B$  instances
4     Take i.i.d. noise  $\{\epsilon_0, \dots, \epsilon_K\} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ .
5     Compute the norm of noise  $\epsilon_k^{\text{norm}} = \|\epsilon_k\|$  for  $k \in \{0, \dots, K\}$ .
6     for  $t \leftarrow 1$  to  $i$  do
7        $x_0 \leftarrow \epsilon_0$ 
8       for  $k \leftarrow 1$  to  $K$  do
9          $x_k \leftarrow \mathcal{S}_\theta(x_{k-1}, \epsilon_k)$ 
10      Set a pseudo target  $\hat{x}^*$  by the pseudo target model with input  $x_K$ .
11      for  $k \leftarrow 0$  to  $K$  do
12         $\hat{\epsilon}_k \leftarrow \epsilon_k + \alpha(\hat{x}^* - x_K)$ 
13         $\epsilon_k \leftarrow \frac{\hat{\epsilon}_k}{\|\hat{\epsilon}_k\|} \epsilon_k^{\text{norm}}$ 
14      Query black-box objective score  $y \leftarrow f(x_K)$ .
15      Increment dataset  $\mathcal{D} \leftarrow \mathcal{D} \cup \{(x_K, y)\}$ .
16    Update pseudo target model with dataset  $\mathcal{D}$ .
```

where $\mathbf{X}^n = [x^1, \dots, x^n]$ and $\mathbf{y} = [y^1, \dots, y^n]$ denotes the collected data and its corresponding score value in the historical set \mathcal{D} in Algorithm 2, respectively. When the historical set \mathcal{D} is empty, we simply set the gradient $\nabla \hat{f}(x_K; \emptyset) = 0$.

In this work, we employ shift-invariant kernels that can be rewritten in the form as $k(z_1, z_2) = g(\|z_1 - z_2\|_2)$, e.g., Gaussian kernel and Matérn kernel. We show that the gradient of the GP surrogate when employing these shift-invariant kernels lies on the low-dimensional data manifold. We present this property in Proposition 1. Detailed proof of Proposition 1 can be found in Appendix E. This property enables us to generate meaningful data on the data manifold instead of degenerated data, which is important for high-dimensional diffusion-model target generation problems, e.g., high-resolution (1024×1024) target image generation tasks.

Proposition 1. *Given prior data $\mathbf{X}^n = [x^1, \dots, x^n]$ and its corresponding score $\mathbf{y} = [y^1, \dots, y^n]$, let $\hat{f}(x; \mathbf{X}^n)$ denotes the posterior mean of the GP model. For shift-invariant kernels $k(z_1, z_2) = g(\|z_1 - z_2\|_2)$, $\nabla \hat{f}(x; \mathbf{X}^n)$ lies on a subspace spanned by $[x, x^1, \dots, x^n]$, where x^i denotes the i^{th} data sample in \mathbf{X}^n .*

Remark: When setting the pseudo target \hat{x}^* as Eq.(2), we know the pseudo target \hat{x}^* lies on the low-dimensional subspace spanned by data according to Proposition 1. Usually, the gradient descent update in Eq.(2) will lead to a new pseudo target \hat{x}^* with a lower score. In addition, based on the empirical observation from Algorithm 1, we know even the pseudo target \hat{x}^* is noisy, Algorithm 1 can still generate similar data that is meaningful lies on the data manifold. These two properties enable Algorithm 2 to generate data with lower and lower scores along the data manifold.

Set pseudo target \hat{x}^* by historical optimal update: For the historical optimal update case, we set pseudo target \hat{x}^* as the empirical optimal point from the historical dataset \mathcal{D} in Algorithm 2 as:

$$\hat{x}^* = \arg \min_{x \in \mathcal{D}} f(x). \quad (4)$$

Set pseudo target \hat{x}^* as Eq.(4) is an empirical approximation of the optimal target $x^* = \arg \min_{x \in \mathcal{X}} f(x)$. In Algorithm 2, at each iteration w.r.t. the batch query i , we employ the generated data that achieves the current best score as the pseudo target \hat{x}^* to guide the generation process. Intuitively, given this pseudo target \hat{x}^* , Algorithm 2 will generate a batch of data that is similar to \hat{x}^* on the low-dimensional data manifold, which is promising to further improve the score. By looping this procedure, we will get data generation with lower and lower scores.

4 EXPERIMENTS

In this section, we evaluate our algorithm in two domains, images and molecules, and compare it against four baseline methods: DDPO (Black et al., 2023), DPOK (Fan et al., 2024), D3PO (Yang et al., 2024), DNO (Tang et al., 2024). [Moreover, we further evaluate our algorithm in compressibility, incompressibility, and aesthetic quality tasks in Appendix A.](#) Our source code will be made public available upon publication.

4.1 IMAGE BLACK-BOX TARGET GENERATION TASK

Problem: Prompt Alignment. We consider the image-prompt alignment problem. While the current state-of-the-art image generative models excel at generating highly realistic images, they sometime struggle to faithfully generate images that accurately aligned with the input prompts, especially those complex prompts involving rare object combinations, object counting, or specific object positioning.

Pre-trained Model: SDXL-Lightning. We use SDXL (Podell et al., 2023) diffusion model as the backbone text-to-image model. It is able to generate 1024×1024 high resolution realistic image. In our experiment, we use the distilled version, SDXL-Lightning (Lin et al., 2024), for its high sampling efficiency, which can generate image with comparable quality with just $K = 8$ steps. We use the official implementation²

Objective Function: Gemini 1.5. We leverage Gemini 1.5 (Reid et al., 2024), an advanced multi-modal LLM service, as our black-box objective function to evaluate the alignment between input prompts and generated images. To avoid confusion, the term *query* refers to the input to Gemini, while *prompt* refers to the text used for image generation.

The query to Gemini 1.5 composes of the generated images with the question like: "Does the prompt $\$prompt$ accurately describe the image? Rate from 1 to 5". We state the complete query in Appendix G.

Because it is a closed-source paid service, we limit the number of batch queries in our experiments, referred to as the *batch queries budget*. We use the Gemini 1.5 Flash model (code: gemini-1.5-flash-001) for its cost-efficiency, and set the temperature as 0 for experiments consistency and reproducibility. For conciseness, we call it simply *Gemini* in the following section.

Experiment Procedure. We identify 12 prompts that the pre-trained model SDXL-Lightning struggles to generate, and refer these as the 12 tasks in our experiment, where the goal is to generate images that accurately aligned with the input prompts. We compare our Fast Direct algorithm against each baselines methods, and each experiment is constrained with 50 of batch queries budget.

We perform inference-time guidance using Fast Direct (Algorithm 2) [to maximize the Gemini rating](#). We use the GP model in Eq. 2 for the pseudo-target, and set the kernel as Gaussian kernel, and we follow (Hvarfner et al., 2024) to set the lengthscale as $\lambda = \sqrt{d}$, where $d = 4 \times 128 \times 128$ is the latents dimensionality of SDXL. We use $N = 50$ iterations to utilize 50 batch queries budget, and set the batch size as $B = 32$ and step size as $\alpha = 80$. [We use the EulerDiscreteScheduler \(Karras et al., 2022\) sampler as suggested by the SDXL-Lightning implementation \(Lin et al., 2024\), and the DDIMScheduler \(Song et al., 2020\) \(DDIM\) sampler for a fair comparison with the baselines.](#)

[For all baseline methods, we perform 50 iterations, each iteration utilizes one batch query for updating. We set the batch size as 32, and leave the rest of the hyper-parameters as default. We state more details for baselines in the Appendix D. Note that for DNO, in which each experiment trial can only produce one image, we run 16 independent experiment trials and report the average results³.](#)

Experiment Results. We present the results for the "deer-elephant"⁴ task in this section and defer the results for the other 11 tasks to Appendix I due to space constraints. We compare the images generated by our method with those generated by the baseline methods at each 10 batch queries intervals in Fig. 2. We can observe that our algorithm successfully satisfies the visual objectives [for both samplers](#) within the 50 batch queries budget, while the baseline methods fall short within

²<https://huggingface.co/ByteDance/SDXL-Lightning>

³Note that DNO spent $16 \times 50 = 800$ batch queries for each task.

⁴The complete prompt is: A yellow reindeer and a blue elephant.

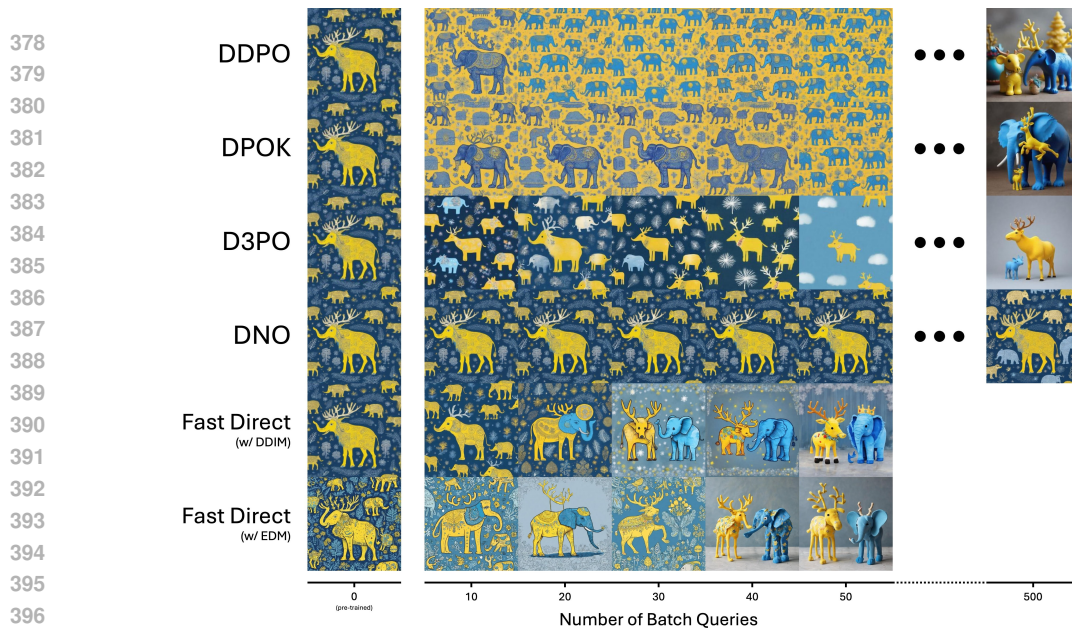


Figure 2: The generated images over each number of batch queries on the prompt "deer-elephant" ⁴, extra batch queries budgets (until 500) are given to the baselines methods for demonstration.

the same budget. Similar phenomenon are observed for the other 11 tasks in Appendix I. We grant baseline methods with extra batch queries budget, and we observed that DDPO and DPOK almost achieves the visual objectives, whereas D3PO and DNO remain unsuccessful in this particular task.

We collect the 32 randomly generated images using our algorithm (by utilizing 50 batch queries budget) and present in Fig. 3. We can observe that most images are well satisfies the visual objective. Similar phenomenon are observed on the other 11 tasks as shown in Appendix I.

For the quantitative results, we report the average objective values over each batch query across the 12 tasks in Fig. 4. For DNO, which only produces one image each experiment, we average the results from 16 independent experiments. We observe that the Fast Direct achieves nearly full score within the 50 batch queries budget in all tasks, which is significant higher than the baseline methods. Moreover, we observe that Fast Direct achieve similar score for both samplers, suggesting it is invariant to sampler. We further perform ablation study in Appendix B, and show that Fast Direct is insensitive to the hyper-parameters: step size α and batch size B .

We further employ tasks 1 to 3 and assigns extra batch queries budget of 200 to the baseline methods. We report the *accumulated objective values*, which represents the best objective values achieved up to each number of batch queries, in Appendix Fig. 11. We observe that our algorithm, utilizing only 50 batch queries budget, consistently outperforms all baseline methods even when they are assigned with expanded batch queries budget of 200.

To quantify this advantage, in Table 1, we identify the minimum number of batch queries budget (denotes as N^*) our algorithm requires to surpass each baseline methods with their 200 batch query budget. We further calculate ours *batch-query-efficiency gain* compared to baselines as $\frac{200}{N^*}$ and report (in the bracket) in Table 1. We can observe that our algorithm is at least 667% and up to 1000% more batch-queries-efficient than the baseline methods.

4.2 MOLECULE BLACK-BOX TARGET GENERATION TASK

Problem: Drug Discovery. We consider the drug discovery problem. One of the key problem is to find drug molecules that has a strong binding affinity with the target protein receptor. The binding affinity quantify the strength of the interaction between two bio-molecules, so a stronger binding affinity indicates a higher drug efficacy.

A traditional approach to this problem requires human-design molecules, which are highly time-consuming. Recent advancements in AI poses a great potential to speedup this process by generating molecules using advanced generative model such as diffusion model. Although the cur-

432
433
434
435
436
437
438
439
440
441
442
443
444
445



Figure 3: The 32 randomly generated image for the prompt "deer-elephant"⁴ guided by Fast Direct (w/ EDM) by utilizing 50 batch queries budget.

449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464

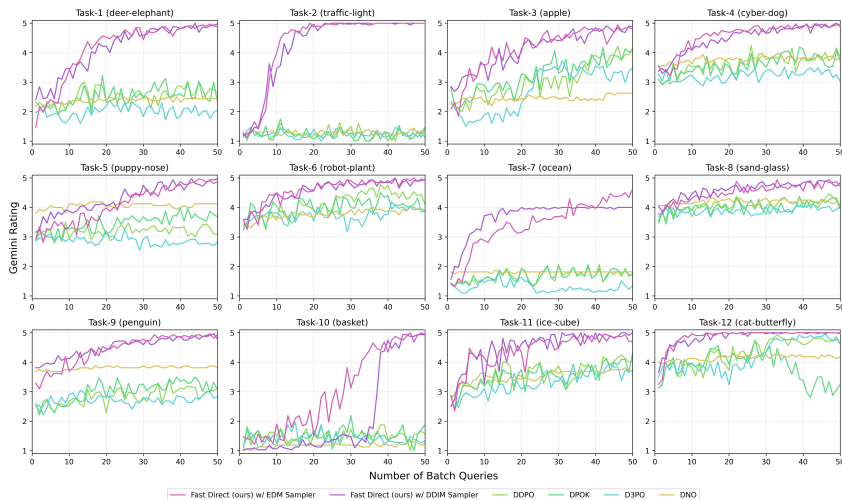


Figure 4: The average Gemini rating (from 1 to 5, higher is better) of the generated images over each number of batch queries on the 12 different tasks (the prompts abbreviation shown in bracket, see the complete prompts in Appendix 3).

rent pre-trained molecules models, such as TargetDiff, are capable of generating relatively realistic molecules, it lacks the ability of generating molecules with user desired high binding affinity.

Pre-trained Model: TargetDiff. We use TargetDiff (Guan et al., 2023) as the backbone molecules generative model. It is pre-trained on the CrossDocked2020 dataset (Francoeur et al., 2020), and can generate realistic 3D molecular structures conditioned on the given protein receptor. For sampling efficiency, we use the DDIMScheduler sampler with $K = 200$ sampling steps as we find that it can generate comparable results. During the generation process, we fix the atoms type according to the dataset reference and only allow the atoms position to be varied. We use the official implementation⁵ in our experiment.

Objective Function: Molecules Binding Affinity. Measurement of the binding affinity requires real-world experiment which is highly expensive (David et al., 2020). In our experiment, we leverage the Vina score (kcal/mol) calculated by the AutoDock Vina simulation software (Eberhardt et al., 2021) to estimate the binding affinity. A lower (more negative) Vina score indicates a stronger estimated binding affinity. Therefore, our objective is to minimize the Vina score.

Experiment Procedure. We conduct experiments on 6 tasks, where the goal of each task is to optimize the Vina score on the protein receptors of ID from 1 to 6, respectively. Similar to images tasks, each algorithm is constrained with 50 of batch queries budget.

⁵<https://github.com/guanjq/targetdiff>

Table 1: The minimum number of batch queries budget (denotes as N^*) for Fast Direct (ours) to outperform each baselines (when each baseline is granted with 200 batch queries budget) on images and molecules tasks. The values in bracket show the batch-query-efficiency gain ($\frac{200}{N^*}$) of Fast Direct (ours) over the baselines.

	Minimum number of batch queries (N^*) for Fast Direct to outperform each baselines:		
	DDPO	DPOK	D3PO
Image Tasks			
Task 1	15 (13.34 \times)	15 (13.34 \times)	10 (20.00 \times)
Task 2	8 (25.00 \times)	8 (25.00 \times)	4 (50.00 \times)
Task 3	37 (5.41 \times)	46 (4.34 \times)	16 (12.50 \times)
Average	20 (10.00 \times)	23 (8.70 \times)	30 (6.67 \times)
Molecules Tasks			
Task 1	6 (33.33 \times)	5 (40.00 \times)	25 (8.00 \times)
Task 2	9 (22.22 \times)	4 (50.00 \times)	11 (18.18 \times)
Average	7.5 (26.67 \times)	4.5 (44.44 \times)	18 (11.11 \times)

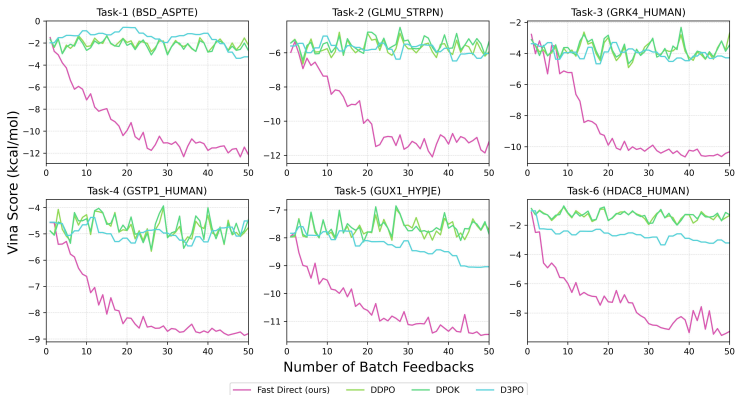


Figure 5: The Vina score (lower is better) of the generated molecules for each number batch queries on the six protein receptors.

We perform inference-time guidance using our Algorithm 2 to minimize the Vina score. We use the historical optimal in Eq. 4 for the pseudo-target. We use $N = 50$ iterations to utilize 50 batch queries budget, and set the batch size as $B = 32$, and the step size as $\alpha = 10^{-2}$. We then conduct experiments on the baseline methods for 50 fine-tuning iteration. For DDPO and DPOK, the batch size is set as the same 32; while D3PO requires binary rewards, so we double the batch size as 64⁶. We keep the hyper-parameters as default for all baseline methods.

Experiment Results. We report the average Vina score (recall, lower is better) for each batch queries across 6 tasks in Fig. 5. We observe that our algorithm consistently achieves much lower Vina score compared to the baseline methods across all tasks. Similar to the images task, we grant extra batch queries budget of 200 to the baselines for task 1 and 2, and report the *accumulated objective values* in Appendix Fig 12, and list the *batch-query-efficiency gain* in Table 1. We observe that our algorithm is at least 1111% and up to 4444% more batch-query-efficient compared to the baseline methods.

5 CONCLUSION

In this work, we proposed **Fast Direct** for diffusion model target generation, which effectively addresses the challenges of limited batch query budgets and black-box objective functions, demonstrating its potential for various applications, including image generation and drug discovery. Fast Direct is highly practical, as it is easy to implement, supports any type of SDE solver, and has only one hyper-parameter to tune (step size α). Our algorithm is based on the surprising empirical observation that the *universal update direction* (i.e., $\mathbf{x}^* - \mathbf{x}_K$ in Algorithm 1) can efficiently guide the diffusion trajectory toward the target, even when the target \mathbf{x}^* is extremely noisy. This phenomenon suggests an intriguing robustness in the diffusion inference process. Future research could investigate for its underlying theoretical principles.

⁶However, unlike image tasks where Gemini can evaluate rewards for two images in a single query, the Vina simulation lacks this capability, requiring 64 individual queries for a batch of size 64.

REFERENCES

- 540
541
542 Arpit Bansal, Hong-Min Chu, Avi Schwarzschild, Soumyadip Sengupta, Micah Goldblum, Jonas
543 Geiping, and Tom Goldstein. Universal guidance for diffusion models. In *Proceedings of the*
544 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 843–852, 2023.
- 545
546 Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion
547 models with reinforcement learning. *arXiv preprint arXiv:2305.13301*, 2023.
- 548
549 Minshuo Chen, Song Mei, Jianqing Fan, and Mengdi Wang. An overview of diffusion models: Ap-
550 plications, guided generation, statistical rates and optimization. *arXiv preprint arXiv:2404.07771*,
551 2024.
- 552
553 Hyungjin Chung, Jeongsol Kim, Michael T Mccann, Marc L Klasky, and Jong Chul Ye. Diffusion
554 posterior sampling for general noisy inverse problems. *arXiv preprint arXiv:2209.14687*, 2022.
- 555
556 Kevin Clark, Paul Vicol, Kevin Swersky, and David J Fleet. Directly fine-tuning diffusion models
557 on differentiable rewards. *arXiv preprint arXiv:2309.17400*, 2023.
- 558
559 Gabriele Corso, Hannes Stärk, Bowen Jing, Regina Barzilay, and Tommi Jaakkola. Diffdock: Dif-
560 fusion steps, twists, and turns for molecular docking. *arXiv preprint arXiv:2210.01776*, 2022.
- 561
562 Florinel-Alin Croitoru, Vlad Hondru, Radu Tudor Ionescu, and Mubarak Shah. Diffusion models
563 in vision: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(9):
564 10850–10869, 2023. doi: 10.1109/TPAMI.2023.3261988.
- 565
566 Laurianne David, Amol Thakkar, Rocío Mercado, and Ola Engkvist. Molecular representations in
567 ai-driven drug discovery: a review and practical guide. *Journal of Cheminformatics*, 12(1):56,
568 2020.
- 569
570 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances*
571 *in neural information processing systems*, 34:8780–8794, 2021.
- 572
573 Jerome Eberhardt, Diogo Santos-Martins, Andreas F Tillack, and Stefano Forli. Autodock vina
574 1.2. 0: New docking methods, expanded force field, and python bindings. *Journal of chemical*
575 *information and modeling*, 61(8):3891–3898, 2021.
- 576
577 Luca Eyring, Shyamgopal Karthik, Karsten Roth, Alexey Dosovitskiy, and Zeynep Akata. Reno:
578 Enhancing one-step text-to-image models through reward-based noise optimization. *arXiv*
579 *preprint arXiv:2406.04312*, 2024.
- 580
581 Ying Fan, Olivia Watkins, Yuqing Du, Hao Liu, Moonkyung Ryu, Craig Boutilier, Pieter Abbeel,
582 Mohammad Ghavamzadeh, Kangwook Lee, and Kimin Lee. Reinforcement learning for fine-
583 tuning text-to-image diffusion models. *Advances in Neural Information Processing Systems*, 36,
584 2024.
- 585
586 Paul G Francoeur, Tomohide Masuda, Jocelyn Sunseri, Andrew Jia, Richard B Iovanisci, Ian Snyder,
587 and David R Koes. Three-dimensional convolutional neural networks and a cross-docked data set
588 for structure-based drug design. *Journal of chemical information and modeling*, 60(9):4200–
589 4215, 2020.
- 590
591 Giorgio Giannone, Akash Srivastava, Ole Winther, and Faez Ahmed. Aligning optimization trajec-
592 tories with diffusion models for constrained design generation. *Advances in Neural Information*
593 *Processing Systems*, 36:51830–51861, 2023.
- 594
595 Jiaqi Guan, Wesley Wei Qian, Xingang Peng, Yufeng Su, Jian Peng, and Jianzhu Ma. 3d equiv-
596 ariant diffusion for target-aware molecule generation and affinity prediction. *arXiv preprint*
597 *arXiv:2303.03543*, 2023.
- 598
599 Yutong He, Naoki Murata, Chieh-Hsin Lai, Yuhta Takida, Toshimitsu Uesaka, Dongjun Kim, Wei-
600 Hsiang Liao, Yuki Mitsufuji, J Zico Kolter, Ruslan Salakhutdinov, et al. Manifold preserving
601 guided diffusion. *arXiv preprint arXiv:2311.16424*, 2023.

- 594 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
595 *neural information processing systems*, 33:6840–6851, 2020.
- 596
- 597 Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J
598 Fleet. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633–
599 8646, 2022.
- 600 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,
601 and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint*
602 *arXiv:2106.09685*, 2021.
- 603
- 604 Carl Hvarfner, Erik Orm Hellsten, and Luigi Nardi. Vanilla bayesian optimization performs great in
605 high dimension. *arXiv preprint arXiv:2402.02229*, 2024.
- 606 Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-
607 based generative models. *Advances in neural information processing systems*, 35:26565–26577,
608 2022.
- 609
- 610 Korrawe Karunratanakul, Konpat Preechakul, Emre Aksan, Thabo Beeler, Supasorn Suwajanakorn,
611 and Siyu Tang. Optimizing diffusion noise can serve as universal motion priors. In *Proceedings*
612 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1334–1345, 2024.
- 613 Siddarth Krishnamoorthy, Satvik Mehul Mashkaria, and Aditya Grover. Diffusion models for black-
614 box optimization. *arXiv preprint arXiv:2306.07180*, 2023.
- 615
- 616 Xiner Li, Yulai Zhao, Chenyu Wang, Gabriele Scalia, Gokcen Eraslan, Surag Nair, Tommaso Bian-
617 calani, Aviv Regev, Sergey Levine, and Masatoshi Uehara. Derivative-free guidance in continuous
618 and discrete diffusion models with soft value-based decoding. *arXiv preprint arXiv:2408.08252*,
619 2024.
- 620 Shanchuan Lin, Anran Wang, and Xiao Yang. Sdxl-lightning: Progressive adversarial diffusion
621 distillation. *arXiv preprint arXiv:2402.13929*, 2024.
- 622
- 623 Cheng Lu, Huayu Chen, Jianfei Chen, Hang Su, Chongxuan Li, and Jun Zhu. Contrastive energy
624 prediction for exact energy-guided diffusion sampling in offline reinforcement learning. In *Inter-*
625 *national Conference on Machine Learning*, pp. 22825–22855. PMLR, 2023.
- 626
- 627 Yurii Nesterov and Vladimir Spokoiny. Random gradient-free minimization of convex functions.
628 *Foundations of Computational Mathematics*, 17(2):527–566, 2017.
- 629
- 630 Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models.
631 In *International conference on machine learning*, pp. 8162–8171. PMLR, 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
635 synthesis. *arXiv preprint arXiv:2307.01952*, 2023.
- 636
- 637 Mihir Prabhudesai, Anirudh Goyal, Deepak Pathak, and Katerina Fragkiadaki. Aligning text-to-
638 image diffusion models with reward backpropagation. *arXiv preprint arXiv:2310.03739*, 2023.
- 639
- 640 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea
641 Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances*
642 *in Neural Information Processing Systems*, 36, 2024.
- 643
- 644 Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-
645 baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gem-
646 ini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint*
647 *arXiv:2403.05530*, 2024.
- 648
- 649 Christoph Schuhmann. Laion aesthetics. <https://laion.ai/blog/laion-aesthetics/>, August 2022.
- 650
- 651 John Schulman, Xi Chen, and Pieter Abbeel. Equivalence between policy gradients and soft q-
652 learning. *arXiv preprint arXiv:1704.06440*, 2017.

- 648 Matthias Seeger. Gaussian processes for machine learning. *International journal of neural systems*,
649 14(02):69–106, 2004.
- 650
- 651 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv*
652 *preprint arXiv:2010.02502*, 2020.
- 653 Jiaming Song, Qinsheng Zhang, Hongxu Yin, Morteza Mardani, Ming-Yu Liu, Jan Kautz, Yongxin
654 Chen, and Arash Vahdat. Loss-guided diffusion models for plug-and-play controllable generation.
655 In *International Conference on Machine Learning*, pp. 32483–32498. PMLR, 2023.
- 656
- 657 Zhiwei Tang, Jiangweizhi Peng, Jiasheng Tang, Mingyi Hong, Fan Wang, and Tsung-Hui Chang.
658 Tuning-free alignment of diffusion models with direct noise optimization. *arXiv preprint*
659 *arXiv:2405.18881*, 2024.
- 660 Masatoshi Uehara, Yulai Zhao, Tommaso Biancalani, and Sergey Levine. Understanding rein-
661 forcement learning-based fine-tuning of diffusion models: A tutorial and review. *arXiv preprint*
662 *arXiv:2407.13734*, 2024.
- 663
- 664 Nikolaos N Vlassis and WaiChing Sun. Denoising diffusion algorithm for inverse design of mi-
665 crostructures with fine-tuned nonlinear material properties. *Computer Methods in Applied Me-*
666 *chanics and Engineering*, 413:116126, 2023.
- 667 Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou, Aaron Lou, Senthil Purushwalkam,
668 Stefano Ermon, Caiming Xiong, Shafiq Joty, and Nikhil Naik. Diffusion model alignment using
669 direct preference optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
670 *and Pattern Recognition*, pp. 8228–8238, 2024.
- 671 Kai Yang, Jian Tao, Jiafei Lyu, Chunjiang Ge, Jiabin Chen, Weihang Shen, Xiaolong Zhu, and Xiu Li.
672 Using human feedback to fine-tune diffusion models without any reward model. In *Proceedings*
673 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8941–8951, 2024.
- 674
- 675 Hui Yuan, Kaixuan Huang, Chengzhuo Ni, Minshuo Chen, and Mengdi Wang. Reward-directed
676 conditional diffusion: Provable distribution estimation and reward improvement. *Advances in*
677 *Neural Information Processing Systems*, 36, 2024.
- 678
- 679
- 680
- 681
- 682
- 683
- 684
- 685
- 686
- 687
- 688
- 689
- 690
- 691
- 692
- 693
- 694
- 695
- 696
- 697
- 698
- 699
- 700
- 701

A MORE EXPERIMENTS: COMPRESSIBILITY, INCOMPRESSIBILITY, AND AESTHETIC QUALITY

Problem. We follow DDPO (Black et al., 2023) to evaluate our algorithm on the three black-box optimization tasks for images: compressibility, incompressibility, and aesthetic quality.

Objective Function. For compressibility, we aimed to minimize the compressed JPEG size (MB) of the generated images. For incompressibility, it’s simply the inverse. For aesthetic quality, the aesthetic score is evaluated by the pre-trained LAION aesthetics predictor (Schuhmann, 2022), which is trained on human rating of the images aesthetic quality, and we aimed to maximize the aesthetic score.

Experiment Procedure. For each task, we uniformly sample from the 45 common animals (as proposed by DDPO (Black et al., 2023)) as input prompts and optimize the objective using Fast Direct in Algorithm 2. The prompt for each instance within a batch is sampled randomly and independently. We set $N = 50$ batch queries budget for compressibility and incompressibility, and $N = 100$ for aesthetic quality. We use the same pre-trained model, hyper-parameters, and GP settings as in Section 4.1. The 45 common animals prompts used in the experiment is as follows:

cat	dog	horse	monkey	rabbit	zebra	spider	bird	sheep
deer	cow	goat	lion	tiger	bear	raccoon	fox	wolf
lizard	beetle	ant	butterfly	fish	shark	whale	dolphin	squirrel
mouse	rat	snake	turtle	frog	chicken	duck	goose	bee
pig	turkey	fly	llama	camel	bat	gorilla	hedgehog	kangaroo

For DNO, where each experimental trial generates only a single image and the objective score is highly dependent on the specific prompt. To ensure a more accurate evaluation, we run 45 independent experiment trials using each prompt and report the average results. Consequently, DNO requires $45\times$ more batch queries per task compared to other methods.

Experiment Result. We present the objective scores for each task in the left column of Fig. 6 and provide the generated images for the three tasks in the supplementary materials. For compressibility and incompressibility, we observe that Fast Direct achieves significantly better scores than the baselines. For aesthetic quality, Fast Direct with the EDM sampler achieves significantly better scores than with DDIM, likely due to EDM being a more advanced sampler, capable of generating higher-quality images. DNO achieves comparable scores; however, note that each experiment optimizes and generates only a single image.

Generalization Ability. We follow DDPO to evaluate generalization ability. For Fast Direct, we freeze the learned GP model to generate 16 unseen images using distinct unseen animal prompts. Specifically, in this phase, Lines 17 and 18 are removed from Algorithm 2, and our algorithm does not access the objective function. For DDPO, DPOK, and D3PO, the fine-tuned models are frozen to generate images with the 16 unseen animal prompts. DNO is not applicable to this experiment because DNO needs to perform optimization for each image without generalization. For the unseen prompts, since DDPO didn’t publish the unseen prompts, we created our own as follows:

elephant	eagle	pigeon	hippo	hamster	otter	panda	reindeer
owl	penguin	flamingo	seal	koala	giraffe	parrot	cheetah

In Fig.7, we present the images generated using the unseen prompt ”hippo” by each algorithm across the three tasks. The remaining 15 prompts for the three tasks are provided in the supplementary material. A similar phenomenon is observed for all 15 other unseen prompts across the three tasks. For quantitative evaluation, we report the objective values for the unseen prompts in the right column of Fig.6. We can observe that Fast Direct with GP model has generalization ability to unseen prompts.

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

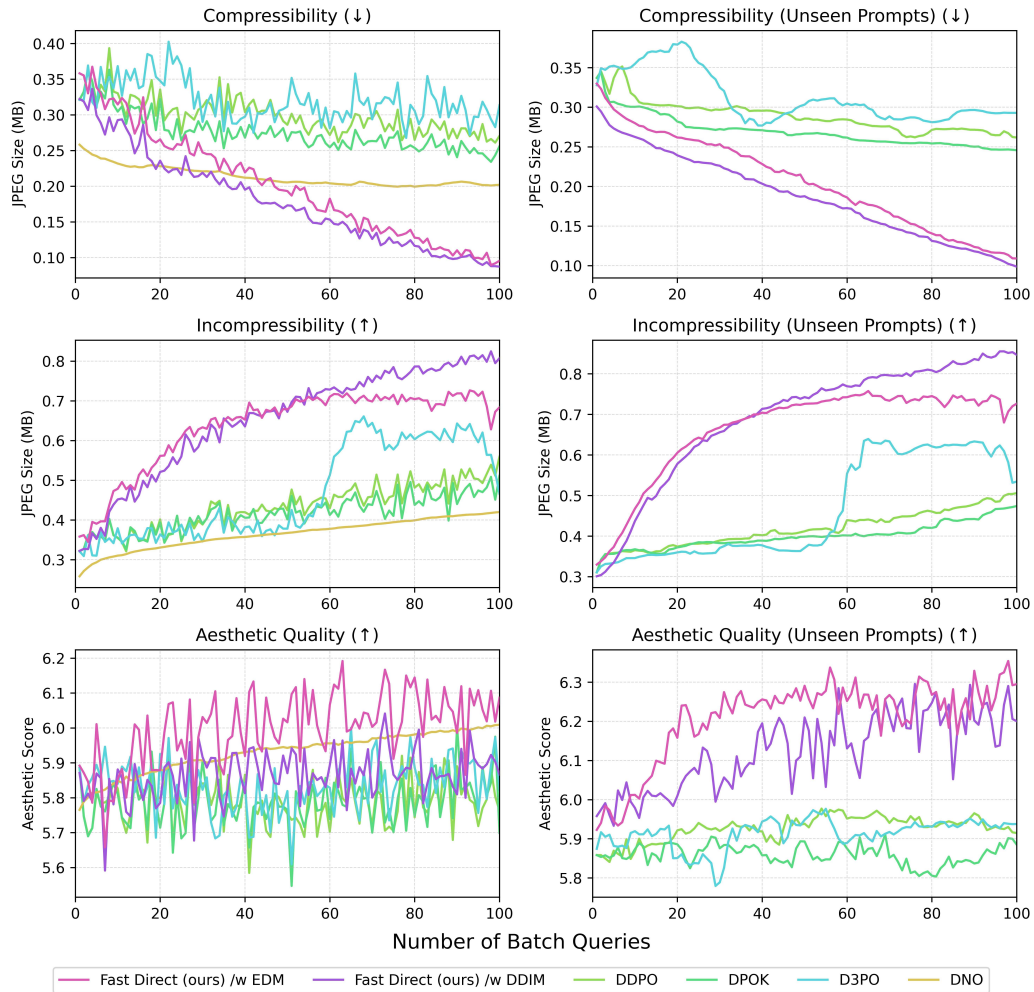


Figure 6: Left column: The average objective score of the generated images over each number of batch queries on the 3 black-box optimization tasks, the images are generated using the 45 common animals that used in DDPO (Tang et al., 2024) as the input prompts. Right column: The objective score of the images generated by **unseen prompts**, it demonstrates the generalization capability. Note that DNO is not applicable to this task.

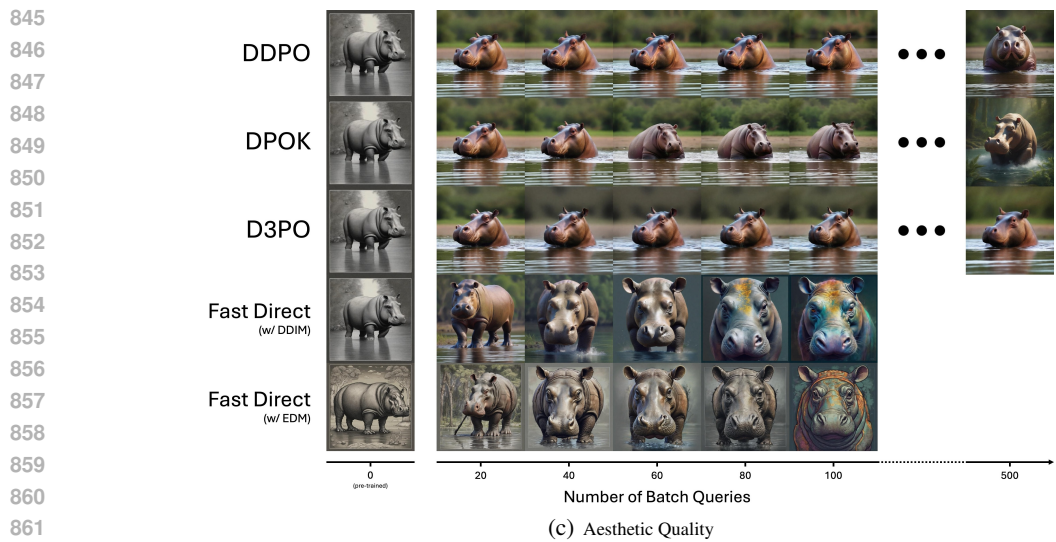
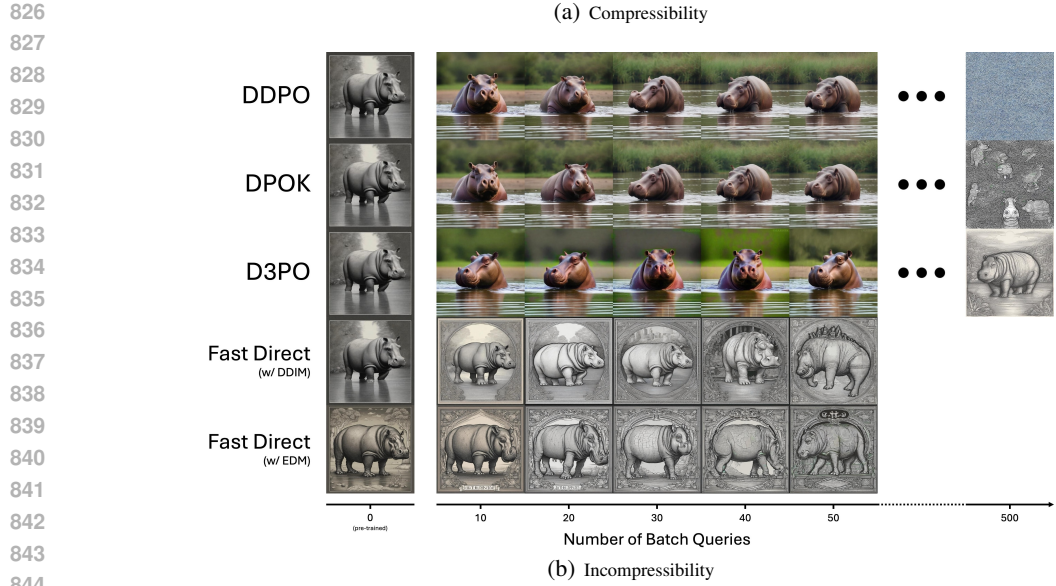
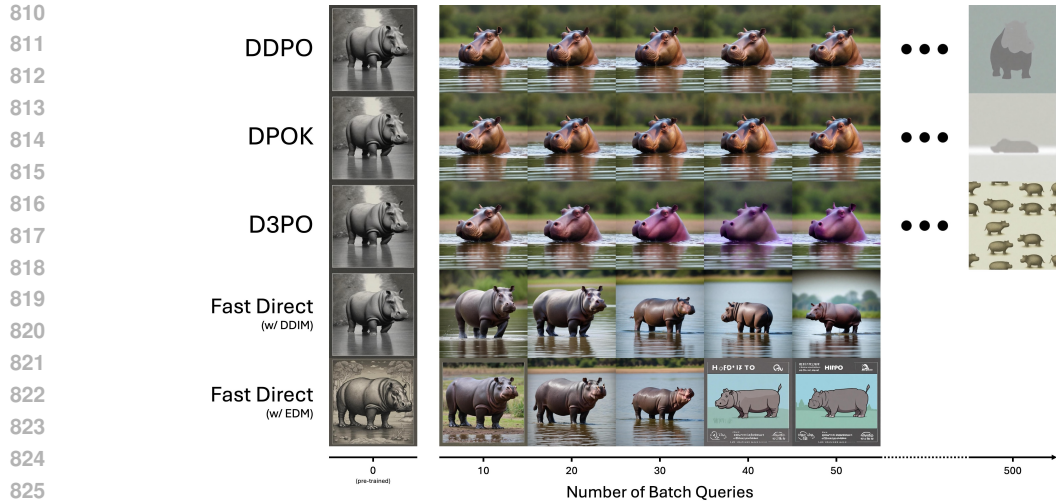


Figure 7: The generated images using the unseen prompt "hippo" in aesthetic quality task, extra batch queries budgets (until 500) are given to the baselines methods for demonstration.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

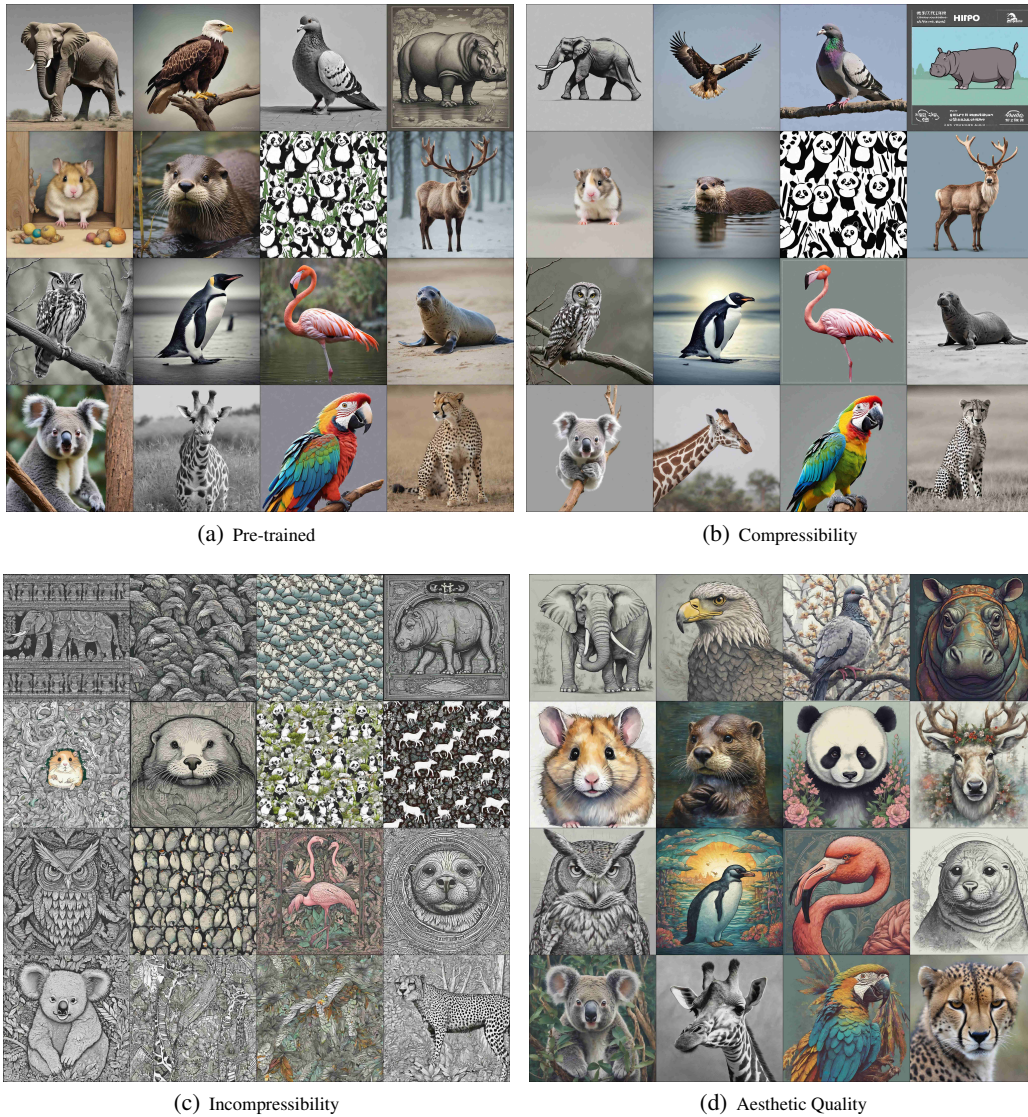


Figure 8: The generated images using the 16 unseen prompts. Top left is the images generated by pre-trained model, the rest are the Fast Direct generated images in the corresponding tasks.

B ABLATION STUDY

We perform an ablation study on the target image generation Task-1. For step size analysis, we fix the batch size as 32, then perform experiment with different step size $\alpha = \{20, 40, 80, 160, 320\}$; for batch size analysis, we fix the step size as 80, then perform experiment with different batch size $B = \{4, 8, 16, 32, 64\}$. Additionally, we report the run time for different batch size. We report the result in Fig. 9.

We can observe that the performance steadily increase for any step size and any batch size, suggests that our algorithm is not sensitive to the hyper-parameters settings. The run time scales linearly with the batch size. In our target image generation experiments in Section 4.1, as the batch size is set as 32, each experiment takes approximately 6.4 hours to process 32 images in parallel.

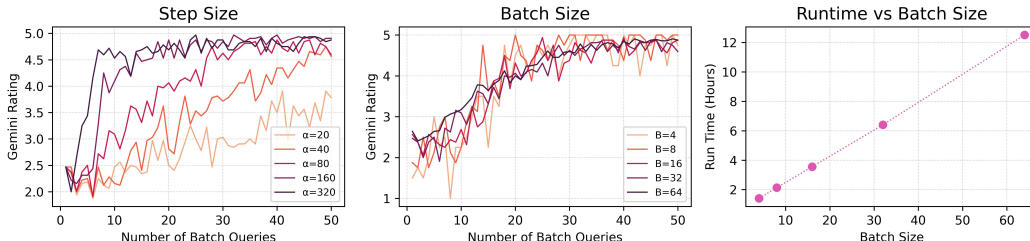


Figure 9: Left: Gemini rating for different step size settings. Middle: Gemini rating for different batch size settings. Right: Run time (hours) for different batch size settings.

C ANALYSIS OF UNIVERSAL DIRECTION

In Fig. 10, we demonstrate Algorithm 1 with update direction $\vec{d} = \mathbf{x}^* - \mathbf{x}_{K'}$ for $K' \in \{K, K/2, K/4, K/8\}$. It shows that the generated image quality decreases as the K' decreases. This is because as K' decreases, the $\mathbf{x}_{K'}$ becomes noisier and may move further away from the data manifold.

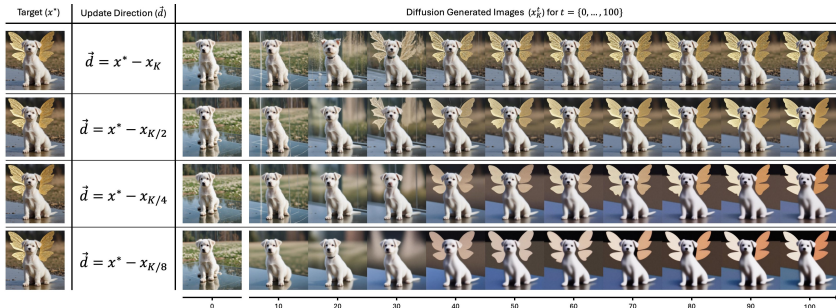


Figure 10: The update direction $\vec{d} = \mathbf{x}^* - \mathbf{x}_{K'}$ for $K' \in \{K, K/2, K/4, K/8\}$, the generated images quality decrease as the $\mathbf{x}_{K'}$ being more more noisy.

D BASELINES DETAILS

For DDPO, DPOK and D3PO, we fine-tune the model to maximize the Gemini rating. We fine-tune the model with 50 epochs; each epoch utilizes one batch query for model updating. For DDPO and DPOK, we set the batch size to 32; for D3PO, it requires a relative reward, so we doubled the batch size to 64.⁷ As these RL-based methods require closed-form expression of the logarithm probabilities, we follow their official implementation to use the DDIMScheduler (Song et al., 2020) sampler.

⁷For D3PO, to evaluate 64 images only requires 32 queries calls, see Appendix G.

For DNO, we optimize the noise sequence w.r.t the Gemini rating. Recall that the Gemini is black-box, so we use the non-differentiable mode, with the number of samples for gradient approximation set as 32, and optimize for 50 iterations; each iteration utilizes one batch query for updating. However, each experiment of DNO only produces one image, where the score is highly dependent on the initial prior. Thus, for a more accurate evaluation, we run 16 independent experiments to generate 16 images and report the average results. Note that this requires $16 \times 50 = 800$ batch queries for each task, which is 16 times compared with our Fast Direct and other baselines.

E PROOF OF PROPOSITION 1

Proof. Let $\boldsymbol{\alpha} = (\mathcal{K}(\mathbf{X}^n, \mathbf{X}^n) + \lambda \mathbf{I})^{-1} \mathbf{y} = [\alpha_1, \dots, \alpha_n]^\top$, for GP with a shift-invariant kernel that can be rewritten as $k(\mathbf{z}_1, \mathbf{z}_2) = g(\|\mathbf{z}_1 - \mathbf{z}_2\|_2)$, the gradient of the GP prediction is

$$\nabla \hat{f}(\mathbf{x}; \mathbf{X}^n) = \nabla \mathbf{k}(\mathbf{x}, \mathbf{X}^n)^\top \boldsymbol{\alpha} \quad (5)$$

$$= \sum_{i=1}^n \alpha_i \nabla k(\mathbf{x}, \mathbf{x}^i) \quad (6)$$

$$= \sum_{i=1}^n \frac{\alpha_i}{\|\mathbf{x} - \mathbf{x}^i\|} g'(\|\mathbf{x} - \mathbf{x}^i\|_2) (\mathbf{x} - \mathbf{x}^i) \quad (7)$$

$$= \sum_{i=1}^n c_i(\mathbf{x}) (\mathbf{x} - \mathbf{x}^i) \quad (8)$$

where \mathbf{x}^i denotes the i^{th} sample in $\mathbf{X}^n = [\mathbf{x}^1, \dots, \mathbf{x}^n]$, and $c_i(\mathbf{x}) = \frac{\alpha_i}{\|\mathbf{x} - \mathbf{x}^i\|} g'(\|\mathbf{x} - \mathbf{x}^i\|_2)$, and $g'(\cdot)$ denotes the derivative of $g(\cdot)$.

From Eq.(8), we can see that the gradient $\nabla \hat{f}(\mathbf{x}; \mathbf{X}^n)$ is a weighted sum of $\mathbf{x} - \mathbf{x}^i$ for $i \in \{1, \dots, n\}$. Thus, we know $\nabla \hat{f}(\mathbf{x}; \mathbf{X}^n)$ lies on a subspace spanned by $[\mathbf{x}, \mathbf{x}^1, \dots, \mathbf{x}^n]$

□

F OVERVIEW OF RELATED WORKS

We present the summary of related works according to its category and supported problem scenarios in Table 2.

Table 2: Overview of existing approaches.

Algorithm Category	Algorithm	Online	Black-box
Training / Fine-tune	DRaFT (Clark et al., 2023)	✓	✗
	DDOM (Krishnamoorthy et al., 2023), RCGDM (Yuan et al., 2024), DPO (Wallace et al., 2024)	✗	✓
	DDPO (Black et al., 2023), DPOK (Fan et al., 2024), D3PO (Yang et al., 2024)	✓	✓
Inference-time Guidance	LGD (Song et al., 2023), MPGD (He et al., 2023), DNO (2024a) (Karunratanakul et al., 2024), ReNO (Eyring et al., 2024)	✓	✗
	CEP (Lu et al., 2023)	✗	✓
	DNO (Tang et al., 2024), Fast Direct (ours)	✓	✓

G EXPERIMENT DETAILS

We use the following query question to the Gemini:

Does the prompt $\$prompt$ accurately describe the image?
 Rate from 1 (inaccurate) to 5 (accurate).
 Answer in the format: Score=(score), Reason=(reason).

where the $\$prompt$ is substituted by the input prompt used for image generation. We extract the integer score as the objective value.

The D3PO takes relative reward, so we specially crafted for its selective query question:

```
Given this two images, which image is better aligned with
the prompt $prompt?
Answer in the format: Choice=(1/2), Reason=(reason).
```

where we extract the choice (1/2) as the selections.

The Gemini 1.5 Flash can faithfully respond with the correct format to ensure our experiment is consistent.

H MORE EXPERIMENT RESULTS

We report the accumulated objective values over number of batch queries for images task in Fig. 11, and for molecules task in Fig. 12.

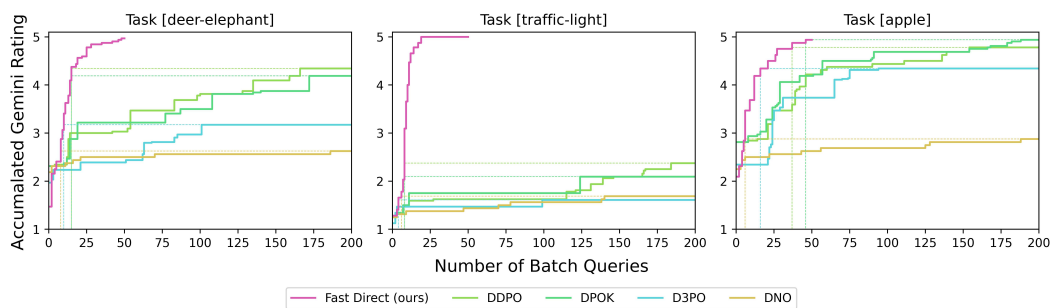


Figure 11: The accumulated Gemini rating (from 1 to 5, higher is better) over number of batch queries. The accumulated plot displays the maximum Gemini rating achieved up to each number of batch queries.

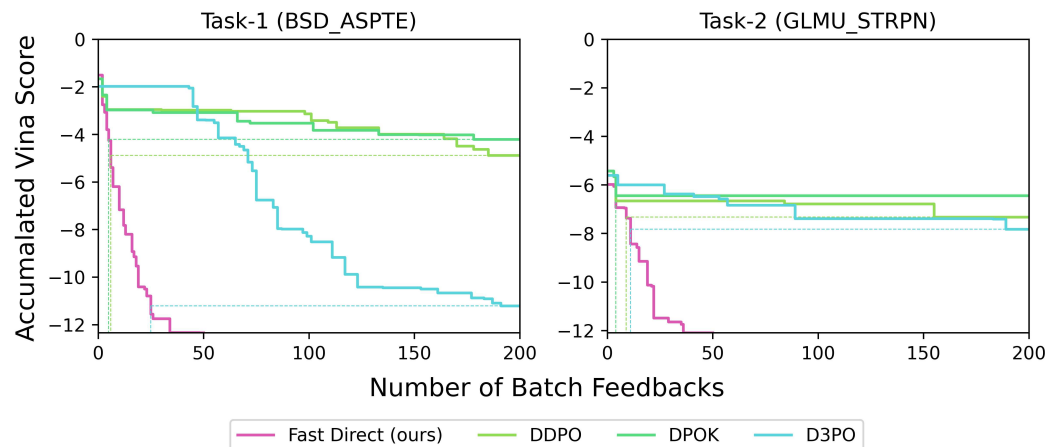


Figure 12: The accumulated Vina score (lower is better) over number of batch queries. The accumulated plot displays the minimum Vina score achieved up to each number of batch queries.

I MORE GENERATED IMAGES BY FAST DIRECT FOR IMAGE-PROMPT ALIGNMENT TASK

Table 3 shows the list the complete prompts used for each task in our image generation experiment. From Fig. 13 to Fig. 23, we present the generated images of each algorithms over each number of batch queries for prompt 2 to 12. From Fig. 24 to Fig. 34, we showcase the 32 randomly generated by pre-trained model followed by the resultant images guided by the Fast Direct for prompt 2 to 12.

Table 3: Target Image Generation Tasks

Task ID	Abbreviated Prompt	Complete Prompt
1	deer-elephant	A yellow reindeer and a blue elephant.
2	traffic-light	A traffic light with yellow at top, green at middle, and red at bottom.
3	apple	Seven red apples arranged in a circle.
4	cyber-dog	A natural fluffy dog talking to a cybertic dog.
5	puppy-nose	Side view of a puppy lying on floor, one butterfly stopping on its nose.
6	robot-plant	A cute cybernetic robot plants a tree in the forest.
7	ocean	A helicopter floating under the ocean.
8	sand-glass	A transparent glass filled with a mixture of water and sand, with one feather floating inside.
9	penguin	A photo realistic photo showing a penguin standing on a very small floating ice, with a tree is on fire in the background.
10	basket	Exactly one orange in a basket of apples.
11	ice-cube	A glass of water with exactly one ice cube.
12	cat-butterfly	A cat with butterfly wings.

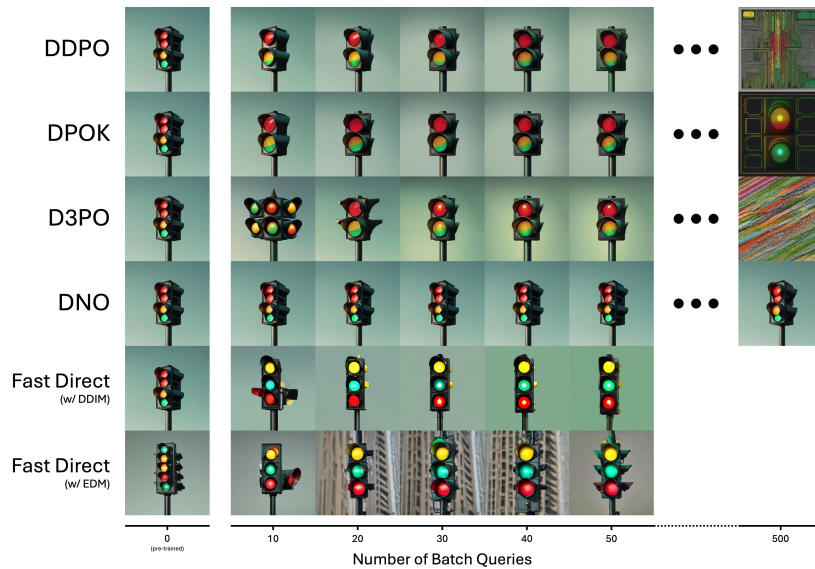
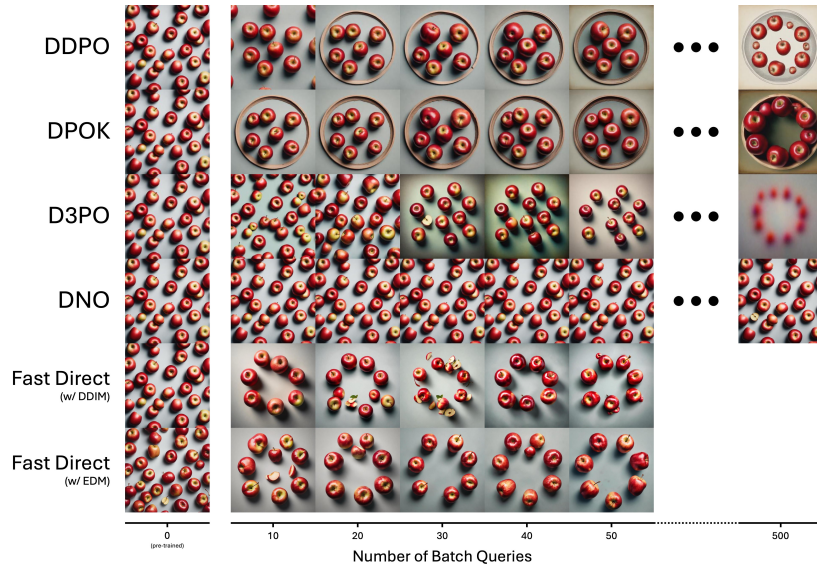


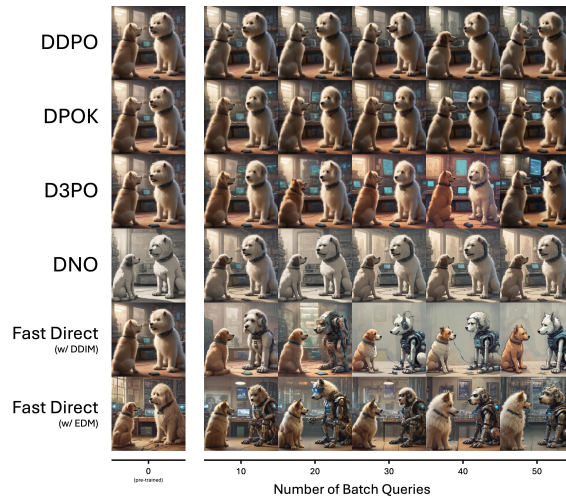
Figure 13: The generated images over each number of batch queries on the prompt "traffic-light", extra batch queries budgets (until 500) are given to the baselines methods for demonstration.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155



1156
1157 Figure 14: The generated images over each number of batch queries on the prompt "apple", extra
1158 batch queries budgets (until 500) are given to the baselines methods for demonstration.

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 Figure 15: The generated images over each number of batch queries on the prompt "cyber-dog" for
1184 each algorithms.

1185
1186
1187

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

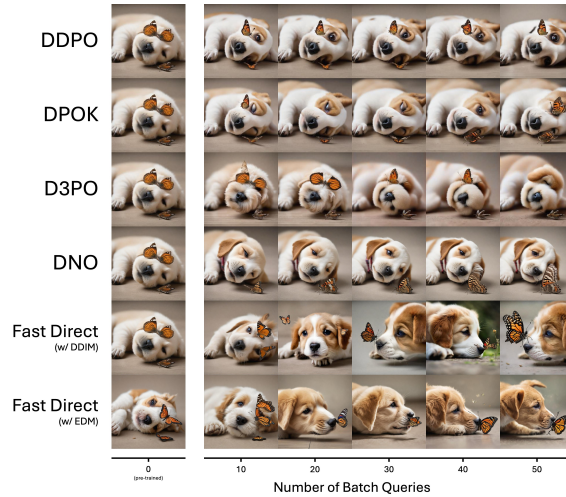


Figure 16: The generated images over each number of batch queries on the prompt "puppy-nose" for each algorithms.

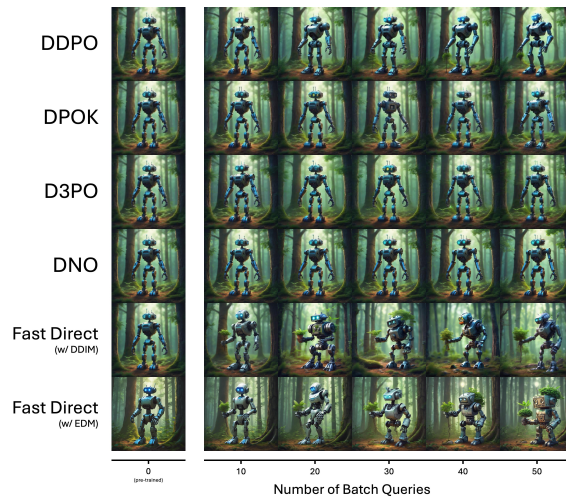


Figure 17: The generated images over each number of batch queries on the prompt "robot-plant" for each algorithms.

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

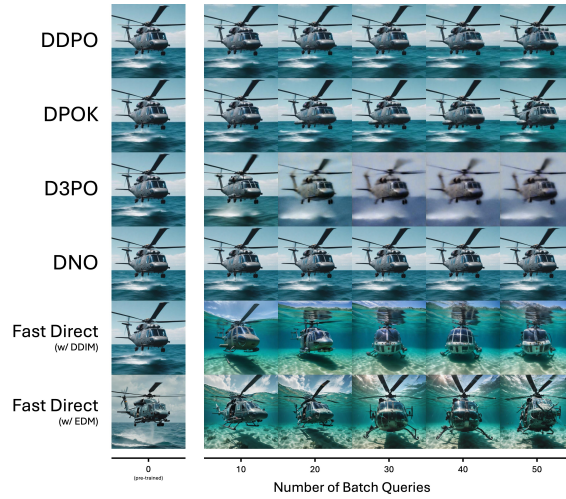


Figure 18: The generated images over each number of batch queries on the prompt "ocean" for each algorithms.

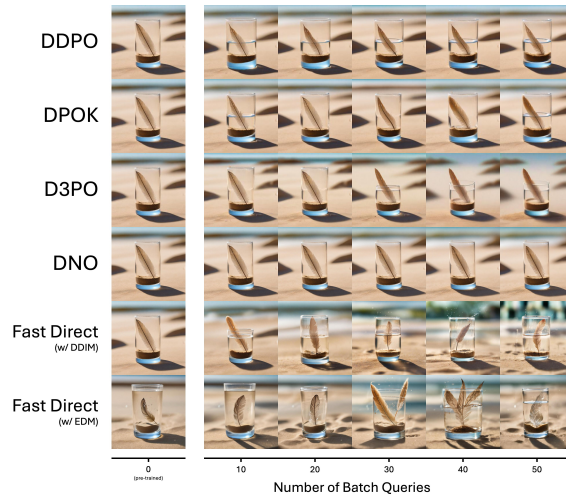


Figure 19: The generated images over each number of batch queries on the prompt "sand-glass" for each algorithms.

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

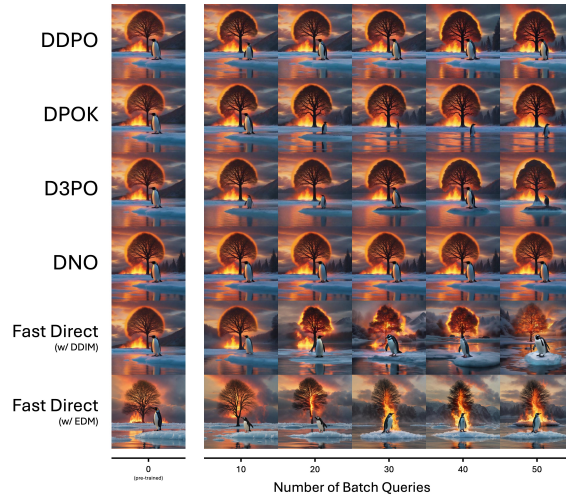


Figure 20: The generated images over each number of batch queries on the prompt "penguin" for each algorithms.

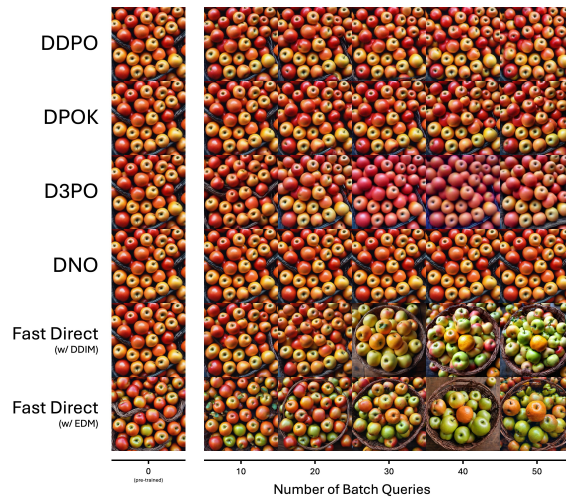


Figure 21: The generated images over each number of batch queries on the prompt "basket" for each algorithms.

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

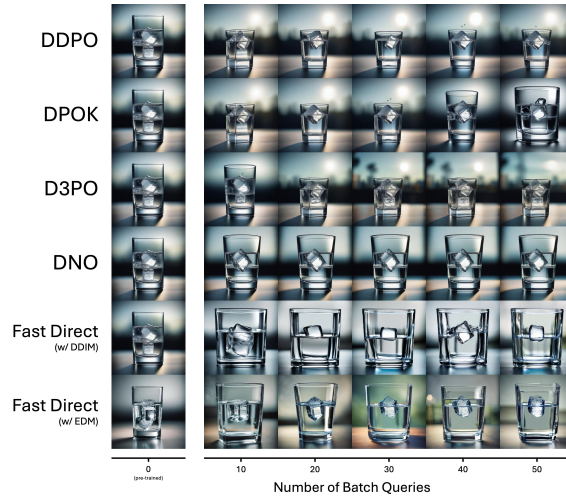


Figure 22: The generated images over each number of batch queries on the prompt "ice-cube" for each algorithms.

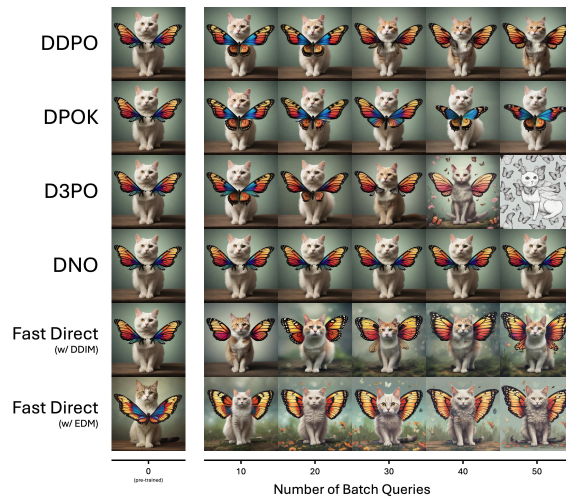
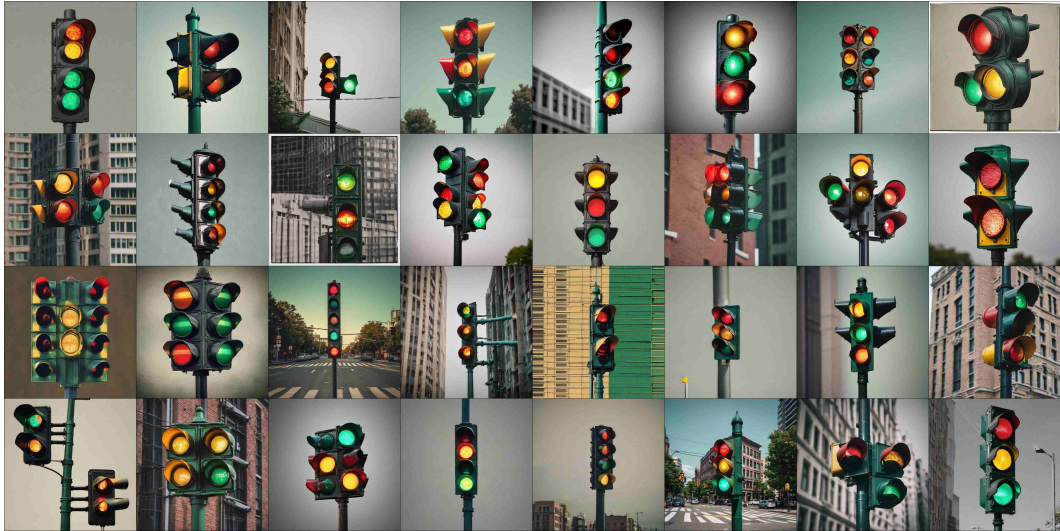
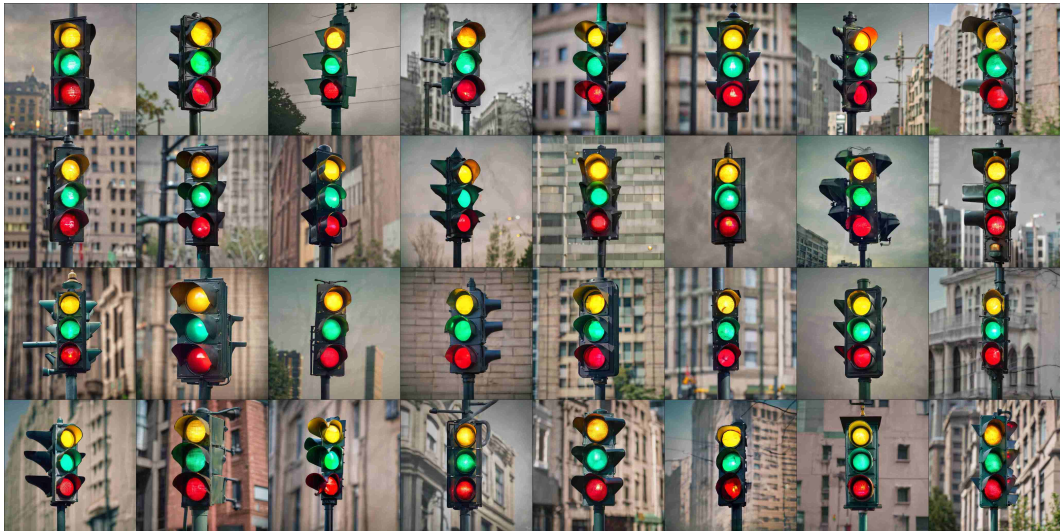


Figure 23: The generated images over each number of batch queries on the prompt "cat-butterfly" for each algorithms.

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457



(a) Pre-trained



(b) Fast Direct

Figure 24: The 32 randomly generated image for the prompt "traffic-light" guided by Fast Direct by utilizing 50 batch queries budget.

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511



(a) Pre-trained



(b) Fast Direct

Figure 25: The 32 randomly generated image for the prompt "apple" guided by Fast Direct by utilizing 50 batch queries budget.

1512
1513
1514
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



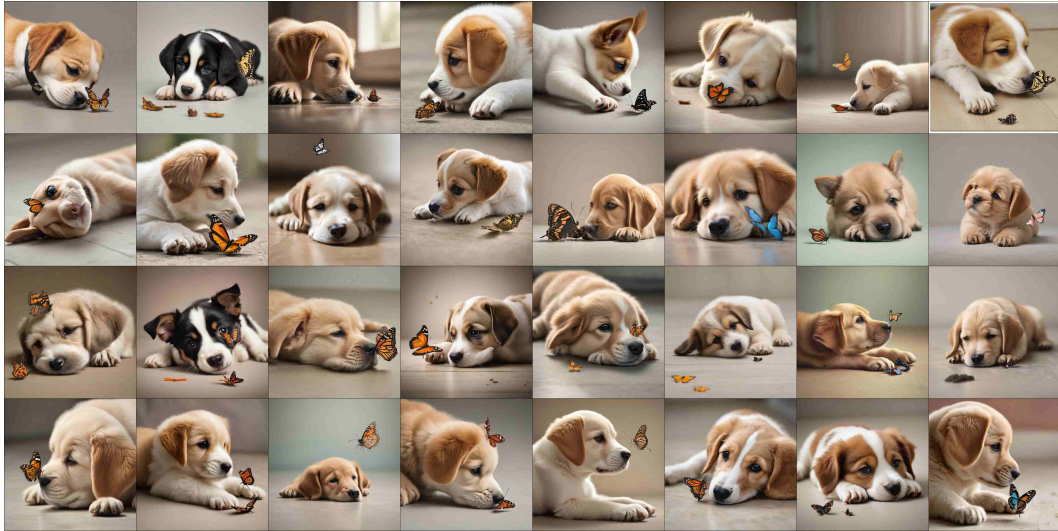
(a) Pre-trained



(b) Fast Direct

Figure 26: The 32 randomly generated image for the prompt "cyber-dog" guided by Fast Direct by utilizing 50 batch queries budget.

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619



(a) Pre-trained



(b) Fast Direct

Figure 27: The 32 randomly generated image for the prompt "puppy-nose" guided by Fast Direct by utilizing 50 batch queries budget.

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673



(a) Pre-trained



(b) Fast Direct

Figure 28: The 32 randomly generated image for the prompt "robot-plant" guided by Fast Direct by utilizing 50 batch queries budget.

1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727



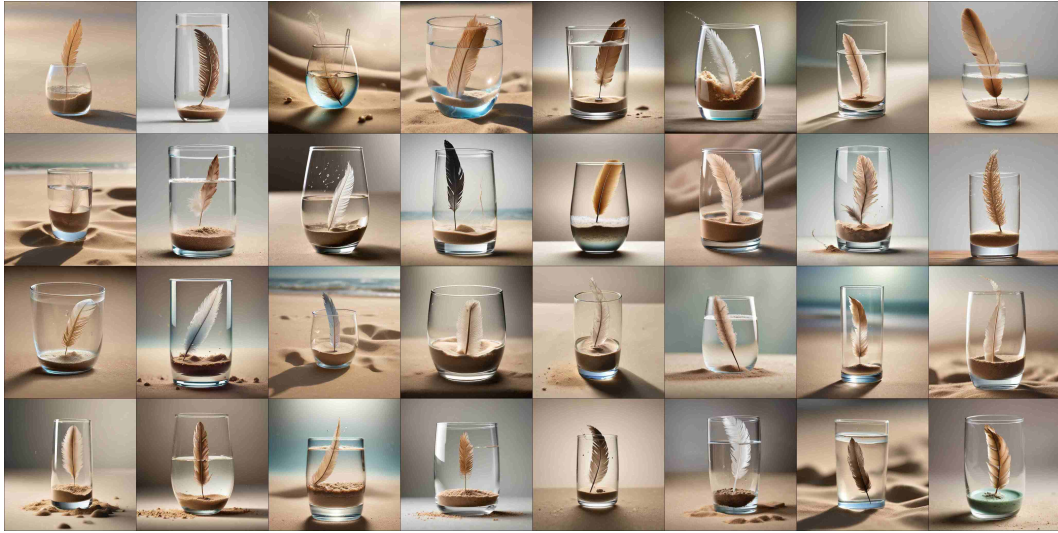
(a) Pre-trained



(b) Fast Direct

Figure 29: The 32 randomly generated image for the prompt "ocean" guided by Fast Direct by utilizing 50 batch queries budget.

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781



(a) Pre-trained



(b) Fast Direct

Figure 30: The 32 randomly generated image for the prompt "sand-glass" guided by Fast Direct by utilizing 50 batch queries budget.

1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835



(a) Pre-trained



(b) Fast Direct

Figure 31: The 32 randomly generated image for the prompt "penguin" guided by Fast Direct by utilizing 50 batch queries budget.

1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889



(a) Pre-trained



(b) Fast Direct

Figure 32: The 32 randomly generated image for the prompt "basket" guided by Fast Direct by utilizing 50 batch queries budget.

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943



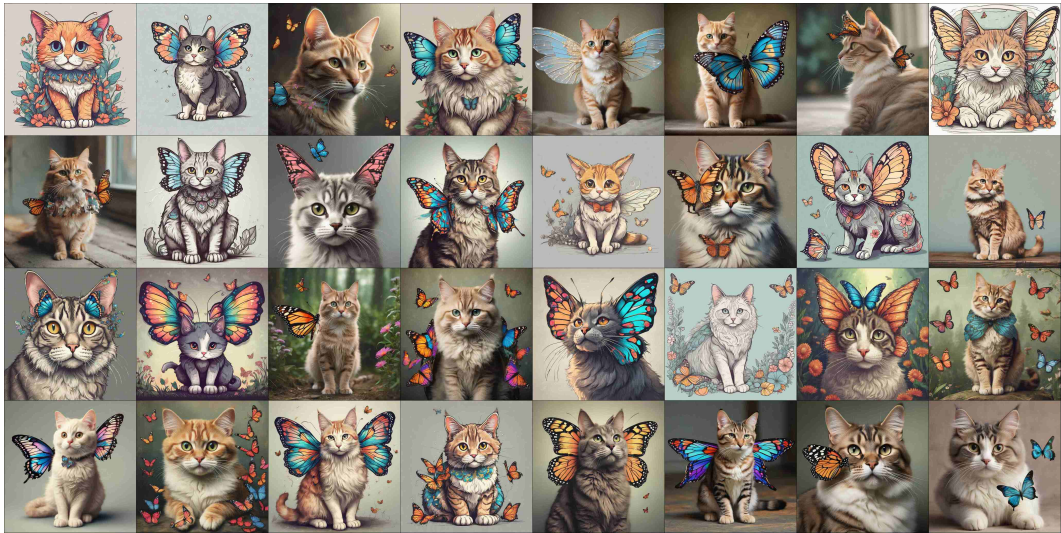
(a) Pre-trained



(b) Fast Direct

Figure 33: The 32 randomly generated image for the prompt "ice-cube" guided by Fast Direct by utilizing 50 batch queries budget.

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997



(a) Pre-trained



(b) Fast Direct

Figure 34: The 32 randomly generated image for the prompt "cat-butterfly" guided by Fast Direct by utilizing 50 batch queries budget.