

Alignment of Diffusion Model and Flow Matching for Text-to-Image Generation

Yidong Ouyang¹, Liyan Xie², Hongyuan Zha³, and Guang Cheng¹

¹Department of Statistics, University of California, Los Angeles

²Department of Industrial and Systems Engineering, University of Minnesota

³School of Data Science, The Chinese University of Hong Kong, Shenzhen

Abstract

Diffusion models and flow matching have demonstrated remarkable success in text-to-image generation. While many existing alignment methods primarily focus on fine-tuning pre-trained generative models to maximize a given reward function, these approaches require extensive computational resources and may not generalize well across different objectives. In this work, we propose a novel alignment framework by leveraging the underlying nature of the alignment problem—sampling from reward-weighted distributions—and show that it applies to both diffusion models (via score guidance) and flow matching models (via velocity guidance). The score function (velocity field) required for the reward-weighted distribution can be decomposed into the pre-trained score (velocity field) plus a conditional expectation of the reward. For the alignment on the diffusion model, we identify a fundamental challenge: the adversarial nature of the guidance term can introduce undesirable artifacts in the generated images. Therefore, we propose a finetuning-free framework that trains a guidance network to estimate the conditional expectation of the reward. We achieve comparable performance to finetuning-based models with one-step generation with at least a 60% reduction in computational cost. For the alignment on flow matching, we propose a training-free framework that improves the generation quality without additional computational cost.

1 Introduction

Diffusion models and flow matching have achieved impressive performance in text-to-image generation, as demonstrated by state-of-the-art models such as Imagen [36], DALL-E 3 [4], and Stable Diffusion [35]. These models have been proven capable of generating high-quality, creative images even from novel and complex text prompts.

Inspired by Reinforcement Learning from Human Feedback (RLHF) [28], many alignment approaches leverage preference pairs to fine-tune models for generating samples that align with task-specific objectives. RLHF-type methods typically learn a reward function and use the policy gradients to update the model [20, 12, 5, 10, 7, 17, 24, 18]. On the other hand, Direct Preference Optimization (DPO)-type methods directly optimize the model to adhere to human preferences, without requiring explicit reward modeling or reinforcement learning [34, 49, 53, 22, 54].

Despite their effectiveness, these approaches require modifying model parameters through fine-tuning, which comes with several potential limitations. For example, fine-tuning for new reward functions is computationally expensive and often requires carefully designed training strategies; otherwise, optimizing on a limited set of input prompts can limit generalization to unseen prompts. More importantly, existing fine-tuning approaches do not fully exploit the structure of the alignment

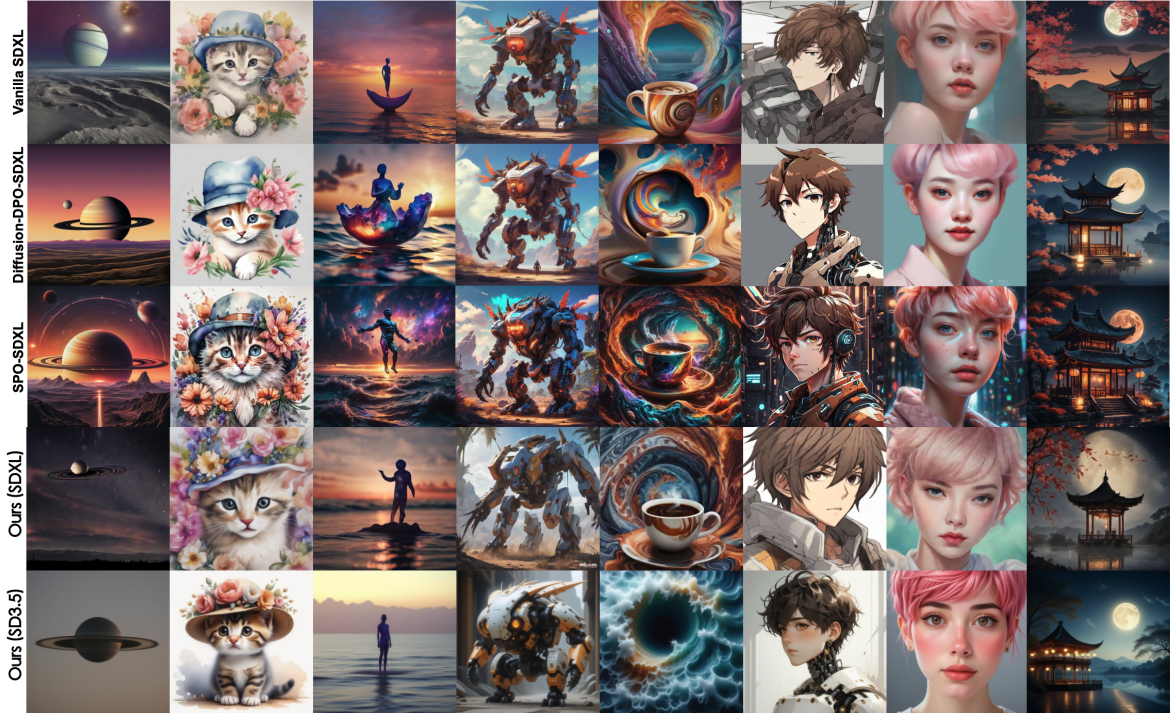


Figure 1: Qualitative comparison with Vanilla SDXL, Diffusion-DPO, and SPO. Our method achieves better aesthetic quality and stronger alignment with the text prompt. Prompts are provided in the Appendix C.3.

problem. Instead, they typically apply Low-Rank Adaptation (LoRA) to optimize model weights for a specific reward function [16], which may not be the most efficient strategy.

In contrast, plug-and-play alignment methods integrate new objectives without modifying the underlying model parameters, significantly reducing computational costs while adapting flexibly to different reward functions. In this paper, we cast alignment for both diffusion models and flow matching models as a unified sampling problem from reward-weighted distributions. Under this formulation, the key object needed for sampling—the new score function for diffusion or the new velocity field for flow matching—can be written as the corresponding pre-trained quantity plus an additional reward-driven guidance term.

For diffusion models, the guidance term admits an adversarial nature flaw, i.e., the guidance is the gradient of the log conditional expectation of the reward. Directly using the gradient of high dimensional input space can lead to undesirable artifacts in the generated images. To address this issue, we propose a finetuning-free alignment method that trains a lightweight guidance network to estimate the required conditional expectation, together with a regularization strategy that stabilizes the guidance landscape. We evaluate the effectiveness of the proposed method on four widely used criteria for text-to-image generation, and the proposed method achieves comparable performance to finetuning-based models in one-step generation while reducing computational cost by at least 60%.

For flow matching, we derive the exact form of velocity guidance and further propose a training-free estimator that directly computes the guidance term without additional model fine-tuning. The proposed method improves the generation quality without additional training overhead.

2 Preliminaries

In this section, we begin with a brief overview of diffusion models and flow matching in Section 2.1 and Section 2.2. We then review existing techniques for aligning pre-trained models with human preferences, and decompose the alignment procedure into two key components: reward learning for modeling human preferences in Section 2.3 and the alignment methods in Section 2.4.

2.1 Diffusion Models

Diffusion generative models are characterized by their forward and backward processes [15, 46]. The forward process gradually injects Gaussian noise into samples \mathbf{x}_0 from the data distribution following the stochastic differential equation:

$$d\mathbf{x}_t = \mathbf{f}(\mathbf{x}_t, t)dt + g(t)d\mathbf{w}, \quad t \in [0, T], \quad (1)$$

where \mathbf{w} is the standard Brownian motion, $\mathbf{f}(\cdot, t) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a drift coefficient, and $g(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$ is a diffusion coefficient. We use $p_t(\mathbf{x})$ to denote the marginal distribution of \mathbf{x}_t at time t . And we can use the time reversal of (1) for generation, which admits the following form [1]:

$$d\mathbf{x}_t = [\mathbf{f}(\mathbf{x}_t, t) - g(t)^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x})] dt + g(t)d\bar{\mathbf{w}}, \quad (2)$$

where $\bar{\mathbf{w}}$ is a standard Brownian motion when time flows backwards from T to 0, and dt is an infinitesimal negative time step. The score function of each marginal distribution $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$ needs to be estimated by the following score matching objective:

$$\min_{\theta} \mathbb{E}_t \left\{ \lambda(t) \mathbb{E}_{p_t(\mathbf{x}_t)} \left[\|\mathbf{s}_{\theta}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)\|_2^2 \right] \right\}, \quad (3)$$

where $\lambda(t) : [0, T] \rightarrow \mathbb{R}_{>0}$ is a positive weighting function, t is uniformly sampled over $[0, T]$. The latent diffusion model [35, 32] further extends diffusion models to text-to-image generation. They use an image encoder \mathcal{E} that maps \mathbf{x} into a latent representation and use a text encoder τ that maps the prompts y into an embedding as the condition.

2.2 Flow Matching

Flow matching models learn a time-dependent velocity field that transports a simple base distribution to the data distribution [23] via the probability flow ODE

$$\frac{d\mathbf{x}_t}{dt} = \mathbf{v}_{\phi}(\mathbf{x}_t, y, t), \quad t \in [0, 1],$$

where $\mathbf{v}_{\phi} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a learnable velocity field. Unlike diffusion models, we denote \mathbf{x}_0 as a sample from a base distribution (e.g., standard Gaussian) and \mathbf{x}_1 as a sample from the data distribution.

The flow matching objective minimizes the discrepancy between the model vector field and the oracle velocity field along the trajectory:

$$\mathcal{L}(\theta) = \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t | \mathbf{x}_0, \mathbf{x}_1), t \sim \mathcal{U}[0, 1]} \left[\|\mathbf{v}_{\phi}(\mathbf{x}_t, t) - \mathbf{v}(\mathbf{x}_t, y, t)\|_2^2 \right], \quad (4)$$

where \mathbf{x}_t is a linear interpolation between \mathbf{x}_0 and \mathbf{x}_1 , and $\mathbf{v}(\mathbf{x}_t, y, t)$ is the oracle velocity field.

2.3 Reward Learning

The Bradley-Terry (BT) model [6], and the more general Plackett-Luce ranking models [31, 26], are commonly used to model preferences. Given a prompt y and a pair of responses $\mathbf{x}_w \succ \mathbf{x}_l \mid y$, where \mathbf{x}_w denotes the winning response and \mathbf{x}_l denotes the losing response under the preference of humans. The BT model depicts the preference distribution as

$$p(\mathbf{x}_w \succ \mathbf{x}_l \mid y) = \frac{\exp(r(\mathbf{x}_w, y))}{\exp(r(\mathbf{x}_w, y)) + \exp(r(\mathbf{x}_l, y))},$$

where $r(\mathbf{x}, y)$ denotes the reward model and can be learned by the following maximum likelihood objective,

$$\min_{\phi} - \mathbb{E}_{(\mathbf{x}_w, \mathbf{x}_l, y) \sim \mathcal{D}} [\log \sigma(r(\mathbf{x}_w, y) - r(\mathbf{x}_l, y))], \quad (5)$$

where $\mathcal{D} = \{\mathbf{x}_w^{(i)}, \mathbf{x}_l^{(i)}, y^{(i)}\}_{i=1}^N$ is the offline preference dataset and σ denotes the logistic function.

2.4 Alignment

Building on the success of alignment techniques for finetuning large pre-trained models, many studies have explored aligning diffusion models and flow matching with human preferences. We review these approaches in the following.

Reinforcement Learning from Human Feedback. This type of works [20, 51, 12, 5, 10] finetune the pre-trained model π_{ref} by policy gradient objectives [17, 18]. In particular, the fine-tuned model π_{θ} is obtained by solving the following optimization problem:

$$\max_{\pi_{\theta}} \mathbb{E}_{y \sim \mathcal{D}_{\text{prompt}}, \mathbf{x} \sim \pi_{\theta}(\mathbf{x} \mid y)} [r(\mathbf{x}, y)] - \beta \mathbb{D}_{\text{KL}} [\pi_{\theta}(\mathbf{x} \mid y) \parallel \pi_{\text{ref}}(\mathbf{x} \mid y)], \quad (6)$$

where $\mathcal{D}_{\text{prompt}}$ denotes the prompt dataset. This type of method requires a pre-trained reward function for policy optimization [40].

Direct Preference Optimization. Rafailov et al. [34] propose not to explicitly learn the reward function. They start with the analytic solution of (6) as the energy-guided form,

$$\pi_r(\mathbf{x} \mid y) = \frac{1}{Z(y)} \pi_{\text{ref}}(\mathbf{x} \mid y) \exp\left(\frac{1}{\beta} r(\mathbf{x}, y)\right), \quad (7)$$

where $Z(y) = \int \pi_{\text{ref}}(\mathbf{x} \mid y) \exp\left(\frac{1}{\beta} r(\mathbf{x}, y)\right) d\mathbf{x}$ is the partition function. Therefore, they can reparameterize the reward function $r(\mathbf{x}, y)$ as

$$r(\mathbf{x}, y) = \beta \log \frac{\pi_r(\mathbf{x} \mid y)}{\pi_{\text{ref}}(\mathbf{x} \mid y)} + \beta \log Z(y). \quad (8)$$

Plugging (8) into (5) yields the objective of DPO-type methods:

$$\min - \mathbb{E}_{(\mathbf{x}_w, \mathbf{x}_l, y) \sim \mathcal{D}} \left[\log \sigma \left(\beta \log \frac{\pi_{\theta}(\mathbf{x}_w \mid y)}{\pi_{\text{ref}}(\mathbf{x}_w \mid y)} - \beta \log \frac{\pi_{\theta}(\mathbf{x}_l \mid y)}{\pi_{\text{ref}}(\mathbf{x}_l \mid y)} \right) \right]. \quad (9)$$

3 Proposed Framework for Diffusion Model

In this section, we introduce a finetuning-free frameworks to directly sample from the reward-guided distribution for diffusion models. We begin by introducing the methodology formulation in Section 3.1. We then provide an in-depth analysis of several vanilla methods for calculating the guidance in Section 3.2. We highlight that these vanilla guidance methods exhibit adversarial guidance, which generates undesirable artifacts and worsens performance, particularly in text-to-image generation. Then, we present an enhanced method in Section 3.3 that alleviates the problem.

3.1 Methodology Formulation

Inspired by previous works from transfer learning [29], we consider preference learning in terms of transferring a pre-trained diffusion model to adapt to the given preference data. To this end, we propose a finetuning-free alignment method for the diffusion models. Instead of using RLHF-type (like (6)) or DPO-type (like (9)) alignments, we propose to directly sample from the reward-weighted distribution $\pi_r(\mathbf{x}|y)$ in (7) leveraging the relationships between score functions in the following Theorem.



Figure 2: Illustration of the Adversarial Nature of Guidance. When the strength of the guidance is too small, there is little difference between the generated images with or without guidance. However, as the magnitude of the guidance increases (from left to right), undesirable artifacts become more pronounced. The prompt is "A 3D Rendering of a cockatoo wearing sunglasses. The sunglasses have a deep black frame with bright pink lenses. Fashion photography, volumetric lighting, CG rendering".

Theorem 3.1. *Let the conditional distribution of reference diffusion model $\pi_{ref}(\mathbf{x}|y)$ be denoted as distribution p and the reward-weighted distribution $\pi_r(\mathbf{x}|y)$ defined in (7) as distribution q . Under some mild assumption of the forward noising process detailed in Appendix A, let ϕ^* be the optimal solution for the conditional diffusion model trained on target domain $q(\mathbf{x}_0, y)$, i.e.,*

$$\phi^* = \arg \min_{\phi} \mathbb{E}_t \left\{ \lambda(t) \mathbb{E}_{q_t(\mathbf{x}_t, y)} \left[\left\| \mathbf{s}_{\phi}(\mathbf{x}_t, y, t) - \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t|y) \right\|_2^2 \right] \right\},$$

then

$$\mathbf{s}_\phi^*(\mathbf{x}_t, y, t) = \underbrace{\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|y)}_{\text{pre-trained conditional model on source}} + \underbrace{\nabla_{\mathbf{x}_t} \log \mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) \right]}_{\text{conditional guidance}}. \quad (10)$$

The proof can be found in Appendix A. Based on (10), we can calculate the additional guidance term rather than finetuning the text-to-image generative model. In general, the guidance term in (10) is not straightforward to compute as we need to sample from $p(\mathbf{x}_0|\mathbf{x}_t, y)$ for each \mathbf{x}_t in the generation process. In the following, we first discuss some existing ways to calculate the guidance term.

3.2 Vanilla Method to Compute the Guidance Term

M1: Direct backpropagate through diffusion process. The first method directly backpropagates through diffusion process to calculate $\nabla_{\mathbf{x}_t} \log \mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} [\exp(r(\mathbf{x}_0, y)/\beta)]$ for fine-tuning the diffusion model. In [45], the authors propose an unbiased Monte Carlo estimation:

$$\nabla_{\mathbf{x}_t} \log \mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\exp \left(\frac{1}{\beta} r(\mathbf{x}_0, y) \right) \right] \approx \nabla_{\mathbf{x}_t} \log \frac{1}{n} \sum_{i=1}^n \exp \left(\frac{1}{\beta} r(\mathbf{x}_0^i, y) \right),$$

where \mathbf{x}_0^i denotes the i -th sample drawn from $p(\mathbf{x}_0|\mathbf{x}_t, y)$. However, this Monte Carlo estimation significantly increases memory costs, especially in text-to-image generation. Inspired by recent studies [10], we can borrow the same techniques, e.g., accumulated gradients along the diffusion process using techniques such as low-rank adaptation (LoRA) [16] and truncation or gradient checkpointing [33, 10], to alleviate the memory cost of backpropagating through the diffusion process for calculating the guidance term. We can further reduce the memory cost by using the few-step diffusion model as the reference model. Despite these techniques, the memory requirements remain higher compared to the proposed approach.

M2: Approximate and apply Tweedie’s formula. The second method first approximates the guidance term by [9]:

$$\nabla_{\mathbf{x}_t} \log \mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\exp \left(\frac{1}{\beta} r(\mathbf{x}_0, y) \right) \right] \approx \frac{1}{\beta} \nabla_{\mathbf{x}_t} r(\mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} [\mathbf{x}_0], y). \quad (11)$$

Then, Tweedie’s formula is further applied by [2, 9, 56]:

$$\mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t, y] = \mathbf{x}_t + \sigma_t^2 \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|y).$$

However, as noted in [25, 45], the approximation used in (11) is biased, leading to an incorrect calculation of the guidance term.

In the following, we empirically evaluate the effectiveness of these methods for aligning text-to-image generation tasks. We first identify a previously overlooked issue that contributes to suboptimal alignment performance. Figure 2 illustrates the performance of two vanilla methods under the guidance of PickScore [19], a reward function that evaluates whether the generated images align with human aesthetic and semantic preferences. The x-axis represents the strength of the guidance term, denoted by α ¹. Our experiments reveal that tuning this hyperparameter presents significant

¹Although there is no α in (10), many guidance methods [25, 45] add this hyperparameter in practice to balance the strength of the guidance term with the score.

Table 1: Comparison of finetuning-free alignment algorithms on diffusion models. Our method uniquely provides theoretical guarantees for the correct form for guidance with a step size guarantee.

Method	Classifier Guidance	Direct backpropagate (M1)	Tweedie’s formula (M2)	Ours
Formulation	$\frac{1}{\beta} \nabla_{\mathbf{x}_t} r(\mathbf{x}_t, y)$	$\nabla_{\mathbf{x}_t} \log \frac{1}{n} \sum_{i=1}^n \exp\left(\frac{1}{\beta} r(\mathbf{x}_0^i, y)\right)$	$\frac{1}{\beta} \nabla_{\mathbf{x}_t} r(\mathbb{E}_{p(\mathbf{x}_0 \mathbf{x}_t, y)}[\mathbf{x}_0], y)$	$\nabla_{\mathbf{x}_t} \log h_{\psi^*}(\mathbf{x}_t, y, t)$
Unbiased	✗	✓	✗	✓
Step size guarantee	✗	✗	✗	✓

challenges. Insufficient values of α produce results indistinguishable from unguided generation, while excessive values introduce substantial artifacts that degrade image quality.

We attribute this phenomenon to the adversarial nature of the guidance mechanism, as observed in prior work [43]. In (10), the guidance term is directly added to the estimated score. If the landscape is not smooth or does not behave well², the adversarial nature of the guidance can lead to undesirable artifacts in the generated images. To address these limitations, our proposed framework provides theoretical guarantees for generating properly aligned distributions with a fixed strength parameter $\alpha = 1$. Furthermore, we develop an additional regularization technique for training the guidance network that mitigates these instability issues.

3.3 Proposed Finetuning-free Guidance for Diffison Models

We first utilize the following trick to calculate the conditional expectation, similar to previous works [29, 25].

Lemma 3.2. *For a neural network $h_{\psi}(\mathbf{x}_t, y, t)$ parameterized by ψ , define the objective*

$$\mathcal{L}_{\text{guidance}}(\psi) := \mathbb{E}_{p(\mathbf{x}_0, \mathbf{x}_t, y)} \left[\left\| h_{\psi}(\mathbf{x}_t, y, t) - \exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) \right\|_2^2 \right], \quad (12)$$

then its minimizer $\psi^* = \arg \min_{\psi} \mathcal{L}_{\text{guidance}}(\psi)$ satisfies:

$$h_{\psi^*}(\mathbf{x}_t, y, t) = \mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) \right].$$

By Lemma 3.2, we can instead estimate the value $\mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)}[\exp(r(\mathbf{x}_0, y)/\beta)]$ using the guidance network h_{ψ^*} obtained by minimizing the objective function $\mathcal{L}_{\text{guidance}}(\psi)$, which can be approximated by easy sampling from the joint distribution $p(\mathbf{x}_0, \mathbf{x}_t, y)$. Then, the estimated score function for the aligned diffusion model can be calculated as follows:

$$\mathbf{s}_{\phi^*}(\mathbf{x}_t, y, t) = \underbrace{\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|y)}_{\text{pre-trained model on source}} + \underbrace{\nabla_{\mathbf{x}_t} \log h_{\psi^*}(\mathbf{x}_t, y, t)}_{\text{guidance network}}. \quad (13)$$

To alleviate the adversarial nature of the guidance, we can adopt the consistency regularization $\mathcal{L}_{\text{consistence}}$ to learn the guidance network h_{ψ^*} better, i.e., the gradient of $\mathcal{L}_{\text{consistence}}(\mathbf{x}_t, y, t)$ with respect to \mathbf{x}_t should match the score in preferred data. The key point of this regularization is that we cannot easily change the landscape of a given predetermined reward function, but we can

²We use landscape to describe the change of reward given the change of images.

regularize the landscape of the learned guidance network to ensure the generation of high-quality images.

$$\begin{aligned}\psi^* &= \arg \min_{\psi} \mathcal{L}_{\text{consistence}} \\ &:= \mathbb{E}_{q(\mathbf{x}_0, y)} \mathbb{E}_{q(\mathbf{x}_t | \mathbf{x}_0)} \left[\left\| \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{x}_0, y) + \nabla_{\mathbf{x}_t} \log h_{\psi}(\mathbf{x}_t, y, t) - \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t | \mathbf{x}_0, y) \right\|_2^2 \right].\end{aligned}\quad (14)$$

Combining the consistency regularization terms together with the original guidance loss in (12), the final learning objective for the guidance network can be described as follows:

$$\psi^* = \arg \min_{\psi} \{ \mathcal{L}_{\text{guidance}} + \eta \mathcal{L}_{\text{consistence}} \}, \quad (15)$$

where $\eta \geq 0$ are hyperparameters that control the strength of additional regularization, which also enhances the flexibility of our solution scheme.

3.4 Further Improvement to One-step Generation

The training objectives in (12) and (14) are agnostic to the reference model, indicating that we can use any pre-trained diffusion model with any reward function, whether differentiable or not. Motivated by the computational efficiency of one-step generative models in practical applications, we further present a straightforward approach for applying our proposed finetuning-free guidance to one-step text-to-image models.

Specifically, instead of sampling t uniformly from $[0, T]$, we can simply set $t = T$. This small modification offers several advantages. First, while one-step diffusion models may not perform as well as few-step (2–4 step) models [37], we empirically find that with additional guidance, their performance improves significantly, as presented in Section 5.3. Second, as the guidance network h_{ψ} now becomes time-independent, we empirically observe that h_{ψ} is easy to train—with ten training epochs on the Pick-a-Pic V1 dataset, our guidance network produces high-quality images, which can be found in Section 5.2. We summarize the overall learning pipeline in Algorithm 1 in the Appendix. And we leave another two gradient-free designs for diffusion models in Appendix B.

4 Proposed Framework for Flow Matching

Training-free Alignment Framework for Flow Matching Given that state-of-the-art models are grounded in Diffusion Transformers [30] and flow matching [23], we present the exact form of flow-matching guidance in the theorem below.

Theorem 4.1. *Let ϕ_q^* be the optimal solution for the conditional flow matching model trained on target domain $q(\mathbf{x}_1, y)$ (where \mathbf{x}_1 are sampled from data distribution, $\mathbf{v}_q(\mathbf{x}_t, y, t)$ denotes the oracle velocity field on target distribution), i.e., ϕ_q^* equals*

$$\arg \min_{\phi} \mathbb{E}_t \left\{ \mathbb{E}_{q_t(\mathbf{x}_t, y)} \left[\left\| \mathbf{v}_{\phi}(\mathbf{x}_t, y, t) - \mathbf{v}_q(\mathbf{x}_t, y, t) \right\|_2^2 \right] \right\},$$

then

$$\mathbf{v}_{\phi_q^*}(\mathbf{x}_t, y, t) = \mathbf{v}_{\phi_p}(\mathbf{x}_t, y, t) + \mathbb{E}_{\mathbf{x}_1 \sim p_{1|t}(\mathbf{x}_1 | \mathbf{x}_t, y)} \left[(R(\mathbf{x}_1, \mathbf{x}_t, y) - 1) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_1, y) \right], \quad (16)$$

where

$$R(\mathbf{x}_1, \mathbf{x}_t, y) = \frac{\exp\left(\frac{1}{\beta} r(\mathbf{x}_1, y)\right)}{\mathbb{E}_{\mathbf{x}'_1 \sim p_{1|t}(\mathbf{x}_1 | \mathbf{x}_t, y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}'_1, y)\right) \right]}.$$

Table 2: Benchmark comparison of different methods on text-to-image alignment. Results are grouped by base model.

Type	Method	PickScore	HPSV2	ImageReward	Aesthetic	Training GPU Hour
Base Model: SDXL						
Baseline	SDXL	21.95	26.95	0.5380	5.950	–
Training-free	Direct backpropagate	21.84	27.53	0.5870	5.922	–
	Tweedie’s formula	22.34	28.76	0.9501	6.002	–
Finetuning-based	Diff.-DPO	22.64	29.31	0.9436	6.015	4800
	SPO	23.06	31.80	1.0803	6.364	234
Finetuning-free	Ours	23.08	32.12	1.0625	6.452	92
Base Model: SD3.5 Large Turbo						
Baseline	SD3.5 Large Turbo	22.30	30.29	1.0159	6.5190	–
Finetuning-free	Ours	23.14	32.31	1.1025	6.5280	–

Estimation of the Guidance Term for Flow Matching Different from the guidance term of diffusion models in (10), the guidance for flow matching in (16) does not have the adversarial problem. The guidance term is a conditional expectation without the gradient operator.

To enable fast sampling, we would like to use importance sampling to convert the conditional expectation under $p(\mathbf{x}_1 | \mathbf{x}_t, y)$ into an expectation under $p(\mathbf{x}_1 | y)$. After the detailed derivation in A.2, we can calculate the guidance term of flow matching by

$$\mathbb{E}_{\mathbf{x}_1 \sim p(\mathbf{x}_1 | y)} \left[\left(\frac{\exp(\frac{1}{\beta} r(\mathbf{x}_1, y))}{\mathbb{E}_{\mathbf{x}_1} \left[\exp(\frac{1}{\beta} r(\mathbf{x}_1, y)) \frac{p_{t|1}(\mathbf{x}_t | \mathbf{x}_1, y)}{\mathbb{E}_{\mathbf{x}_1} [p_{t|1}(\mathbf{x}_t | \mathbf{x}_1, y)]} \right]} - 1 \right) v_t(\mathbf{x}_t | \mathbf{x}_1, y) \frac{p_{t|1}(\mathbf{x}_t | \mathbf{x}_1, y)}{\mathbb{E}_{\mathbf{x}_1} [p_{t|1}(\mathbf{x}_t | \mathbf{x}_1, y)]} \right].$$

Therefore, we do not need to sample \mathbf{x}_1 with multiple function evaluations, but just sample from the marginal data distribution. Compared with the finetuning-free method proposed in (13), this formulation is training-free and offers greater computational efficiency.

5 Experimental Results

In this section, we present a comprehensive experimental evaluation, demonstrating the effectiveness of our two frameworks for sampling directly from reward-guided distributions. We first outline our experimental setup and evaluation criteria in Section 5.1, followed by benchmark results against state-of-the-art methods in Section 5.2. Finally, we provide an in-depth ablation study that validates our key theoretical claims and demonstrates the superior performance of our guidance network in Section 5.3.

5.1 Experimental Setup

For the experiments on diffusion models, we follow the official configurations recommended for SPO [22], Diffusion-DPO [49], and MAPO [42]. Diffusion-DPO and MAPO are fine-tuned on the Pick-a-Pic V2 dataset, which contains over 800k image preference pairs. In contrast, SPO is fine-tuned online using 4k text prompts (without images) randomly selected from Pick-a-Pic V1. Our method trains the guidance network offline using 583k image preference pairs from Pick-a-Pic V1. Overall, our method and the competing models in the text-to-image alignment benchmark

are trained on comparable datasets, allowing for a fair comparison. We adopt Stable Diffusion XL (SDXL)-Turbo [38] as the reference model for one-step text-to-image generation. For the experiments on flow matching, we adopt the state-of-the-art SD3.5 Large Turbo [11] as the backbone. The official recommendation for the number of sampling steps is four to eight, and we use four steps for all experiments.

Implementation Details. In the following, we provide the training details for the guidance network of the diffusion model. Since the guidance network takes noisy images \mathbf{x}_T and prompts y as input and outputs a scalar value, we adopt the same variational autoencoder (VAE), tokenizer, and text encoder from the reference diffusion model for encoding image and text. Consequently, the trainable parameters of our guidance network are quite small. In practice, we adopt two convolutional layers for processing VAE-encoded feature maps and a five-layer multi-layer perceptron (MLP) to project the image and text embedding to a scalar. The total parameter size of the guidance network is only 72 MB, making it lightweight and easy to train. We train the guidance network on the Pick-a-Pic training dataset for 10 epochs with batch size 32, Adam optimizer, learning rate $1e-3$, and hyperparameters $\eta = 1$.

Evaluation Criterion. Following established evaluation protocols [49, 22], we report quantitative results using 500 validation prompts from the validation unique split of Pick-a-Pic. We adopt four evaluation criteria to evaluate different aspects of image quality. PickScore [19] measures overall human preference by aggregating judgments on aesthetic appeal, coherence, and realism. HPSV2 [50] assesses prompt adherence, ensuring the generated image accurately reflects the given textual description. ImageReward [51] quantifies human preference based on fine-grained attributes such as composition, detail preservation, and semantic relevance. Lastly, the aesthetic evaluation model from LAION [39] focuses on visual appeal, capturing factors such as color harmony, style, and artistic quality.

5.2 Experimental Results

As shown in Table 2, our method surpasses baseline approaches across four evaluation criteria, demonstrating the effectiveness of the two proposed frameworks in enhancing text-to-image alignment. The improvements are observed in both perceptual quality and semantic coherence, indicating that our guidance network successfully refines image generation to better match textual descriptions. This performance gain highlights the advantages of our lightweight architecture and the optimization strategy used during training. Figure 1 provides a qualitative comparison with baseline methods, further illustrating the superior visual fidelity and text alignment achieved by our approach.

5.3 Ablation study

In this section, we first verify the advantages of our proposed method against other finetuning-free guidance methods as summarized in Table 1. We then analyze the impact of few-step (2–4 step) generation compared to one-step generation, highlighting how our guidance term significantly enhances performance.

As illustrated in Figure 3, vanilla guidance methods struggle to induce meaningful improvements in generated images, even with carefully tuned guidance strength. Increasing the guidance parameter α often leads to undesirable artifacts rather than quality improvements. In contrast, our method effectively enhances image generation by leveraging a regularized guidance network, demonstrating its ability to refine scene details and improve alignment with input prompts.



Figure 3: Effectiveness of the proposed method for diffusion models: The results demonstrate that 2-step and 3-step generation significantly improve the quality of the generated images compared to one-step generation. While two vanilla guidance methods (Tweedie’s formula or directly backpropagation summarized in Section 3.2) fail to produce meaningful changes in the scene despite appropriate guidance strength, our method successfully achieves this enhancement. The prompt is “A photo of a frog holding an apple while smiling in the forest”.

Table 3: Ablation study comparing the performance of our method with no guidance and two vanilla guidance methods under one-step and multi-step generation. Our method outperforms all baselines, which demonstrates the effectiveness of our guidance network in refining image quality and prompt alignment.

Method	PickScore
Ours (1 step)	23.08
No guidance (1 step)	22.14
Tweedie’s (1 step)	22.34
Backpropagate (1 step)	21.84
No guidance (2 steps)	22.64
No guidance (3 steps)	22.56

To further explore this, we examine the performance of our method against two vanilla guidance techniques, Tweedie’s and Backpropagate, as well as the no guidance baseline, all under a one-step sampling condition. As shown in Table 3, our method achieves the highest PickScore. This demonstrates that our regularized guidance network provides a substantial improvement over no guidance scenario and traditional methods. Consistent with prior studies, increasing the number of steps from one to two or three results in improved image quality, as shown in Figure 3 and Table 3. However, our method enables one-step generation to achieve performance even better than 2- or 3-step generation, highlighting the power of our guidance network. In Appendix C, we include the sensitive analysis of the regularization strength.

6 Related Work

Existing alignment methods can be broadly categorized into two approaches: RLHF-based method that uses policy gradient to update the diffusion models and flow matching, and DPO-based methods that use a parametrization trick to update the diffusion models without explicitly learning the

reward function.

RLHF-based alignment of diffusion model and flow matching. Lee et al. [20] first train a reward model to predict human feedback and adopt a reward-weighted finetuning objective to align the diffusion model. In [12, 5], diffusion models are updated using policy gradient algorithms under Kullback–Leibler (KL) constraints. Clark et al. [10] propagate gradients of the reward function through the full sampling procedure, and reduce memory costs by adopting low-rank adaptation (LoRA) [16] and gradient checkpointing [8]. In [24, 21, 52, 14], the authors improve GRPO [41] for the alignment of flow matching.

DPO-based alignment of diffusion model. A line of work [49, 53] directly applies DPO [34] to align the diffusion model with human preference. Liang et al. [22] propose a step-aware preference model and a step-wise resampler to align the preference optimization target with the denoising performance at each timestep. Yang et al. [54] take on a finer dense reward perspective and derive a tractable alignment objective that emphasizes the initial steps.

Training-free guidance. This line of work [9, 13, 25, 45, 2, 56, 43, 55] explores the use of diffusion models as plug-and-play priors for solving inverse problems. Some work [43, 47, 48, 27, 44] study inference-time optimization for alignment. However, to the best of our knowledge, there has been limited exploration of applying guidance on diffusion models to address the challenge of text-to-image alignment in the context of one-step generation. Also, there has been limited exploration of training-free guidance on flow matching of text-to-image alignment. This gap motivates our work.

7 Conclusion

In this paper, we introduced two novel framework for aligning text-to-image diffusion models and flow matching models with human preferences. By formulating alignment as sampling from a reward-weighted distribution, our approach leverages a plug-and-play guidance mechanism. Specifically, we decomposed the score function (velocity field) of the reward-weighted distribution into the pre-trained score (velocity field) plus a reward-driven guidance term. For diffusion models, we identify that the adversarial nature of the guidance can introduce undesirable artifacts, and we propose a finetuning-free approach that trains a lightweight guidance network to estimate the conditional expectation of the reward, together with a regularization strategy that stabilizes the guidance landscape. Empirically, our method achieves performance comparable to finetuning-based approaches for one-step generation while reducing computational cost by at least 60%. For flow matching, we derive the exact form of velocity guidance and propose a training-free estimator that improves generation quality without additional training.

References

- [1] Brian. D. O. Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12:313–326, 1982.
- [2] Arpit Bansal, Hong-Min Chu, Avi Schwarzschild, Soumyadip Sengupta, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Universal guidance for diffusion models. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 843–852, 2023.

- [3] Georgios Batzolis, Jan Stanczuk, Carola-Bibiane Schönlieb, and Christian Etmann. Conditional image generation with score-based diffusion models. *arXiv preprint arXiv:2111.13606*, 2021.
- [4] James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2023.
- [5] Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion models with reinforcement learning. In *The Twelfth International Conference on Learning Representations*, 2024.
- [6] Ralph Allan Bradley and Milton E. Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- [7] Souradip Chakraborty, Jiahao Qiu, Hui Yuan, Alec Koppel, Dinesh Manocha, Furong Huang, Amrit Bedi, and Mengdi Wang. MaxMin-RLHF: Alignment with diverse human preferences. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 6116–6135. PMLR, 2024.
- [8] Tianqi Chen, Bing Xu, Chiyuan Zhang, and Carlos Guestrin. Training deep nets with sublinear memory cost. *arXiv preprint arXiv:1604.06174*, 2016.
- [9] Hyungjin Chung, Jeongsol Kim, Michael Thompson Mccann, Marc Louis Klasky, and Jong Chul Ye. Diffusion posterior sampling for general noisy inverse problems. In *The Eleventh International Conference on Learning Representations*, 2023.
- [10] Kevin Clark, Paul Vicol, Kevin Swersky, and David J. Fleet. Directly fine-tuning diffusion models on differentiable rewards. In *The Twelfth International Conference on Learning Representations*, 2024.
- [11] Patrick Esser, Sumith Kulal, A. Blattmann, Rahim Entezari, Jonas Muller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English, Kyle Lacey, Alex Goodwin, Yannik Marek, and Robin Rombach. Scaling rectified flow transformers for high-resolution image synthesis. *ICML*, 2024.
- [12] Ying Fan, Olivia Watkins, Yuqing Du, Hao Liu, Moonkyung Ryu, Craig Boutilier, Pieter Abbeel, Mohammad Ghavamzadeh, Kangwook Lee, and Kimin Lee. Dpok: Reinforcement learning for fine-tuning text-to-image diffusion models. In *Advances in Neural Information Processing Systems*, volume 36, pages 79858–79885, 2023.
- [13] Alexandros Graikos, Nikolay Malkin, Nebojsa Jojic, and Dimitris Samaras. Diffusion models as plug-and-play priors. In *Advances in Neural Information Processing Systems*, volume 35, pages 14715–14728, 2022.
- [14] Xiaoxuan He, Siming Fu, Yuke Zhao, Wanli Li, Jian Yang, Dacheng Yin, Fengyun Rao, and Bo Zhang. Tempflow-grpo: When timing matters for grpo in flow models. *ArXiv*, abs/2508.04324, 2025.
- [15] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, volume 33, pages 6840–6851, 2020.

- [16] Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022.
- [17] Natasha Jaques, Shixiang Shane Gu, Dzmitry Bahdanau, José Miguel Hernández-Lobato, Richard E. Turner, and Douglas Eck. Sequence tutor: Conservative fine-tuning of sequence generation models with kl-control. In *International Conference on Machine Learning*, 2016.
- [18] Natasha Jaques, Judy Hanwen Shen, Asma Ghandeharioun, Craig Ferguson, Àgata Lapedriza, Noah Jones, Shixiang Shane Gu, and Rosalind Picard. Human-centric dialog training via offline reinforcement learning. In *Proceedings of the 2020 conference on empirical methods in natural language processing (EMNLP)*, pages 3985–4003, 2020.
- [19] Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-a-pic: An open dataset of user preferences for text-to-image generation. In *Advances in Neural Information Processing Systems*, volume 36, pages 36652–36663, 2023.
- [20] Kimin Lee, Hao Liu, Moonkyung Ryu, Olivia Watkins, Yuqing Du, Craig Boutilier, P. Abbeel, Mohammad Ghavamzadeh, and Shixiang Shane Gu. Aligning text-to-image models using human feedback. *arXiv preprint arXiv:2302.12192*, 2023.
- [21] Junzhe Li, Yutao Cui, Tao Huang, Yinpeng Ma, Chun Fan, Miles Yang, and Zhao Zhong. Mixgrpo: Unlocking flow-based grpo efficiency with mixed ode-sde. *ArXiv*, abs/2507.21802, 2025.
- [22] Zhanhao Liang, Yuhui Yuan, Shuyang Gu, Bohan Chen, Tiankai Hang, Mingxi Cheng, Ji Li, and Liang Zheng. Aesthetic post-training diffusion models from generic preferences with step-by-step preference optimization. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 13199–13208, 2025.
- [23] Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. In *The Eleventh International Conference on Learning Representations*, 2023.
- [24] Jie Liu, Gongye Liu, Jiajun Liang, Yangguang Li, Jiaheng Liu, Xintao Wang, Pengfei Wan, Di Zhang, and Wanli Ouyang. Flow-grpo: Training flow matching models via online rl. *NeurIPS*, 2025.
- [25] Cheng Lu, Huayu Chen, Jianfei Chen, Hang Su, Chongxuan Li, and Jun Zhu. Contrastive energy prediction for exact energy-guided diffusion sampling in offline reinforcement learning. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202, pages 22825–22855, 2023.
- [26] R Duncan Luce. *Individual Choice Behavior: A Theoretical Analysis*. Wiley New York, 1959.
- [27] Nanye Ma, Shangyuan Tong, Haolin Jia, Hexiang Hu, Yu-Chuan Su, Mingda Zhang, Xuan Yang, Yandong Li, T. Jaakkola, Xuhui Jia, and Saining Xie. Inference-time scaling for diffusion models beyond scaling denoising steps. *arXiv preprint arXiv:2501.09732*, 2025.
- [28] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano,

- Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*, volume 35, pages 27730–27744, 2022.
- [29] Yidong Ouyang, Liyan Xie, Hongyuan Zha, and Guang Cheng. Transfer learning for diffusion models. In *Advances in Neural Information Processing Systems*, volume 37, pages 136962–136989, 2024.
- [30] William S. Peebles and Saining Xie. Scalable diffusion models with transformers. *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4172–4182, 2022.
- [31] Robin L. Plackett. The analysis of permutations. *Journal of The Royal Statistical Society Series C-applied Statistics*, 24:193–202, 1975.
- [32] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. SDXL: Improving latent diffusion models for high-resolution image synthesis. In *The Twelfth International Conference on Learning Representations*, 2024.
- [33] Mihir Prabhudesai, Anirudh Goyal, Deepak Pathak, and Katerina Fragkiadaki. Aligning text-to-image diffusion models with reward backpropagation. *arXiv preprint arXiv:2310.03739*, 2023.
- [34] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36, 2024.
- [35] Robin Rombach, A. Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10674–10685, 2021.
- [36] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. In *Advances in Neural Information Processing Systems*, volume 35, pages 36479–36494, 2022.
- [37] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. In *International Conference on Learning Representations*, 2022.
- [38] Axel Sauer, Dominik Lorenz, A. Blattmann, and Robin Rombach. Adversarial diffusion distillation. In *European Conference on Computer Vision*, 2023.
- [39] Christoph Schuhmann. Laion-aesthetics, 2022.
- [40] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [41] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Jun-Mei Song, Mingchuan Zhang, Y. K. Li, Yu Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *ArXiv*, abs/2402.03300, 2024.
- [42] Shuaijie She, Shujian Huang, Wei Zou, Wenhao Zhu, Xiang Liu, Xiang Geng, and Jiajun Chen. MAPO: Advancing multilingual reasoning through multilingual alignment-as-preference optimization. In *Annual Meeting of the Association for Computational Linguistics*, 2024.

- [43] Yifei Shen, Xinyang Jiang, Yifan Yang, Yezhen Wang, Dongqi Han, and Dongsheng Li. Understanding and improving training-free loss-based diffusion guidance. In *Advances in Neural Information Processing Systems*, volume 37, pages 108974–109002, 2024.
- [44] Raghav Singhal, Zachary Horvitz, Ryan Teehan, Mengye Ren, Zhou Yu, Kathleen McKeown, and Rajesh Ranganath. A general framework for inference-time scaling and steering of diffusion models. In *Forty-second International Conference on Machine Learning*, 2025.
- [45] Jiaming Song, Qincheng Zhang, Hongxu Yin, Morteza Mardani, Ming-Yu Liu, Jan Kautz, Yongxin Chen, and Arash Vahdat. Loss-guided diffusion models for plug-and-play controllable generation. In *International Conference on Machine Learning*, 2023.
- [46] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021.
- [47] Zhiwei Tang, Jiangweizhi Peng, Jiasheng Tang, Mingyi Hong, Fan Wang, and Tsung-Hui Chang. Inference-time alignment of diffusion models with direct noise optimization. In *Forty-second International Conference on Machine Learning*, 2025.
- [48] Masatoshi Uehara, Yulai Zhao, Tommaso Biancalani, and Sergey Levine. Understanding reinforcement learning-based fine-tuning of diffusion models: A tutorial and review. *arXiv preprint arXiv:2407.13734*, 2024.
- [49] Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou, Aaron Lou, Senthil Purushwalkam, Stefano Ermon, Caiming Xiong, Shafiq R. Joty, and Nikhil Naik. Diffusion model alignment using direct preference optimization. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8228–8238, 2023.
- [50] Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis. *arXiv preprint arXiv:2306.09341*, 2023.
- [51] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation. In *Advances in Neural Information Processing Systems*, volume 36, pages 15903–15935, 2023.
- [52] Zeyue Xue, Jie Wu, Yu Gao, Fangyuan Kong, Lingting Zhu, Mengzhao Chen, Zhiheng Liu, Wei Liu, Qiushan Guo, Weilin Huang, and Ping Luo. Dancegrpo: Unleashing grpo on visual generation. *ArXiv*, abs/2505.07818, 2025.
- [53] Kai Yang, Jian Tao, Jiafei Lyu, Chunjiang Ge, Jiabin Chen, Qimai Li, Wei Han Shen, Xiaolong Zhu, and Xiu Li. Using human feedback to fine-tune diffusion models without any reward model. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8941–8951, 2023.
- [54] Shentao Yang, Tianqi Chen, and Mingyuan Zhou. A dense reward view on aligning text-to-image diffusion with preference. In *Forty-first International Conference on Machine Learning*, 2024.
- [55] Haotian Ye, Haowei Lin, Jiaqi Han, Minkai Xu, Sheng Liu, Yitao Liang, Jianzhu Ma, James Zou, and Stefano Ermon. TFG: Unified training-free guidance for diffusion models. In *Advances in Neural Information Processing Systems*, volume 37, pages 22370–22417, 2024.

- [56] Jiwen Yu, Yinhuai Wang, Chen Zhao, Bernard Ghanem, and Jian Zhang. Freedom: Training-free energy-guided conditional diffusion model. *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 23117–23127, 2023.

A Theoretical Details for Section 3

A.1 Proof of Theorem 3.1

We first restate the complete theorem as follows:

Theorem A.1. *Let the conditional distribution of reference diffusion model $\pi_{\text{ref}}(\mathbf{x}|y)$ be denoted as distribution p and the reward-weighted distribution $\pi_r(\mathbf{x}|y)$ defined in (7) as distribution q . Assume \mathbf{x}_t and y are conditionally independent given \mathbf{x}_0 in the forward process, i.e., $p(\mathbf{x}_t|\mathbf{x}_0, y) = p(\mathbf{x}_t|\mathbf{x}_0)$, $\forall t \in [0, T]$. Additionally, assume the forward process on the reward-weighted distribution is identical to that on the reference distribution $q(\mathbf{x}_t|\mathbf{x}_0) = p(\mathbf{x}_t|\mathbf{x}_0)$ ³, and ϕ^* is the optimal solution for the conditional diffusion model trained on target domain $q(\mathbf{x}_0, y)$, i.e.,*

$$\phi^* = \arg \min_{\phi} \mathbb{E}_t \left\{ \lambda(t) \mathbb{E}_{q_t(\mathbf{x}_t, y)} \left[\left\| \mathbf{s}_{\phi}(\mathbf{x}_t, y, t) - \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t|y) \right\|_2^2 \right] \right\}, \quad (17)$$

then

$$\mathbf{s}_{\phi^*}(\mathbf{x}_t, y, t) = \underbrace{\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|y)}_{\text{pre-trained conditional model on source}} + \underbrace{\nabla_{\mathbf{x}_t} \log \mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) \right]}_{\text{conditional guidance}}. \quad (18)$$

Proof. The proof is based on the theoretical framework of [29]. For the ease of readers, we incorporate the relevant conclusion from their work as lemmas below. To prove (18), we first build the connection between the Conditional Score Matching on the target domain and Importance Weighted Conditional Denoising Score Matching on the source domain in the following Lemma:

Lemma A.2. *Conditional Score Matching on the target domain is equivalent to Importance Weighted Denoising Score Matching on the source domain, i.e.,*

$$\begin{aligned} \phi^* &= \arg \min_{\phi} \mathbb{E}_t \left\{ \lambda(t) \mathbb{E}_{q_t(\mathbf{x}_t, y)} \left[\left\| \mathbf{s}_{\phi}(\mathbf{x}_t, y, t) - \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t|y) \right\|_2^2 \right] \right\} \\ &= \arg \min_{\phi} \mathbb{E}_t \left\{ \lambda(t) \mathbb{E}_{p(\mathbf{x}_0, y)} \mathbb{E}_{p(\mathbf{x}_t|\mathbf{x}_0)} \left[\left\| \mathbf{s}_{\phi}(\mathbf{x}_t, y, t) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{x}_0) \right\|_2^2 \frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)} \right] \right\}. \end{aligned}$$

Proof of Lemma A.2. We first connect the Conditional Score Matching objective in the target domain to the Conditional Denoising Score Matching objective in target distribution, which is proven by [3], i.e.,

$$\begin{aligned} \phi^* &= \arg \min_{\phi} \mathbb{E}_t \left\{ \lambda(t) \mathbb{E}_{q_t(\mathbf{x}_t, y)} \left[\left\| \mathbf{s}_{\phi}(\mathbf{x}_t, y, t) - \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t|y) \right\|_2^2 \right] \right\} \\ &= \arg \min_{\phi} \mathbb{E}_t \left\{ \lambda(t) \mathbb{E}_{q(\mathbf{x}_0, y)} \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} \left[\left\| \mathbf{s}_{\phi}(\mathbf{x}_t, y, t) - \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t|\mathbf{x}_0) \right\|_2^2 \right] \right\}. \end{aligned}$$

Then we split the mean squared error of the Conditional Denoising Score Matching objective on the target distribution into three terms as follows:

$$\begin{aligned} &\mathbb{E}_{q(\mathbf{x}_0, y)} \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} \left[\left\| \mathbf{s}_{\phi}(\mathbf{x}_t, y, t) - \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t|\mathbf{x}_0) \right\|_2^2 \right] \\ &= \mathbb{E}_{q(\mathbf{x}_0, \mathbf{x}_t, y)} \left[\left\| \mathbf{s}_{\phi}(\mathbf{x}_t, y, t) \right\|_2^2 \right] - 2 \mathbb{E}_{q(\mathbf{x}_0, \mathbf{x}_t, y)} \left[\langle \mathbf{s}_{\phi}(\mathbf{x}_t, y, t), \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t|\mathbf{x}_0) \rangle \right] + C_1, \end{aligned} \quad (19)$$

³These two assumptions are mild since \mathbf{x}_0 contains all information about y and $p(\mathbf{x}_t|\mathbf{x}_0)$ and $q(\mathbf{x}_t|\mathbf{x}_0)$ are forward noising process, which is easy to control.

where $C_1 = \mathbb{E}_{q(\mathbf{x}_0, \mathbf{x}_t, y)} [\|\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t | \mathbf{x}_0)\|_2^2]$ is a constant independent with ϕ , and $q(\mathbf{x}_t | \mathbf{x}_0, y) = q(\mathbf{x}_t | \mathbf{x}_0)$ because of conditional independent of \mathbf{x}_t and y given \mathbf{x}_0 by assumption. We can similarly split the mean squared error of Denoising Score Matching on the source domain into three terms as follows:

$$\begin{aligned} & \mathbb{E}_{p(\mathbf{x}_0, y)} \mathbb{E}_{p(\mathbf{x}_t | \mathbf{x}_0)} \left[\|\mathbf{s}_\phi(\mathbf{x}_t, y, t) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{x}_0)\|_2^2 \frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)} \right] \\ &= \mathbb{E}_{p(\mathbf{x}_0, \mathbf{x}_t, y)} \left[\|\mathbf{s}_\phi(\mathbf{x}_t, y, t)\|_2^2 \frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)} \right] - 2 \mathbb{E}_{p(\mathbf{x}_0, \mathbf{x}_t, y)} \left[\langle \mathbf{s}_\phi(\mathbf{x}_t, y, t), \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{x}_0) \rangle \frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)} \right] \quad (20) \\ &+ C_2, \end{aligned}$$

where C_2 is a constant independent with ϕ .

It is obvious to show that the first term in (19) is equal to the first term in (20), i.e.,

$$\begin{aligned} & \mathbb{E}_{p(\mathbf{x}_0, \mathbf{x}_t, y)} \left[\|\mathbf{s}_\phi(\mathbf{x}_t, y, t)\|_2^2 \frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)} \right] \\ &= \int_{\mathbf{x}_0} \int_{\mathbf{x}_t} \int_y p(\mathbf{x}_0, y) p(\mathbf{x}_t | \mathbf{x}_0) \|\mathbf{s}_\phi(\mathbf{x}_t, y, t)\|_2^2 \frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)} d\mathbf{x}_0 d\mathbf{x}_t dy \\ &= \int_{\mathbf{x}_0} \int_{\mathbf{x}_t} \int_y p(\mathbf{x}_0, y) q(\mathbf{x}_t | \mathbf{x}_0) \|\mathbf{s}_\phi(\mathbf{x}_t, y, t)\|_2^2 \frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)} d\mathbf{x}_0 d\mathbf{x}_t dy \\ &= \int_{\mathbf{x}_0} \int_{\mathbf{x}_t} \int_y q(\mathbf{x}_0, \mathbf{x}_t, y) \|\mathbf{s}_\phi(\mathbf{x}_t, y, t)\|_2^2 d\mathbf{x}_0 d\mathbf{x}_t dy \\ &= \mathbb{E}_{q(\mathbf{x}_0, \mathbf{x}_t, y)} [\|\mathbf{s}_\phi(\mathbf{x}_t, y, t)\|_2^2]. \end{aligned}$$

And the second term is also equivalent:

$$\begin{aligned} & \mathbb{E}_{p(\mathbf{x}_0, \mathbf{x}_t, y)} \left[\langle \mathbf{s}_\phi(\mathbf{x}_t, y, t), \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{x}_0) \rangle \frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)} \right] \\ &= \int_{\mathbf{x}_0} \int_{\mathbf{x}_t} \int_y p(\mathbf{x}_0, \mathbf{x}_t, y) \langle \mathbf{s}_\phi(\mathbf{x}_t, y, t), \frac{\nabla_{\mathbf{x}_t} p(\mathbf{x}_t | \mathbf{x}_0)}{p(\mathbf{x}_t | \mathbf{x}_0)} \rangle \frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)} d\mathbf{x}_0 d\mathbf{x}_t dy \\ &= \int_{\mathbf{x}_0} \int_{\mathbf{x}_t} \int_y p(\mathbf{x}_0, \mathbf{x}_t, y) \langle \mathbf{s}_\phi(\mathbf{x}_t, y, t), \frac{\nabla_{\mathbf{x}_t} q(\mathbf{x}_t | \mathbf{x}_0)}{p(\mathbf{x}_t | \mathbf{x}_0)} \rangle \frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)} d\mathbf{x}_0 d\mathbf{x}_t dy \\ &= \int_{\mathbf{x}_0} \int_{\mathbf{x}_t} \int_y \langle \mathbf{s}_\phi(\mathbf{x}_t, y, t), \nabla_{\mathbf{x}_t} q(\mathbf{x}_t | \mathbf{x}_0) \rangle q(\mathbf{x}_0, y) d\mathbf{x}_0 d\mathbf{x}_t dy \\ &= \int_{\mathbf{x}_0} \int_{\mathbf{x}_t} \int_y \langle \mathbf{s}_\phi(\mathbf{x}_t, y, t), \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t | \mathbf{x}_0) \rangle q(\mathbf{x}_t | \mathbf{x}_0) q(\mathbf{x}_0, y) d\mathbf{x}_0 d\mathbf{x}_t dy \\ &= \mathbb{E}_{q(\mathbf{x}_0, \mathbf{x}_t, y)} [\langle \mathbf{s}_\phi(\mathbf{x}_t, y, t), \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t | \mathbf{x}_0) \rangle]. \end{aligned}$$

□

Lemma A.3. Assume \mathbf{x}_t and y are conditional independent given \mathbf{x}_0 in the forward process, i.e., $p(\mathbf{x}_t | \mathbf{x}_0, y) = p(\mathbf{x}_t | \mathbf{x}_0)$, $\forall t \in [0, T]$, and let the forward process on the target domain be identical to that on the source domain $q(\mathbf{x}_t | \mathbf{x}_0) = p(\mathbf{x}_t | \mathbf{x}_0)$, and ϕ^* is the optimal solution for the conditional diffusion model trained on target domain $q(\mathbf{x}_0, y)$, i.e.,

$$\phi^* = \arg \min_{\phi} \mathbb{E}_t \left\{ \lambda(t) \mathbb{E}_{q_t(\mathbf{x}_t, y)} [\|\mathbf{s}_\phi(\mathbf{x}_t, y, t) - \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t | y)\|_2^2] \right\}, \quad (21)$$

then

$$\mathbf{s}_{\phi^*}(\mathbf{x}_t, y, t) = \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|y) + \nabla_{\mathbf{x}_t} \log \mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)} \right]. \quad (22)$$

Proof of Lemma A.3. According to Lemma A.2, the optimal solution satisfies

$$\phi^* = \arg \min_{\phi} \mathbb{E}_t \left\{ \lambda(t) \mathbb{E}_{p(\mathbf{x}_0, y)} \mathbb{E}_{p(\mathbf{x}_t|\mathbf{x}_0)} \left[\left\| \mathbf{s}_{\phi}(\mathbf{x}_t, y, t) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{x}_0) \right\|_2^2 \frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)} \right] \right\}$$

where $Z(y) = \int p(\mathbf{x}_0, y) \exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) d\mathbf{x}$. Then, we use Importance Weighted Conditional Denoising Score Matching on the source domain to get the analytic form of \mathbf{s}_{ϕ^*} as follows:

$$\mathbf{s}_{\phi^*}(\mathbf{x}_t, y, t) = \frac{\mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{x}_0) \frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)} \right]}{\mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)} \right]}.$$

Moreover, the RHS of (22) can be rewritten as:

$$\begin{aligned} \text{RHS} &= \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|y) + \nabla_{\mathbf{x}_t} \log \mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)} \right] \\ &= \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|y) + \frac{\nabla_{\mathbf{x}_t} \mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)} \right]}{\mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)} \right]} \\ &= \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|y) + \frac{\mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)} \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_0|\mathbf{x}_t, y) \right]}{\mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)} \right]}. \end{aligned}$$

Since

$$\begin{aligned} \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_0|\mathbf{x}_t, y) &= \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{x}_0, y) + \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_0|y) - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|y) \\ &= \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{x}_0, y) - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|y), \\ &= \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{x}_0) - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|y), \end{aligned}$$

we can further simplify the RHS of (22) as follows:

$$\begin{aligned} \text{RHS} &= \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|y) + \frac{\mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)} \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{x}_0) \right]}{\mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)} \right]} - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|y) \\ &= \frac{\mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{x}_0) \frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)} \right]}{\mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)} \right]} \\ &= \mathbf{s}_{\phi^*}(\mathbf{x}_t, t). \end{aligned}$$

Thereby, we finish the proof. \square

According to the lemma A.3, we replace the density ratio $\frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)}$ by $\frac{\exp\left(\frac{1}{\beta}r(\mathbf{x}_0, y)\right)}{Z(y)}$, we get

$$\begin{aligned} \mathbf{s}_{\phi^*}(\mathbf{x}_t, y, t) &= \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t | y) + \nabla_{\mathbf{x}_t} \log \mathbb{E}_{p(\mathbf{x}_0 | \mathbf{x}_t, y)} \left[\frac{q(\mathbf{x}_0, y)}{p(\mathbf{x}_0, y)} \right] \\ &= \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t | y) + \nabla_{\mathbf{x}_t} \log \mathbb{E}_{p(\mathbf{x}_0 | \mathbf{x}_t, y)} \left[\frac{\exp\left(\frac{1}{\beta}r(\mathbf{x}_0, y)\right)}{Z(y)} \right] \\ &= \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t | y) + \nabla_{\mathbf{x}_t} \log \mathbb{E}_{p(\mathbf{x}_0 | \mathbf{x}_t, y)} \left[\exp\left(\frac{1}{\beta}r(\mathbf{x}_0, y)\right) \right] \end{aligned}$$

Thereby, we finish the proof. □

A.2 Proof of Theorem 4.1

We provide a detailed discussion about training-free guidance of flow matching in this subsection.

Proof of Theorem 4.1. Denote $\mathbf{v}_t(\mathbf{x}_t, y)$ and $\mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_1, y)$ as the marginal and conditional velocities,

respectively. Then we have

$$\begin{aligned}
\mathbf{v}_t^q(\mathbf{x}_t, y) &= \mathbb{E}_{\mathbf{x}_1 \sim q_{1|t}(\mathbf{x}_1|\mathbf{x}_t, y)} [\mathbf{v}_t(\mathbf{x}_t \mid \mathbf{x}_1, y)] \\
&= \mathbb{E}_{\mathbf{x}_1 \sim p_{1|t}(\mathbf{x}_1|\mathbf{x}_t, y)} \left[\mathbf{v}_t(\mathbf{x}_t \mid \mathbf{x}_1, y) \frac{q_{1|t}(\mathbf{x}_1 \mid \mathbf{x}_t, y)}{p_{1|t}(\mathbf{x}_1 \mid \mathbf{x}_t, y)} \right] \\
&= \mathbb{E}_{\mathbf{x}_1 \sim p_{1|t}(\mathbf{x}_1|\mathbf{x}_t, y)} \left[\mathbf{v}_t(\mathbf{x}_t \mid \mathbf{x}_1, y) \frac{\frac{q_{t|1}(\mathbf{x}_t|\mathbf{x}_1, y) q_1(\mathbf{x}_1)}{q_t(\mathbf{x}_t, y)}}{\frac{p_{t|1}(\mathbf{x}_t|\mathbf{x}_1, y) p_1(\mathbf{x}_1)}{p_t(\mathbf{x}_t, y)}} \right] \\
&= \mathbb{E}_{\mathbf{x}_1 \sim p_{1|t}(\mathbf{x}_1|\mathbf{x}_t, y)} \left[\mathbf{v}_t(\mathbf{x}_t \mid \mathbf{x}_1, y) \frac{q_{t|1}(\mathbf{x}_t \mid \mathbf{x}_1, y) q_1(\mathbf{x}_1) p_t(\mathbf{x}_t, y)}{p_{t|1}(\mathbf{x}_t \mid \mathbf{x}_1, y) p_1(\mathbf{x}_1) q_t(\mathbf{x}_t, y)} \right] \\
&= \mathbb{E}_{\mathbf{x}_1 \sim p_{1|t}(\mathbf{x}_1|\mathbf{x}_t, y)} \left[\mathbf{v}_t(\mathbf{x}_t \mid \mathbf{x}_1, y) \frac{q_1(\mathbf{x}_1)}{p_1(\mathbf{x}_1)} \cdot \frac{p_t(\mathbf{x}_t, y)}{q_t(\mathbf{x}_t, y)} \right] \quad (\text{because } q_{t|1}(\mathbf{x}_t \mid \mathbf{x}_1, y) = p_{t|1}(\mathbf{x}_t \mid \mathbf{x}_1, y)) \\
&= \mathbb{E}_{\mathbf{x}_1 \sim p_{1|t}(\mathbf{x}_1|\mathbf{x}_t, y)} \left[\mathbf{v}_t(\mathbf{x}_t \mid \mathbf{x}_1, y) \frac{\frac{q_1(\mathbf{x}_1)}{p_1(\mathbf{x}_1)}}{\frac{q_t(\mathbf{x}_t, y)}{p_t(\mathbf{x}_t, y)}} \right] \\
&= \mathbb{E}_{\mathbf{x}_1 \sim p_{1|t}(\mathbf{x}_1|\mathbf{x}_t, y)} \left[\mathbf{v}_t(\mathbf{x}_t \mid \mathbf{x}_1, y) \frac{\frac{q_1(\mathbf{x}_1)}{p_1(\mathbf{x}_1)}}{\sum_{\mathbf{x}'_1} p_{1|t}(\mathbf{x}'_1 \mid \mathbf{x}_t, y) \frac{q_1(\mathbf{x}'_1)}{p_1(\mathbf{x}'_1)}} \right] \\
&= \mathbb{E}_{\mathbf{x}_1 \sim p_{1|t}(\mathbf{x}_1|\mathbf{x}_t, y)} \left[\mathbf{v}_t(\mathbf{x}_t \mid \mathbf{x}_1, y) \frac{\frac{q_1(\mathbf{x}_1)}{p_1(\mathbf{x}_1)}}{\mathbb{E}_{\mathbf{x}'_1 \sim p_{1|t}(\mathbf{x}_1|\mathbf{x}_t, y)} \left[\frac{q_1(\mathbf{x}'_1)}{p_1(\mathbf{x}'_1)} \right]} \right] \\
&= \mathbb{E}_{\mathbf{x}_1 \sim p_{1|t}(\mathbf{x}_1|\mathbf{x}_t, y)} \left[\mathbf{v}_t(\mathbf{x}_t \mid \mathbf{x}_1, y) \frac{\exp\left(\frac{1}{\beta} r(\mathbf{x}_1, y)\right)}{\mathbb{E}_{\mathbf{x}'_1 \sim p_{1|t}(\mathbf{x}_1|\mathbf{x}_t, y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}'_1, y)\right) \right]} \right] \\
&= \mathbb{E}_{\mathbf{x}_1 \sim p_{1|t}(\mathbf{x}_1|\mathbf{x}_t, y)} \left[\mathbf{v}_t(\mathbf{x}_t \mid \mathbf{x}_1, y) \frac{\exp\left(\frac{1}{\beta} r(\mathbf{x}_1, y)\right)}{\mathbb{E}_{\mathbf{x}'_1 \sim p_{1|t}(\mathbf{x}_1|\mathbf{x}_t, y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}'_1, y)\right) \right]} \right] \\
&= \mathbf{v}_t^p(\mathbf{x}_t, y) + \mathbb{E}_{\mathbf{x}_1 \sim p_{1|t}(\mathbf{x}_1|\mathbf{x}_t, y)} \left[\left(\frac{\exp\left(\frac{1}{\beta} r(\mathbf{x}_1, y)\right)}{\mathbb{E}_{\mathbf{x}'_1 \sim p_{1|t}(\mathbf{x}_1|\mathbf{x}_t, y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}'_1, y)\right) \right]} - 1 \right) \mathbf{v}_t(\mathbf{x}_t \mid \mathbf{x}_1, y) \right].
\end{aligned}$$

The above derivation is the training-based guidance for flow matching, where we need to train the first guidance network ψ_1^* satisfies:

$$h_{\psi_1^*}(\mathbf{x}_t, y, t) = \mathbb{E}_{\mathbf{x}_1 \sim p_{1|t}(\mathbf{x}_1|\mathbf{x}_t, y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}_1, y)\right) \right]$$

by minimizing the objective

$$\mathcal{L}_{\text{guidance}}(\psi_1) := \mathbb{E}_{p(\mathbf{x}_1, \mathbf{x}_t, y)} \left[\left\| h_{\psi_1}(\mathbf{x}_t, y, t) - \exp\left(\frac{1}{\beta} r(\mathbf{x}_1, y)\right) \right\|_2^2 \right].$$

And then we need the second guidance network ψ_2^* satisfies:

$$h_{\psi_2^*}(\mathbf{x}_t, y, t) = \mathbb{E}_{\mathbf{x}_1 \sim p_{1|t}(\mathbf{x}_1 | \mathbf{x}_t, y)} \left[\left(\frac{\exp\left(\frac{1}{\beta} r(\mathbf{x}_1, y)\right)}{\mathbb{E}_{\mathbf{x}'_1 \sim p_{1|t}(\mathbf{x}'_1 | \mathbf{x}_t, y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}'_1, y)\right) \right]} - 1 \right) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_1, y) \right]$$

by minimizing the objective

$$\mathcal{L}_{\text{guidance}}(\psi_2) := \mathbb{E}_{p(\mathbf{x}_1, \mathbf{x}_t, y)} \left[\left\| h_{\psi_2}(\mathbf{x}_t, y, t) - \left(\frac{\exp\left(\frac{1}{\beta} r(\mathbf{x}_1, y)\right)}{h_{\psi_1}(\mathbf{x}_t, y, t)} - 1 \right) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_1, y) \right\|_2^2 \right].$$

The guidance network for flow matching is more complex than that used in diffusion models. The estimation errors from two guidance networks may accumulate and ultimately degrade generation performance. To address this limitation, we propose a training-free guidance method for flow matching that mitigates these issues.

$$\mathbf{v}_t^q(\mathbf{x}_t, y)$$

$$\begin{aligned} &= \mathbf{v}_t^p(\mathbf{x}_t, y) + \mathbb{E}_{\mathbf{x}_1 \sim p_{1|t}(\mathbf{x}_1 | \mathbf{x}_t, y)} \left[\left(\frac{\exp\left(\frac{1}{\beta} r(\mathbf{x}_1, y)\right)}{\mathbb{E}_{\mathbf{x}'_1 \sim p_{1|t}(\mathbf{x}'_1 | \mathbf{x}_t, y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}'_1, y)\right) \right]} - 1 \right) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_1, y) \right] \\ &= \mathbf{v}_t^p(\mathbf{x}_t, y) + \int_{\mathbf{x}_1} \left(\frac{\exp\left(\frac{1}{\beta} r(\mathbf{x}_1, y)\right)}{\mathbb{E}_{\mathbf{x}'_1 \sim p_{1|t}} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}'_1, y)\right) \right]} - 1 \right) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_1, y) p_{1|t}(\mathbf{x}_1 | \mathbf{x}_t, y) d\mathbf{x}_1 \\ &= \mathbf{v}_t^p(\mathbf{x}_t, y) + \int_{\mathbf{x}_1} \left(\frac{\exp\left(\frac{1}{\beta} r(\mathbf{x}_1, y)\right)}{\mathbb{E}_{\mathbf{x}_1 \sim p_{1|t}} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}_1, y)\right) \right]} - 1 \right) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_1, y) \frac{p_{t|1}(\mathbf{x}_t | \mathbf{x}_1, y) p(\mathbf{x}_1 | y)}{p_t(\mathbf{x}_t | y)} d\mathbf{x}_1 \\ &= \mathbf{v}_t^p(\mathbf{x}_t, y) + \mathbb{E}_{\mathbf{x}_1 \sim p(\mathbf{x}_1 | y)} \left[\left(\frac{\exp\left(\frac{1}{\beta} r(\mathbf{x}_1, y)\right)}{\mathbb{E}_{\mathbf{x}_1 \sim p_{1|t}} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}_1, y)\right) \right]} - 1 \right) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_1, y) \frac{p_{t|1}(\mathbf{x}_t | \mathbf{x}_1, y)}{p_t(\mathbf{x}_t | y)} \right] \\ &= \mathbf{v}_t^p(\mathbf{x}_t, y) + \mathbb{E}_{\mathbf{x}_1 \sim p(\mathbf{x}_1 | y)} \left[\left(\frac{\exp\left(\frac{1}{\beta} r(\mathbf{x}_1, y)\right)}{\mathbb{E}_{\mathbf{x}_1 \sim p_{1|t}} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}_1, y)\right) \right]} - 1 \right) \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_1, y) \frac{p_{t|1}(\mathbf{x}_t | \mathbf{x}_1, y)}{\mathbb{E}_{\mathbf{x}_1 \sim p(\mathbf{x}_1 | y)} [p_{t|1}(\mathbf{x}_t | \mathbf{x}_1, y)]} \right] \\ &= \mathbf{v}_t^p(\mathbf{x}_t, y) + \mathbb{E}_{\mathbf{x}_1 \sim p(\mathbf{x}_1 | y)} \left[\left(\frac{\exp\left(\frac{1}{\beta} r(\mathbf{x}_1, y)\right)}{\mathbb{E}_{\mathbf{x}_1 \sim p(\mathbf{x}_1 | y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}_1, y)\right) \frac{p_{t|1}(\mathbf{x}_t | \mathbf{x}_1, y)}{\mathbb{E}_{\mathbf{x}_1 \sim p(\mathbf{x}_1 | y)} [p_{t|1}(\mathbf{x}_t | \mathbf{x}_1, y)]} \right]} - 1 \right) \right. \\ &\quad \left. \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_1, y) \frac{p_{t|1}(\mathbf{x}_t | \mathbf{x}_1, y)}{\mathbb{E}_{\mathbf{x}_1 \sim p(\mathbf{x}_1 | y)} [p_{t|1}(\mathbf{x}_t | \mathbf{x}_1, y)]} \right]. \end{aligned}$$

□

A.3 Proof of Lemma 3.2

Proof. The proof is straightforward and we include it below for completeness. Note that the objective function can be rewritten as

$$\begin{aligned}
& \mathcal{L}_{\text{guidance}}(\psi) \\
& := \mathbb{E}_{p(\mathbf{x}_0, \mathbf{x}_t, y)} \left[\left\| h_{\psi}(\mathbf{x}_t, y, t) - \exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) \right\|_2^2 \right] \\
& = \int_{\mathbf{x}_t} \int_y \left\{ \int_{\mathbf{x}_0} p(\mathbf{x}_0 | \mathbf{x}_t, y) \left\| h_{\psi}(\mathbf{x}_t, y, t) - \exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) \right\|_2^2 d\mathbf{x}_0 \right\} p(\mathbf{x}_t | y) p(y) dy d\mathbf{x}_t \\
& = \int_{\mathbf{x}_t} \int_y \left\{ \|h_{\psi}(\mathbf{x}_t, y, t)\|_2^2 - 2 \langle h_{\psi}(\mathbf{x}_t, y, t), \int_{\mathbf{x}_0} p(\mathbf{x}_0 | \mathbf{x}_t, y) \exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) d\mathbf{x}_0 \rangle \right\} p(\mathbf{x}_t | y) p(y) dy d\mathbf{x}_t + C \\
& = \int_{\mathbf{x}_t} \int_y \left\| h_{\psi}(\mathbf{x}_t, y, t) - \mathbb{E}_{p(\mathbf{x}_0 | \mathbf{x}_t, y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) \right] \right\|_2^2 p(\mathbf{x}_t | y) p(y) dy d\mathbf{x}_t,
\end{aligned}$$

where C is a constant independent of ψ . Thus we have the minimizer $\psi^* = \arg \min_{\psi} \mathcal{L}_{\text{guidance}}(\psi)$

satisfies $h_{\psi^*}(\mathbf{x}_t, y, t) = \mathbb{E}_{p(\mathbf{x}_0 | \mathbf{x}_t, y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) \right]$. \square

B Gradient-free Designs for Diffusion Models

After we learn the guidance network by Algorithm 1, we can adopt Eq (13) for inference. Although we can easily calculate the gradient of the guidance network with respect to \mathbf{x}_t by autograd, a question is whether we can avoid the gradient calculation. In this section, we propose two additional designs for gradient-free guidance of diffusion models.

B.1 Training-free guidance for Diffusion Models

The first design is converting the gradient of the log expectation to the expectation under reward weighted distribution, and then we can apply a similar trick as training-free guidance that uses importance sampling to approximate the conditional expectation through Monte Carlo sampling under the marginal data distribution.

Theorem B.1 (Reward-Weighted Score Gradient). *Let $p(\mathbf{x}_0 | \mathbf{x}_t)$ be the reverse diffusion posterior, $r(\mathbf{x}_0, y)$ be a reward function, and $\beta > 0$ be a temperature parameter. Define the reward-weighted distribution*

$$\tilde{p}(\mathbf{x}_0 | \mathbf{x}_t, y) = \frac{p(\mathbf{x}_0 | \mathbf{x}_t, y) \exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right)}{\mathbb{E}_{p(\mathbf{x}_0 | \mathbf{x}_t, y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) \right]}.$$

Then the gradient of the log-partition function satisfies

$$\nabla_{\mathbf{x}_t} \log \mathbb{E}_{p(\mathbf{x}_0 | \mathbf{x}_t, y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) \right] = \mathbb{E}_{\tilde{p}(\mathbf{x}_0 | \mathbf{x}_t, y)} [\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_0 | \mathbf{x}_t, y)].$$

Furthermore, this gradient can be approximated via importance sampling: given samples $\{\mathbf{x}_0^{(k)}\}_{k=1}^K \sim p(\mathbf{x}_0 | y)$ from a proposal distribution,

$$\nabla_{\mathbf{x}_t} \log \mathbb{E}_{p(\mathbf{x}_0 | \mathbf{x}_t, y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) \right] \approx \sum_{k=1}^K \tilde{w}_k \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_0^{(k)} | \mathbf{x}_t, y),$$

where the normalized importance weights are

$$\tilde{w}_k = \frac{u_k}{\sum_{j=1}^K u_j}, \quad u_k = \frac{p(\mathbf{x}_t|\mathbf{x}_0^{(k)}, y) \exp\left(\frac{1}{\beta} r(\mathbf{x}_0^{(k)}, y)\right)}{\mathbb{E}_{p(\mathbf{x}_0|y)} [p(\mathbf{x}_t|\mathbf{x}_0, y)]}.$$

Proof of Theorem B.1. We first convert the gradient of the log expectation to the expectation under reward weighted distribution:

$$\begin{aligned} & \nabla_{\mathbf{x}_t} \log \mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) \right] \\ &= \frac{\nabla_{\mathbf{x}_t} \mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) \right]}{\mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) \right]} && \text{(chain rule)} \\ &= \frac{\nabla_{\mathbf{x}_t} \int p(\mathbf{x}_0|\mathbf{x}_t, y) \exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) d\mathbf{x}_0}{\mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) \right]} && \text{(definition of expectation)} \\ &= \frac{\int \nabla_{\mathbf{x}_t} p(\mathbf{x}_0|\mathbf{x}_t, y) \exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) d\mathbf{x}_0}{\mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) \right]} && \text{(interchange } \nabla \text{ and } \int) \\ &= \frac{\int p(\mathbf{x}_0|\mathbf{x}_t, y) \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_0|\mathbf{x}_t, y) \exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) d\mathbf{x}_0}{\mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) \right]} && \text{(log-derivative trick)} \\ &= \frac{\mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_0|\mathbf{x}_t, y) \exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) \right]}{\mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_t, y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) \right]} && \text{(definition of expectation)} \\ &= \mathbb{E}_{\tilde{p}(\mathbf{x}_0|\mathbf{x}_t, y)} [\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_0|\mathbf{x}_t, y)], && \text{(definition of } \tilde{p}) \end{aligned}$$

where the last equality follows from the fact that the ratio of expectations defines precisely the expectation under the normalized distribution $\tilde{p}(\mathbf{x}_0|\mathbf{x}_t, y)$.

For the importance sampling approximation, we rewrite the expectation using proposal distribution $q(\mathbf{x}_0) = p(\mathbf{x}_0|y)$:

$$\begin{aligned} & \mathbb{E}_{\tilde{p}(\mathbf{x}_0|\mathbf{x}_t, y)} [\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_0|\mathbf{x}_t, y)] \\ &= \frac{\int \tilde{p}(\mathbf{x}_0|\mathbf{x}_t, y) \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_0|\mathbf{x}_t, y) d\mathbf{x}_0}{1} \\ &= \frac{\int \frac{\tilde{p}(\mathbf{x}_0|\mathbf{x}_t, y)}{p(\mathbf{x}_0|y)} p(\mathbf{x}_0|y) \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_0|\mathbf{x}_t, y) d\mathbf{x}_0}{\int \frac{\tilde{p}(\mathbf{x}_0|\mathbf{x}_t, y)}{p(\mathbf{x}_0|y)} p(\mathbf{x}_0|y) d\mathbf{x}_0} \\ &= \frac{\mathbb{E}_{p(\mathbf{x}_0|y)} \left[\frac{\tilde{p}(\mathbf{x}_0|\mathbf{x}_t, y)}{p(\mathbf{x}_0|y)} \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_0|\mathbf{x}_t, y) \right]}{\mathbb{E}_{p(\mathbf{x}_0|y)} \left[\frac{\tilde{p}(\mathbf{x}_0|\mathbf{x}_t, y)}{p(\mathbf{x}_0|y)} \right]} \\ &\approx \frac{\sum_{k=1}^K u_k \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_0^{(k)}|\mathbf{x}_t, y)}{\sum_{k=1}^K u_k} \\ &= \sum_{k=1}^K \tilde{w}_k \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_0^{(k)}|\mathbf{x}_t, y), \end{aligned} \tag{23}$$

where the approximation uses Monte Carlo sampling with $\mathbf{x}_0^{(k)} \sim p(\mathbf{x}_0|y)$ and the unnormalized importance weights u_k are

$$\begin{aligned}
u_k &= \frac{p(\mathbf{x}_0^{(k)}|\mathbf{x}_t, y) \exp\left(\frac{1}{\beta} r(\mathbf{x}_0^{(k)}, y)\right)}{p(\mathbf{x}_0^{(k)}|y)} \\
&= \frac{p(\mathbf{x}_0^{(k)}|\mathbf{x}_t, y)p(\mathbf{x}_t|y) \exp\left(\frac{1}{\beta} r(\mathbf{x}_0^{(k)}, y)\right)}{p(\mathbf{x}_0^{(k)}|y)p(\mathbf{x}_t|y)} && \text{(multiply by } p(\mathbf{x}_t|y)\text{)} \\
&= \frac{p(\mathbf{x}_t|\mathbf{x}_0^{(k)}, y) \exp\left(\frac{1}{\beta} r(\mathbf{x}_0^{(k)}, y)\right)}{p(\mathbf{x}_t|y)} && \text{(Bayes' rule)} \\
&= \frac{p(\mathbf{x}_t|\mathbf{x}_0^{(k)}, y) \exp\left(\frac{1}{\beta} r(\mathbf{x}_0^{(k)}, y)\right)}{\int p(\mathbf{x}_t, \mathbf{x}_0|y) d\mathbf{x}_0} && \text{(marginalization)} \\
&= \frac{p(\mathbf{x}_t|\mathbf{x}_0^{(k)}, y) \exp\left(\frac{1}{\beta} r(\mathbf{x}_0^{(k)}, y)\right)}{\int p(\mathbf{x}_t|\mathbf{x}_0, y)p(\mathbf{x}_0|y) d\mathbf{x}_0} && \text{(chain rule)} \\
&= \frac{p(\mathbf{x}_t|\mathbf{x}_0^{(k)}, y) \exp\left(\frac{1}{\beta} r(\mathbf{x}_0^{(k)}, y)\right)}{\mathbb{E}_{p(\mathbf{x}_0|y)} [p(\mathbf{x}_t|\mathbf{x}_0, y)]}. && \text{(definition of expectation)}
\end{aligned}$$

And the $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_0^{(k)}|\mathbf{x}_t, y)$ can be easily computed through Bayesian rule.

$$\begin{aligned}
\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_0^{(k)}|\mathbf{x}_t, y) &= \nabla_{\mathbf{x}_t} \log \frac{p(\mathbf{x}_t|\mathbf{x}_0^{(k)}, y)p(\mathbf{x}_0^{(k)}|y)}{p(\mathbf{x}_t|y)} \\
&= \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{x}_0^{(k)}, y) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|y)
\end{aligned}$$

□

B.2 Gradient-free Finetuning-free Guidance for Diffusion Models

Another method is directly fitting a neural network to estimate the guidance term.

Theorem B.2. *For a neural network $h_\psi(\mathbf{x}_t, y, t)$ parameterized by ψ , define the objective*

$$\begin{aligned}
\mathcal{L}_{\text{guidance}}^*(\psi) &:= \mathbb{E}_{p(\mathbf{x}_0, \mathbf{x}_T, y)} \left[\frac{1}{\mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_T, y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) \right]} \right. \\
&\quad \left. \frac{1}{p(\mathbf{x}_0|\mathbf{x}_T, y)} \left\| h_\psi(\mathbf{x}_T, y) - \nabla_{\mathbf{x}_T} p(\mathbf{x}_0|\mathbf{x}_T, y) \exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) \right\|_2^2 \right],
\end{aligned}$$

then its minimizer $\psi^ = \arg \min_{\psi} \mathcal{L}_{\text{guidance}}(\psi)$ satisfies:*

$$h_{\psi^*}(\mathbf{x}_T, y) = \nabla_{\mathbf{x}_T} \log \mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_T, y)} \left[\exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) \right].$$

Proof of Theorem B.2.

$$\begin{aligned}
& \int_{\mathbf{x}_T} \left\| h_\psi(\mathbf{x}_T, y) - \nabla_{\mathbf{x}_T} \log \mathbb{E}_{p(\mathbf{x}_0|\mathbf{x}_T)} \left[\exp \left(\frac{1}{\beta} r(\mathbf{x}_0, y) \right) \right] \right\|_2^2 p(\mathbf{x}_T) d\mathbf{x}_T \\
&= \int_{\mathbf{x}_T} \left\{ \|h_\psi(\mathbf{x}_T, y)\|_2^2 - 2 \langle h_\psi(\mathbf{x}_T, y), \nabla_{\mathbf{x}_T} \log \int_{\mathbf{x}_0} p(\mathbf{x}_0|\mathbf{x}_T) \exp \left(\frac{1}{\beta} r(\mathbf{x}_0, y) \right) d\mathbf{x}_0 \rangle \right\} p(\mathbf{x}_T) d\mathbf{x}_T + C \\
&= \int_{\mathbf{x}_T} \left\{ \|h_\psi(\mathbf{x}_T, y)\|_2^2 - 2 \langle h_\psi(\mathbf{x}_T, y), \frac{\nabla_{\mathbf{x}_T} \int_{\mathbf{x}_0} p(\mathbf{x}_0|\mathbf{x}_T) \exp \left(\frac{1}{\beta} r(\mathbf{x}_0, y) \right) d\mathbf{x}_0}{\int_{\mathbf{x}_0} p(\mathbf{x}_0|\mathbf{x}_T) \exp \left(\frac{1}{\beta} r(\mathbf{x}_0, y) \right) d\mathbf{x}_0} \rangle \right\} p(\mathbf{x}_T) d\mathbf{x}_T + C \\
&= \int_{\mathbf{x}_T} \frac{1}{\int_{\mathbf{x}_0} p(\mathbf{x}_0|\mathbf{x}_T) \exp \left(\frac{1}{\beta} r(\mathbf{x}_0, y) \right) d\mathbf{x}_0} \left\{ \|h_\psi(\mathbf{x}_T, y)\|_2^2 - 2 \left\langle h_\psi(\mathbf{x}_T, y), \int_{\mathbf{x}_0} \nabla_{\mathbf{x}_T} p(\mathbf{x}_0|\mathbf{x}_T) \exp \left(\frac{1}{\beta} r(\mathbf{x}_0, y) \right) d\mathbf{x}_0 \right\rangle \right\} \\
&\quad p(\mathbf{x}_T) d\mathbf{x}_T + C \\
&= \int_{\mathbf{x}_T} \left\{ \int_{\mathbf{x}_0} p(\mathbf{x}_0|\mathbf{x}_T) \frac{1}{\int_{\mathbf{x}_0} p(\mathbf{x}_0|\mathbf{x}_T) \exp \left(\frac{1}{\beta} r(\mathbf{x}_0, y) \right) d\mathbf{x}_0} \right. \\
&\quad \left. \frac{1}{p(\mathbf{x}_0|\mathbf{x}_T)} \left\| h_\psi(\mathbf{x}_T, y) - \nabla_{\mathbf{x}_T} p(\mathbf{x}_0|\mathbf{x}_T) \exp \left(\frac{1}{\beta} r(\mathbf{x}_0, y) \right) \right\|_2^2 d\mathbf{x}_0 \right\} p(\mathbf{x}_T) d\mathbf{x}_T + C \\
&= \mathbb{E}_{p(\mathbf{x}_0, \mathbf{x}_T)} \left[\frac{1}{\int_{\mathbf{x}_0} p(\mathbf{x}_0|\mathbf{x}_T) \exp \left(\frac{1}{\beta} r(\mathbf{x}_0, y) \right) d\mathbf{x}_0} \frac{1}{p(\mathbf{x}_0|\mathbf{x}_T)} \left\| h_\psi(\mathbf{x}_T, y) - \nabla_{\mathbf{x}_T} p(\mathbf{x}_0|\mathbf{x}_T) \exp \left(\frac{1}{\beta} r(\mathbf{x}_0, y) \right) \right\|_2^2 \right] + C \\
&:= \mathcal{L}_{\text{guidance}}(\psi) + C,
\end{aligned}$$

where $\nabla_{\mathbf{x}_T} p(\mathbf{x}_0|\mathbf{x}_T)$ can be computed by $\nabla_{\mathbf{x}_T} \log p(\mathbf{x}_0|\mathbf{x}_T) \frac{1}{p(\mathbf{x}_0|\mathbf{x}_T)}$. \square

C More Details on Experiments

C.1 Algorithms for Training the Guidance Network

Algorithm 1 is the algorithm for training the guidance network.

C.2 Ablation Study on Hyperparameter

In this subsection, we provide the ablation study of the strength of the regularization η and the strength of the reward function β in the following table.

Table 4: Ablation study of hyperparameter on PickScore.

η	$\beta = 10$	$\beta = 15$	$\beta = 20$
0.1	22.82	22.79	22.72
0.5	22.78	23.01	22.79
1	22.76	23.08	22.84

C.3 Prompts for Figure in Main Paper

Algorithm 1 Algorithm for Training a Guidance Network

Require: Samples from alignment dataset, pre-trained one-step diffusion model $s(\mathbf{x}_T, y, T)$, pre-determined reward function $r(\mathbf{x}_0, y)$, hyperparameters η, β , and initial weights of guidance network ψ .

1: **repeat**

2: Sample mini-batch data from alignment dataset with batch size b .

3: Perturb \mathbf{x}_0 using forward transition $p(\mathbf{x}_T|\mathbf{x}_0)$.

4: Compute guidance loss:

5:

$$\mathcal{L}_{\text{guidance}}(\psi) = \frac{1}{b} \sum_{\mathbf{x}_0, \mathbf{x}_T, y} \left\| h_{\psi}(\mathbf{x}_T, y) - \exp\left(\frac{1}{\beta} r(\mathbf{x}_0, y)\right) \right\|_2^2.$$

6: Sample mini-batch from winning responses (\mathbf{x}', y) with batch size b .

7: Perturb \mathbf{x}'_0 using forward transition $q(\mathbf{x}'_T|\mathbf{x}'_0)$.

8: Compute consistency loss:

9:

$$\mathcal{L}_{\text{consistence}} = \frac{1}{b} \sum_{\mathbf{x}'_0, \mathbf{x}'_T, y} \left\| s(\mathbf{x}'_T, y, T) + \nabla_{\mathbf{x}'_T} \log h_{\psi}(\mathbf{x}'_T, y) - \nabla_{\mathbf{x}'_T} \log q(\mathbf{x}'_T|\mathbf{x}'_0, y) \right\|_2^2.$$

10: Update ψ via gradient descent:

$$\nabla_{\psi} (\mathcal{L}_{\text{guidance}} + \eta \mathcal{L}_{\text{consistence}}).$$

11: **until** convergence

12: **return** weights of guidance network ψ .

Table 5: Prompts used to generate Figure 1.

Image	Prompt
Col1	Saturn rises on the horizon.
Col2	a watercolor painting of a super cute kitten wearing a hat of flowers
Col3	A galaxy-colored figurine floating over the sea at sunset, photorealistic.
Col4	fireclaw machine mecha animal beast robot of horizon forbidden west horizon zero dawn bioluminescence, behance hd by jesper ejasing, by rhads, makoto shinkai and lois van baarle, ilya kuvshinov, rossdraws global illumination
Col5	A swirling, multicolored portal emerges from the depths of an ocean of coffee, with waves of the rich liquid gently rippling outward. The portal engulfs a coffee cup, which serves as a gateway to a fantastical dimension. The surrounding digital art landscape reflects the colors of the portal, creating an alluring scene of endless possibilities.
Col6	A profile picture of an anime boy, half robot, brown hair
Col7	Detailed Portrait of a cute woman vibrant pixie hair by Yanjun Cheng and Hsiao-Ron Cheng and Ilya Kuvshinov, medium close up, portrait photography, rim lighting, realistic eyes, photorealism pastel, illustration
Col8	On the Mid-Autumn Festival, the bright full moon hangs in the night sky. A quaint pavilion is illuminated by dim lights, resembling a beautiful scenery in a painting. Camera type: close-up. Camera lens type: telephoto. Time of day: night. Style of lighting: bright. Film type: ancient style. HD.