
Inference-Time Alignment Control for Diffusion Models with Reinforcement Learning Guidance

Anonymous Authors¹

Abstract

Pluralistic alignment calls for generative models that can accommodate diverse—and often conflicting—preferences rather than commit to a single fixed objective. Although denoising-based generative models, particularly diffusion and flow matching algorithms, have achieved remarkable success, existing reinforcement learning (RL) fine-tuning methods typically lock them into a single post-training operating point. In this work, we view RL fine-tuning for diffusion models through the lens of stochastic differential equations and implicit reward conditioning. We introduce *Reinforcement Learning Guidance* (RLG), an inference-time method that adapts Classifier-Free Guidance (CFG) by combining the outputs of the base and RL fine-tuned models via a geometric average. Our theoretical analysis shows that RLG’s guidance scale is mathematically equivalent to adjusting the KL-regularization coefficient in standard RL objectives, enabling dynamic control over the alignment-quality trade-off without further training. Extensive experiments demonstrate that RLG consistently improves the performance of RL fine-tuned models across various architectures, RL algorithms, and downstream tasks, including human preference alignment, compositional control, compressibility, text rendering, and multi-objective blending. Overall, RLG provides a simple inference-time mechanism for pluralistic alignment in diffusion models. The source code for RLG is available in the anonymous repository.

¹.

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

¹<https://anonymous.4open.science/r/Reinforcement-learning-guidance-7B5A/>

1. Introduction

Aligning AI systems with diverse, and sometimes conflicting, human preferences is a central challenge in pluralistic alignment. For denoising-based generative models—primarily diffusion (Ho et al., 2020; Rombach et al., 2022) and flow matching (Lipman et al., 2022; Esser et al., 2024) algorithms—this challenge appears in the need to satisfy heterogeneous downstream objectives such as human preferences (Kirstain et al., 2023), compositional correctness (Ghosh et al., 2023), text rendering (Liu et al., 2025b), or data compressibility (Black et al., 2023). Existing approaches include reward-weighted regression (Peng et al.; Lee et al., 2023; Fan et al., 2025), direct reward fine-tuning (Xu et al., 2023; Prabhudesai et al., 2023; Clark et al., 2023), and reinforcement learning (RL) fine-tuning.

Owing to significant advancements in Reinforcement Learning from Human Feedback (RLHF) (Black et al., 2023; Lee et al., 2023) for Large Language Models (LLMs), RL has been adapted to diffusion models by formulating denoising as a multi-step decision-making process, enabling algorithms like REINFORCE (Williams, 1992; Mohamed et al., 2020; Black et al., 2023), Direct Preference Optimization (DPO) (Rafailov et al., 2023; Wallace et al., 2024), and Group Relative Policy Optimization (GRPO) (Shao et al., 2024; Liu et al., 2025b). However, current RL methods for diffusion models still exhibit several limitations, primarily in two respects. First, the exact probability of a sampled image is intractable due to the nature of diffusion algorithms, which undermines the effectiveness of existing RL algorithms (Black et al., 2023; Gong et al., 2025). Second, the degree to which the base model aligns with downstream objectives remains fixed after RL fine-tuning and is sensitive to hyperparameter choices, such as the Kullback–Leibler (KL) coefficient. This rigidity is especially problematic in pluralistic settings, where different users, communities, or deployment contexts may prefer different trade-offs between reward, fidelity, and competing objectives, yet standard RL exposes only a single post-training alignment strength.

In this work, we draw inspiration from the stochastic differential equation (SDE) nature of denoising-based generative models (Song et al., 2020b), which motivates us to interpret RL fine-tuning of diffusion models as a special case of



Figure 1. Selected qualitative results for the human preference alignment task using SD3.5-M with GRPO and our RLG. The PickScore is displayed on each image. As the RLG scale increases, the images generally become more detailed, aesthetically pleasing, which is corroborated by the rising PickScores.

generation conditioned on implicit rewards learned through reinforcement learning objectives (Rafailov et al., 2024; Zhu et al., 2025; Cui et al., 2025). Building upon this perspective, we introduce an inference-time enhancement technique, *Reinforcement Learning Guidance* (RLG), which adapts the established controlling approach, Classifier-Free Guidance (CFG) (Ho & Salimans, 2022; Zheng et al., 2023), by computing a weighted geometric average of the outputs from the base model and the RL fine-tuned model. We theoretically demonstrate that this weighted averaging has the same effect as modifying the KL coefficient in RL fine-tuning, but crucially, it requires no additional training. As a result, RLG is well suited to pluralistic alignment settings, where one may wish to move among multiple acceptable alignment regimes without retraining separate models for each preference profile.

Empirical results on downstream tasks demonstrate that RLG enhances the performance of RL fine-tuned models across diverse tasks and setups: various model types (diffusion and flow matching), a range of RL methods (policy gradient, DPO, GRPO, etc.), and multiple downstream objectives (image aesthetics, compositional control, compressibility, text rendering, inpainting, and personalized genera-

tion). Notably, our multi-objective experiments show that RLG can combine distinct aligned experts at inference time, highlighting its relevance to settings that require explicit trade-offs among competing values. Furthermore, RLG supports both interpolation and extrapolation, thereby offering substantial flexibility in controlling the degree of alignment with downstream objectives.

Our contributions are summarized as follows:

- We propose *Reinforcement Learning Guidance* (RLG), a novel, training-free approach for enhancing and controlling the inference-time alignment of denoising-based generative models.
- We provide a theoretical foundation for RLG, demonstrating that its guidance scale is mathematically equivalent to adjusting the KL-regularization coefficient in the underlying RL objective. This analysis formally accounts for the effectiveness of extrapolation ($w > 1$).
- We perform extensive experiments on a diverse set of alignment tasks, showing that RLG consistently enhances performance by enabling models to surpass

their original fine-tuned capabilities while also supporting flexible trade-offs between competing objectives.

2. Background

2.1. Diffusion and Flow-based Generative Models

Diffusion (Ho et al., 2020; Song et al., 2020a; Rombach et al., 2022) and flow-based (Lipman et al., 2022; Liu et al., 2022; Esser et al., 2024) models generate data by transforming noise into samples via either stochastic (SDE) or deterministic (ODE) processes. Diffusion models, e.g., DDPM (Ho et al., 2020), DDIM (Song et al., 2020a), and Stable Diffusion (Rombach et al., 2022), corrupt data with noise and train a network to learn the score function $\mathbf{s}(\mathbf{X}_t, t)$ for reverse denoising (Song et al., 2020b). Flow-based approaches such as Flow Matching (Lipman et al., 2022) learn a velocity field $\mathbf{v}(\mathbf{X}_t, t)$ to follow a deterministic path from prior to data (Liu et al., 2022; Esser et al., 2024; Tong et al., 2023; Kong et al., 2024; Liu et al., 2025a; Wan et al., 2025).

A reference flow $(X_t)_{t \in [0,1]}$ interpolates between $X_1 \sim p_1$ and $X_0 \sim p_{\text{data}}$:

$$X_t = \beta_t X_1 + \alpha_t X_0, \quad \alpha_0 = \beta_1 = 0, \quad \alpha_1 = \beta_0 = 1. \quad (1)$$

For ODE-based Flow Matching, the model is trained to match the reference velocity:

$$\mathbf{v}(\mathbf{X}_t, t) \approx \dot{\beta}_t X_0 + \dot{\alpha}_t X_1.$$

Diffusion models solve an SDE with noise schedule $\sigma(t)$, typically parameterized by $\alpha_t = \sqrt{\bar{\alpha}_t}$, $\beta_t = \sqrt{1 - \bar{\alpha}_t}$.

Formally, the velocity field in Flow Matching can be written in terms of the score function as

$$\mathbf{v}(\mathbf{x}, t) = \left(\frac{\dot{\alpha}_t}{\alpha_t} \right) \mathbf{x} + \beta_t \left(\frac{\dot{\alpha}_t}{\alpha_t} \beta_t - \dot{\beta}_t \right) \mathbf{s}(\mathbf{x}, t). \quad (2)$$

The two paradigms unify under the SDE (Song et al., 2020b):

$$d\mathbf{X}_t = \left(\mathbf{v}(\mathbf{X}_t, t) - \frac{1}{2} \sigma(t)^2 \mathbf{s}(\mathbf{X}_t, t) \right) dt + \sigma(t) dw, \quad (3)$$

where w is Brownian motion; diffusion and Flow Matching differ in \mathbf{v} , \mathbf{s} , and $\sigma(t)$.

2.2. Guidance and Control in Generative Models

Controlling generative model outputs is essential for conditional generation tasks. Early work such as Classifier Guidance (CG) steers the generation process using gradients from a separately trained classifier (Dhariwal & Nichol, 2021), but this approach is computationally costly and limited by the need for an external model.

Classifier-Free Guidance (CFG) has become the standard alternative (Ho & Salimans, 2022). At inference, CFG computes two passes: one with the actual condition c (e.g., text-guided) and one with the null condition \emptyset . The guided velocity field $\hat{\mathbf{v}}_\theta$ is a linear interpolation between these two outputs:

$$\hat{\mathbf{v}}_\theta(\mathbf{x}_t, t|c) \triangleq (1 - \omega) \mathbf{v}_\theta(\mathbf{x}_t, t|\emptyset) + \omega \mathbf{v}_\theta(\mathbf{x}_t, t|c), \quad (4)$$

where ω is the guidance scale parameter. Setting $\omega = 1$ recovers conditional generation, while $\omega > 1$ extrapolates beyond the conditional prediction.

The same principle applies to the model’s score function:

$$\hat{\mathbf{s}}_\theta(\mathbf{x}_t, t|c) \triangleq (1 - \omega) \mathbf{s}_\theta(\mathbf{x}_t, t|\emptyset) + \omega \mathbf{s}_\theta(\mathbf{x}_t, t|c), \quad (5)$$

where $\mathbf{s}_\theta(\mathbf{x}_t, t|c) = \nabla_{\mathbf{x}_t} \log p_\theta(\mathbf{x}_t|c)$. This can be equivalently written as:

$$\hat{\mathbf{s}}_\theta(\mathbf{x}_t, t|c) = \nabla_{\mathbf{x}_t} \log (p_\theta(\mathbf{x}_t)^{1-\omega} p_\theta(\mathbf{x}_t|c)^\omega), \quad (6)$$

Although CFG is highly effective, most existing methods focus on adherence to training-time conditions such as text condition, leaving open the possibility of leveraging reinforcement learning rewards as dynamic, flexible forms of guidance.

2.3. Preference Alignment in Generative Models

Preference learning methods from LLMs have been adapted to fine-tune T2I diffusion models for human alignment. A pre-trained model π_{ref} is fine-tuned to maximize a reward $R(\mathbf{x})$ under KL regularization:

$$\pi_\theta^* = \arg \max_{\pi_\theta} \mathbb{E}_{\mathbf{x} \sim \pi_\theta} [R(\mathbf{x})] - \beta D_{\text{KL}}(\pi_\theta \| \pi_{\text{ref}}), \quad (7)$$

where β controls the reward–regularization trade-off.

The optimal solution to this problem shows the aligned policy is a re-weighted reference policy, with weights exponentially proportional to the reward (Peng et al.; Lee et al., 2023; Fan et al., 2025):

$$p^*(\mathbf{x}) \propto p_{\text{ref}}(\mathbf{x}) \exp \left(\frac{1}{\beta} R(\mathbf{x}) \right). \quad (8)$$

This objective can be optimized with policy gradient methods include PPO (Schulman et al., 2017; Black et al., 2023), as well as direct approaches such as DPO (Rafailov et al., 2023; Wallace et al., 2024) and GRPO (Sun et al., 2025; Shao et al., 2024). However, many alignment methods may overlook characteristics of diffusion models. For instance, Diffusion-DPO (Wallace et al., 2024) is upper-bounded by the original DPO loss. In particular, integrating reinforcement learning objectives with diffusion-specific techniques—such as Classifier-Free Guidance (CFG) (Ho &

Salimans, 2022)—remains underexplored, presenting opportunities to design approaches that combine reward-based alignment with the generative priors and guidance capabilities unique to diffusion.

Concurrent Work. While finalizing our paper, two concurrent works, CFGRL (Frans et al., 2025) and Diffusion Blend (Cheng et al., 2025), appeared. Both investigate inference-time manipulation techniques via score interpolation. However, CFGRL focuses solely on offline RL and simple task settings, while Diffusion Blend does not establish a connection between interpolation and implicit reward guidance. In contrast, RLG offers a comprehensive analysis of score interpolation from the perspective of implicit classifier guidance and demonstrates its effectiveness across various image generation models, RL algorithms, and tasks.

3. Methods

Deriving Reinforcement Learning Guidance (RLG)

Let r represent the desired attribute, such as a high preference score. Following Bayes’ rule, the score function of the conditional distribution $p_{\text{ref}}(\mathbf{x}_t|r)$ can be decomposed as:

$$\nabla_{\mathbf{x}_t} \log p_{\text{ref}}(\mathbf{x}_t|r) = \nabla_{\mathbf{x}_t} \log p_{\text{ref}}(\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p(r|\mathbf{x}_t). \quad (9)$$

To relate this to a reward function $R(\mathbf{x}_t)$, following (Zhu et al., 2025), we model $p(r|\mathbf{x}_t)$ via an energy-based form:

$$p(r|\mathbf{x}_t) = \frac{\exp(R(\mathbf{x}_t))}{Z}, \quad Z = \int \exp(R(\mathbf{x}_t)) d\mathbf{x}_t. \quad (10)$$

Substituting this into Equation 9 yields the general formula for reward gradient guidance:

$$\hat{\mathbf{s}}(\mathbf{x}_t, t) = \mathbf{s}_{\text{ref}}(\mathbf{x}_t, t) + \eta \nabla_{\mathbf{x}_t} R(\mathbf{x}_t). \quad (11)$$

Here, η is a guidance scale. Since we lack an explicit, differentiable reward model $R(\mathbf{x}_t)$, we draw from the solution to the KL-regularized RL objective (Equation 8) and from (Rafailov et al., 2024; Zhu et al., 2025) to define an implicit, time-dependent reward function $R_t(\mathbf{x}_t)$ that represents the preference learned by π_θ throughout the generative process:

$$R_t(\mathbf{x}_t) \triangleq \beta \log \frac{p_{\theta,t}(\mathbf{x}_t)}{p_{\text{ref},t}(\mathbf{x}_t)}. \quad (12)$$

$p_{\theta,t}$ and $p_{\text{ref},t}$ are the marginal probability distributions of the noisy sample \mathbf{x}_t under the RL-aligned and reference models, respectively, and β is the KL-coefficient from the original RL fine-tuning objective.

To use this implicit reward for guidance, we take its gradient

Algorithm 1 Sampling with Reinforcement Learning Guidance (RLG)

- 1: **Input:** Pre-trained model velocity \mathbf{v}_{ref} , RL-finetuned model velocity \mathbf{v}_{RL} , condition c , RLG scale w , number of steps N .
- 2: Sample initial noise $\mathbf{x}_1 \sim \mathcal{N}(0, \mathbf{I})$.
- 3: **for** $t = 1, \dots, N$ **do**
- 4: Compute reference velocity: $\mathbf{v}_{\text{ref},t} = \mathbf{v}_{\text{ref}}(\mathbf{x}_t, t|c)$.
- 5: Compute RL-aligned velocity: $\mathbf{v}_{\text{RL},t} = \mathbf{v}_{\text{RL}}(\mathbf{x}_t, t|c)$.
- 6: Compute the guided velocity using RLG:
- 7: $\hat{\mathbf{v}}_{\text{RLG},t} = (1 - w)\mathbf{v}_{\text{ref},t} + w\mathbf{v}_{\text{RL},t}$.
- 8: Update the sample using a chosen ODE solver step:
- 9: $\mathbf{x}_{t+1} = \text{SolverStep}(\mathbf{x}_t, \hat{\mathbf{v}}_{\text{RLG},t})$.
- 10: **end for**
- 11: **Return:** Generated sample \mathbf{x}_{N+1} .

with respect to \mathbf{x}_t , yielding a simple result:

$$\begin{aligned} \nabla_{\mathbf{x}_t} R_t(\mathbf{x}_t) &= \beta [\nabla_{\mathbf{x}_t} \log p_{\theta,t} - \nabla_{\mathbf{x}_t} \log p_{\text{ref},t}] \\ &= \beta [\mathbf{s}_\theta(\mathbf{x}_t, t) - \mathbf{s}_{\text{ref}}(\mathbf{x}_t, t)]. \end{aligned} \quad (13)$$

Substituting this into Equation 9 yields the general formula for reward gradient guidance:

$$\hat{\mathbf{s}}_{\text{RLG}}(\mathbf{x}_t, t) = (1 - w)\mathbf{s}_{\text{ref}} + w\mathbf{s}_\theta, \quad (14)$$

which is a linear interpolation of score functions, interpretable as implicit reward gradient guidance. Using Eq. 2, the same applies to velocity fields, yielding:

$$\hat{\mathbf{v}}_{\text{RLG}}(\mathbf{x}_t, t) = (1 - w)\mathbf{v}_{\text{ref}}(\mathbf{x}_t, t) + w\mathbf{v}_\theta(\mathbf{x}_t, t), \quad (15)$$

where w is the RLG guidance scale. A value of $w = 0$ recovers the original model, $w = 1$ recovers the RL-finetuned model, and $w > 1$ extrapolates the learned alignment. The full sampling procedure is outlined in Algorithm 1.

Theoretical Justification: RLG as KL-Coefficient Control

RLG’s mechanisms can be explained by a complementary theoretical justification. Similar to CFG (Ho & Salimans, 2022), the guided score $\hat{\mathbf{s}}_{\text{RLG}}$ corresponds to sampling from a new time-dependent distribution:

$$\hat{\mathbf{s}}_{\text{RLG}} = \nabla_{\mathbf{x}_t} \log (p_{\text{ref},t}(\mathbf{x}_t)^{1-w} p_{\theta,t}(\mathbf{x}_t)^w). \quad (16)$$

As $t \rightarrow 0$, the noisy sample \mathbf{x}_t approaches the clean data \mathbf{x}_0 . In this limit, the marginal distributions $p_{\text{ref},t}$ and $p_{\theta,t}$ converge to their corresponding final distributions, $p_{\text{ref}}(\mathbf{x}_0)$ and $p_\theta(\mathbf{x}_0)$. Therefore, the score function guiding the final steps of generation points towards a target distribution $\hat{p}_{\text{RLG}}(\mathbf{x}_0)$ of the form: $p_{\text{ref}}(\mathbf{x}_0)^{1-w} p_\theta(\mathbf{x}_0)^w$.

Assuming the RL-finetuned model π_θ has converged to the optimal distribution from (Rafailov et al., 2024; 2023) (i.e., $p_\theta(\mathbf{x}_0) \propto p_{\text{ref}}(\mathbf{x}_0) \exp(\frac{1}{\beta} R(\mathbf{x}_0))$), we can substitute this into the expression for the RLG distribution:

$$\begin{aligned} \hat{p}_{\text{RLG}}(\mathbf{x}_0) &\propto p_{\text{ref}}(\mathbf{x}_0)^{1-w} \left(p_{\text{ref}}(\mathbf{x}_0) \exp\left(\frac{1}{\beta} R(\mathbf{x}_0)\right) \right)^w \\ &\propto p_{\text{ref}}(\mathbf{x}_0) \exp\left(\frac{1}{\beta/w} R(\mathbf{x}_0)\right). \end{aligned} \quad (17)$$

This derivation reveals a crucial insight: RLG with guidance scale w is mathematically equivalent to sampling from the optimal policy of an RL objective with an effective KL-regularization coefficient of β/w .

We empirically validated this result with a small-scale demonstration. Our experimental setting uses a flow matching model defined on the real line, with a pretrained target (base) Gaussian mixture distribution: $p_{\text{base}}(x) \triangleq 0.7\mathcal{N}(-2.5, 0.25) + 0.3\mathcal{N}(2.5, 0.49)$. The reward function is $r(x) = 0.1x$. We fine-tuned the pretrained model using a policy gradient algorithm with a KL coefficient $\beta = 0.3$, a batch size of 64, and a learning rate of 1×10^{-5} . Figure 2 presents sampled distributions under various RLG weights w , alongside the corresponding theoretical distribution curves $p_{\text{rl}}(x) \propto p_{\text{base}}(x) \exp(\frac{1}{\beta} r(x))$. Results show that RLG-sampled distributions closely match theoretically predicted RL targets, corroborating our analysis.

Selecting $w > 1$ dynamically reduces the regularization penalty at inference time, allowing the model to pursue higher rewards more aggressively than the original RL-finetuned model. Conversely, $w < 1$ increases regularization. This provides principled justification for RLG’s capacity to extrapolate or interpolate beyond the learned policy, offering a powerful and theoretically grounded mechanism to control the trade-off between alignment and fidelity.

4. Experiments

This section empirically validates RLG’s effectiveness. Experiments are conducted on various text-to-image (T2I) alignment tasks. In each case, the original pre-trained model serves as \mathbf{v}_{ref} , and the RL-aligned model as \mathbf{v}_θ . RLG is applied as described in Equation 15. For all experiments, sampling steps were set to 20. A comprehensive list of additional hyperparameters can be found in the appendix H.1.

4.1. RLG Universally Enhances Alignment Across Diverse Tasks

A key strength of RLG is its broad applicability. This training-free method consistently enhances model capabilities across diverse alignment tasks, from high-level compo-

sitional understanding to fine-grained subject fidelity.

Structured Generation: Compositionality and Text Rendering. RLG is first evaluated on tasks requiring precise adherence to structured prompts. For compositional generation, the GenEval benchmark (Ghosh et al., 2023) is used with a GRPO-finetuned SD3.5-M model (Liu et al., 2025b). This task tests the model’s ability to correctly render object relationships, counts, and attributes. Following the benchmark’s protocol, a Mask2Former model (Cheng et al., 2022) verifies the presence and properties of objects specified in prompts within the official GenEval test set.

For visual text rendering, an GRPO-finetuned SD3.5-M model on the Optical Character Recognition (OCR) task (Mori et al., 1999) (also from (Liu et al., 2025b)) measures its ability to accurately render text from prompts that contain the exact string to appear in the image. OCR accuracy is calculated based on the normalized edit distance $(1 - d_{\text{norm}})$, where d_{norm} is the Levenshtein distance (Yujian & Bo, 2007) between generated text (extracted using PaddleOCR (Authors, 2020)) and the ground-truth text, normalized by ground-truth text length.

Tables 2, 3 and Figures 13, 10 show that while RL-finetuned models ($w_{\text{RL}} = 1.0$) already substantially improve over their base counterparts, extrapolating with RLG unlocks further significant gains.

Fidelity-Driven Generation: Inpainting and Personalization. We next test RLG on tasks demanding high fidelity to reference content: image inpainting and personalized generation. For image inpainting, we use the PrefPaint (Bui et al., 2025) model, an RL-finetuned model built on stable-diffusion-inpainting (Podell et al., 2023) and designed to fill masked regions according to human preferences. To evaluate quality, we use Preference Reward metrics (Bui et al., 2025) on the dataset detailed in appendix K. For personalized generation, we use PatchDPO (Huang et al., 2025), an RL-finetuned model optimized to maintain subject identity from reference images. Here, the original pre-trained model (IP-Adapter-Plus (Ye et al., 2023)) serves as the base (\mathbf{v}_{ref}), and PatchDPO as the RL-aligned model (\mathbf{v}_θ). Subject fidelity is measured using two standard image-similarity metrics: CLIP-I (Ruiz et al., 2023) and DINO (Caron et al., 2021), evaluated on the DreamBench (Ruiz et al., 2023) benchmark, detailed in appendix L.

The results, summarized in Table 4, again show RLG’s effectiveness. In both cases, RLG provides measurable enhancement over state-of-the-art RL-finetuned models without any additional training.

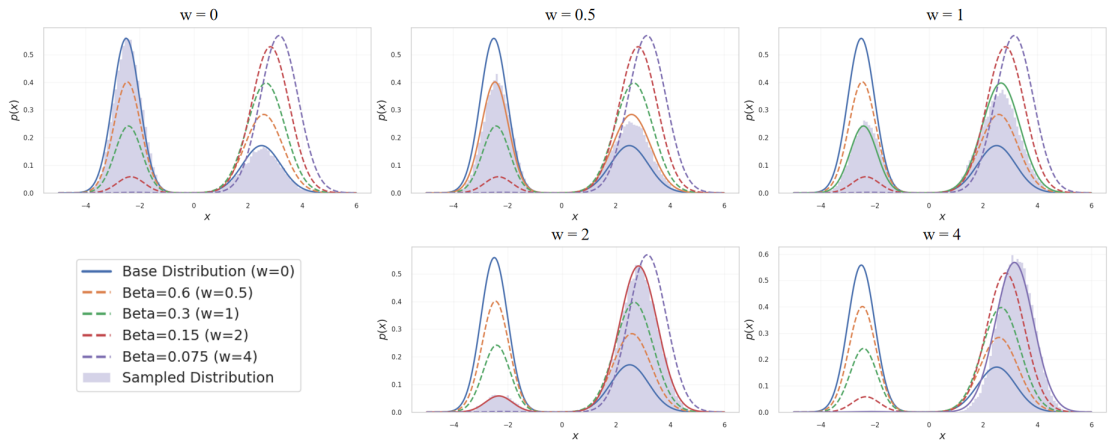


Figure 2. Small-scale demonstration supporting the theoretical justification of RLG. Each subplot shows the sampled distribution under a different RLG weight w , while the curves represent the corresponding theoretically predicted RL-fine-tuned distributions. Here, β denotes the KL regularization coefficient.

Table 1. Quantitative results for Human Preference Alignment. Mean scores for Aesthetic Score, ImageReward, and PickScore are reported. For each metric, values after the slash (/) indicate win rates (%) against the standard RL-finetuned model ($w_{RL} = 1.0$). SD denotes stable-diffusion models (Rombach et al., 2022).

Model	RL	w_{RL}	Aesthetic Score (\uparrow)	ImageReward (\uparrow)	PickScore (\uparrow)
SD1.5	DPO (Lee et al., 2023)	0.0	5.51 / 38.48%	-0.02 / 36.67%	20.03 / 27.34%
		1.0	5.61	0.20	20.39
		1.4	5.63 / 53.56%	0.25 / 53.47%	20.46 / 57.86%
		2.4	5.64 / 56.25%	0.32 / 57.08%	20.56 / 61.23%
SDXL-base	SPO (Liang et al., 2024)	0.0	6.10 / 17.87%	0.72 / 18.02%	21.66 / 7.62%
		1.0	6.42	1.12	22.69
		1.2	6.45 / 59.81%	1.13 / 54.10%	22.71 / 54.64%
		1.4	6.48 / 62.99%	1.14 / 54.15%	22.71 / 54.35%
SD3.5-M	GRPO (Liu et al., 2025b)	0.0	5.97 / 11.33%	0.99 / 17.29%	21.75 / 2.39%
		1.0	6.45	1.40	23.29
		1.4	6.54 / 69.28%	1.40 / 54.44%	23.48 / 74.95%
		2.2	6.64 / 77.39%	1.39 / 53.56%	23.58 / 73.68%

4.2. RLG Naturally Extends to Multi-Objective Blending

We further test whether RLG can combine multiple RL-aligned experts at inference time. Specifically, we fuse an aesthetics-aligned expert and an OCR-aligned expert on top of the same base model, using non-negative weights w_1 and w_2 , and evaluate on a held-out subset of the OCR test prompts.

Table 5 shows that multi-expert RLG yields favorable trade-offs without retraining. Fused configurations recover much of both strengths. For example, the setting $(w_1, w_2) = (1.0, 1.0)$ substantially improves Aesthetic and PickScore over OCR-only while maintaining high OCR performance.

We also study an extrapolation regime by fixing $w_2 = 1.0$ and increasing the aesthetics weight w_1 . As shown in Table 6, Aesthetic and PickScore improve monotonically as w_1 increases, while OCR remains high and stays far above

the aesthetics-only baseline. This shows that multi-expert RLG can strengthen one objective without collapsing the other.

Overall, these results show that RLG naturally extends to multi-objective blending: multiple RL-aligned experts can be combined at inference time to obtain controllable trade-offs between competing objectives.

4.3. RLG is Effective Across Diverse RL Algorithms and Model Architectures

To demonstrate RLG’s broad applicability and model-agnostic nature, we evaluate its consistent enhancement of models differing in generative architecture (i.e., standard diffusion vs. modern flow matching) and the specific reinforcement learning algorithm used for their initial alignment. This also extends to algorithms such as GRPO, whose optimal policy does not necessarily conform to Equation 8.

Table 2. Performance on the GenEval benchmark. We report accuracy (%) for each compositional sub-task and the overall average score across different RLG guidance scales (w_{RL}).

Model	w_{RL}	Single Obj.	Two Objs.	Colors	Color Attr.	Counting	Position	Overall Score
SD3.5-M	0.0	97.81	80.56	80.05	51.75	53.44	23.00	64.44
	1.0	100.00	98.99	89.63	84.00	92.81	93.75	93.20
	1.2	99.69	98.74	90.69	86.11	92.81	94.25	93.72
	1.4	99.69	99.24	91.49	86.50	94.69	94.50	94.35
	1.6	100.00	98.99	91.76	86.00	93.75	95.00	94.25

Table 3. Quantitative results for the visual text rendering task. This table shows the Optical Character Recognition (OCR) accuracy at different RLG guidance scales.

Model	w_{RL}	OCR Acc (\uparrow)	Aesthetic Score (\uparrow)
SD3.5-M	0.0	0.543	5.40
	0.4	0.785	5.28
	0.6	0.838	5.25
	1.0	0.886	5.20
	1.2	0.894	5.17
	1.6	0.910	5.13
	2.2	0.921	5.07
	2.8	0.930	5.00

Table 4. RLG enhances performance on distinct fidelity-driven tasks. The evaluation metrics are presented separately for each task.

Task: Image Inpainting	
Method	Pref. Reward (\uparrow)
Base ($w_{RL} = 0$)	0.080
PrefPaint ($w_{RL} = 1.0$)	0.358
RLG ($w_{RL} = 1.2$)	0.367
RLG ($w_{RL} = 1.4$)	0.368
RLG ($w_{RL} = 1.6$)	0.366

Task: Personalized Generation

Method	DINO (\uparrow)	CLIP-I (\uparrow)
IP-Adapter-Plus ($w_{RL} = 0$)	0.692	0.826
PatchDPO ($w_{RL} = 1.0$)	0.724	0.839
RLG ($w = 1.2$)	0.730	0.841
RLG ($w = 1.8$)	0.730	0.843

Experimental Setup. We analyze RLG on the human preference alignment task, leveraging three distinct, publicly available RL-finetuned models, each representing a unique combination of architecture and alignment method:

- **SD1.5 + DPO:** A Stable Diffusion v1.5 model (Rombach et al., 2022) aligned using Direct Preference Optimization (DPO) (Wallace et al., 2024).
- **SDXL + SPO:** A Stable Diffusion XL model (Podell

Table 5. Multi-objective RLG fusion of aesthetics and OCR experts on a held-out subset of the OCR test prompt set. w_1 and w_2 control the relative strengths of the aesthetics-aligned and OCR-aligned experts, respectively.

Setting	w_1	w_2	Aesthetic \uparrow	PickScore \uparrow	OCR \uparrow
Aesthetics-only	1.0	0.0	5.76	23.76	0.592
OCR-only	0.0	1.0	5.11	22.24	0.936
Fused	1.0	1.0	5.51	23.44	0.926
Fused	1.0	1.8	5.38	22.88	0.940
Fused	2.4	1.0	5.76	23.99	0.894

Table 6. Effect of increasing the aesthetics expert weight w_1 while fixing the OCR expert weight at $w_2 = 1.0$ on a held-out subset of the OCR test prompt set.

w_1	w_2	Aesthetic \uparrow	PickScore \uparrow	OCR \uparrow
0.0	1.0	5.11	22.24	0.936
1.0	1.0	5.51	23.44	0.926
1.8	1.0	5.68	23.86	0.894
2.4	1.0	5.76	23.99	0.894

et al., 2023) aligned using Step-wise Preference Optimization (SPO) (Liang et al., 2024).

- **SD3.5-M + GRPO:** A Stable Diffusion 3.5 Medium flow matching model (Esser et al., 2024) aligned using Group Relative Policy Optimization (GRPO) (Liu et al., 2025b).

For evaluation, we use three established automated reward models: Aesthetic Score, ImageReward, and PickScore, with details provided in the appendix.

Results. Table 1 summarizes quantitative results, unequivocally demonstrating that RLG consistently delivers a significant performance boost across all configurations. The effect is particularly pronounced on the state-of-the-art GRPO-tuned SD3.5-M flow model, where RLG achieves a **74.95%** win rate on PickScore against the original finetuned model ($w_{RL} = 1.0$). As visually confirmed in Figures 1, 7, 8 and 9, increasing the RLG scale consistently enhances image detail and aesthetic appeal.

4.4. RLG Enables Flexible Control Over Alignment Strength

Standard RL fine-tuning fixes alignment strength, offering no inference-time flexibility. In contrast, RLG dynamically controls alignment strength with a powerful, training-free mechanism.

Table 7. RLG provides dynamic control over image compressibility. RLG allows for both interpolation and extrapolation beyond the original RL-tuned model’s capability ($w_{\text{RL}} = 1.0$).

Task	w_{RL}	Compression Ratio
Low Compressibility	0.4	1.14
	0.6	1.22
	1.0	1.35
	1.6	1.43
	3.0	1.37
High Compressibility	0.4	0.75
	0.6	0.64
	1.0	0.45
	1.6	0.18
	2.2	0.17

Controlling a Fundamental Property: Image Compressibility. We first demonstrate RLG’s control over image compressibility, a low-level property. We used two DDPO-finetuned SD1.4 models (Black et al., 2023) to reward either high or low image compressibility. Standard RL produces a model with fixed alignment. RLG transforms this static point into a dynamic spectrum. As shown in Table 7, RLG provides an inference-time ‘slider’ for alignment strength, a capability static fine-tuning lacks.

Balancing Competing Objectives: Text Accuracy vs. Aesthetics. Maximizing one alignment objective often compromises another. Table 3 illustrates this conflict. The standard RL-finetuned model ($w_{\text{RL}} = 1.0$) achieves 88.6% OCR accuracy, but its Aesthetic Score is fixed at 5.20. This trade-off is unalterable. RLG transforms this static outcome into flexible control. For instance, users prioritizing aesthetics over maximum text accuracy can set $w_{\text{RL}} < 1.0$. Conversely, others can push for peak accuracy at the cost of aesthetics by setting $w_{\text{RL}} = 2.8$.

4.5. RLG Outperforms Equivalent KL Training

We further investigate whether inference-time RLG offers practical advantages over training separate models under different KL constraints. Using Flow-GRPO on SD3.5-M for human preference alignment, we train a baseline RL-tuned model ($\beta = 0.001$) and compare it against a series of models independently trained with varying KL coefficients (β).

As shown in Figure 3, we map the KL-trained models to their equivalent RLG scales. The results show that RLG con-

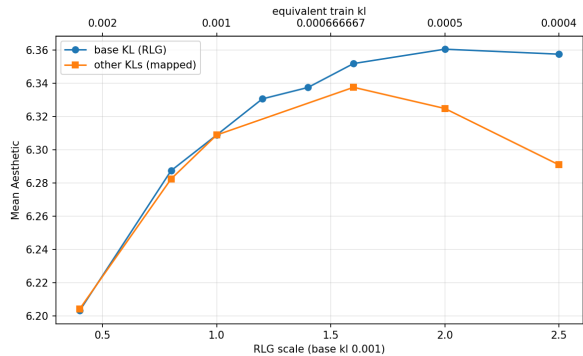


Figure 3. **RLG vs. Equivalent KL Training.** The blue curve represents RLG applied to a fixed base model ($\beta = 0.001$), while the orange curve represents models independently trained with different KL coefficients. RLG consistently yields higher Aesthetic Scores than retraining with equivalent KL constraints.

sistently outperforms the corresponding KL-trained variants in Aesthetic Score across the entire range. Whereas standard RL requires costly retraining to change the alignment strength, RLG turns a single fixed policy into a controllable spectrum at inference time. Crucially, this inference-time interpolation is more sample-efficient, achieving stronger alignment than explicitly relaxing the KL constraint during training. Additional details are provided in Appendix F.

5. Conclusion and Discussion

In this paper, we introduced Reinforcement Learning Guidance (RLG), a training-free method for dynamically controlling generative model alignment at inference time. By interpolating between, or extrapolating beyond, the base and RL-aligned models, RLG effectively adjusts the strength of KL regularization and steers generation toward higher-reward regions. Extensive experiments demonstrate consistent gains across diverse tasks, making RLG a simple yet powerful control layer over learned preferences.

Despite these strong empirical results, RLG has several limitations that motivate future work. First, because RLG inherits the same sampling mechanism as CFG, it also inherits CFG’s core limitation: CFG-based sampling does not guarantee convergence to the target marginal distribution, which also limits the scope of our analysis (Bradley & Nakkiran, 2024; Skreta et al., 2025). In addition, our link between the RLG scale w and the KL coefficient β assumes convergence under a standard reward-KL objective, an idealized setting that may not hold for methods such as GRPO (Vojnovic & Yun, 2025). Future work could explore adaptive timestep-dependent RLG scales and combinations with other orthogonal control methods.

6. Impact Statement

This work introduces Reinforcement Learning Guidance (RLG), a training-free method that improves the alignment and controllability of generative diffusion models at inference time. The main positive impact is practical: users can smoothly tune alignment strength without retraining, reducing compute cost and lowering the barrier to adapting models for diverse downstream objectives. This flexibility can benefit applications that require careful trade-offs between fidelity and alignment, such as accessibility, assistive creativity, and controllable content generation.

Like other alignment techniques, RLG can also be used to steer models toward undesirable objectives if misused. We therefore encourage responsible deployment, including clear usage policies and evaluation on safety-sensitive tasks. Because RLG operates at inference time, it may also be combined with existing safety filters and auditing tools, which can help mitigate potential risks.

References

- Authors, P. Paddleocr, awesome multilingual ocr toolkits based on paddlepaddle. <https://github.com/PaddlePaddle/PaddleOCR>, 2020.
- Black, K., Janner, M., Du, Y., Kostrikov, I., and Levine, S. Training diffusion models with reinforcement learning. *arXiv preprint arXiv:2305.13301*, 2023.
- Bradley, A. and Nakkiran, P. Classifier-free guidance is a predictor-corrector. *arXiv preprint arXiv:2408.09000*, 2024.
- Bui, D.-B., Nguyen, H.-K., and Le, T.-N. Prefpaint: Enhancing image inpainting through expert human feedback, 2025. URL <https://arxiv.org/abs/2506.21834>.
- Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., and Joulin, A. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 9650–9660, 2021.
- Cheng, B., Misra, I., Schwing, A. G., Kirillov, A., and Girdhar, R. Masked-attention mask transformer for universal image segmentation, 2022. URL <https://arxiv.org/abs/2112.01527>.
- Cheng, M., Doudi, F., Kalathil, D., Ghavamzadeh, M., and Kumar, P. R. Diffusion blend: Inference-time multi-preference alignment for diffusion models. *arXiv preprint arXiv:2505.18547*, 2025.
- Clark, K., Vicol, P., Swersky, K., and Fleet, D. J. Directly fine-tuning diffusion models on differentiable rewards. *arXiv preprint arXiv:2309.17400*, 2023.
- Cui, G., Yuan, L., Wang, Z., Wang, H., Li, W., He, B., Fan, Y., Yu, T., Xu, Q., Chen, W., et al. Process reinforcement through implicit rewards. *arXiv preprint arXiv:2502.01456*, 2025.
- Dhariwal, P. and Nichol, A. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.
- Esser, P., Kulal, S., Blattmann, A., Entezari, R., Müller, J., Saini, H., Levi, Y., Lorenz, D., Sauer, A., Boesel, F., Podell, D., Dockhorn, T., English, Z., Lacey, K., Goodwin, A., Marek, Y., and Rombach, R. Scaling rectified flow transformers for high-resolution image synthesis, 2024. URL <https://arxiv.org/abs/2403.03206>.
- Fan, J., Shen, S., Cheng, C., Chen, Y., Liang, C., and Liu, G. Online reward-weighted fine-tuning of flow matching with wasserstein regularization, 2025. URL <https://arxiv.org/abs/2502.06061>.
- Frans, K., Park, S., Abbeel, P., and Levine, S. Diffusion guidance is a controllable policy improvement operator. *arXiv preprint arXiv:2505.23458*, 2025.
- Ghosh, D., Hajishirzi, H., and Schmidt, L. Geneval: An object-focused framework for evaluating text-to-image alignment. *Advances in Neural Information Processing Systems*, 36:52132–52152, 2023.
- Gong, S., Zhang, R., Zheng, H., Gu, J., Jaitly, N., Kong, L., and Zhang, Y. Diffucoder: Understanding and improving masked diffusion models for code generation. *arXiv preprint arXiv:2506.20639*, 2025.
- Ho, J. and Salimans, T. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Huang, Q., Chan, L., Liu, J., He, W., Jiang, H., Song, M., and Song, J. Patchdpo: Patch-level dpo for finetuning-free personalized image generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 18369–18378, 2025.
- Kirstain, Y., Polyak, A., Singer, U., Matiana, S., Penna, J., and Levy, O. Pick-a-pic: An open dataset of user preferences for text-to-image generation, 2023. URL <https://arxiv.org/abs/2305.01569>.

- 495 Kong, W., Tian, Q., Zhang, Z., Min, R., Dai, Z., Zhou, J.,
496 Xiong, J., Li, X., Wu, B., Zhang, J., et al. Hunyuan-
497 video: A systematic framework for large video generative
498 models. *arXiv preprint arXiv:2412.03603*, 2024.
499
- 500 Lee, K., Liu, H., Ryu, M., Watkins, O., Du, Y., Boutilier, C.,
501 Abbeel, P., Ghavamzadeh, M., and Gu, S. S. Aligning text-
502 to-image models using human feedback. *arXiv preprint*
503 *arXiv:2302.12192*, 2023.
504
- 505 Li, J., Li, D., Savarese, S., and Hoi, S. Blip-2: Boot-
506 strapping language-image pre-training with frozen image
507 encoders and large language models, 2023. URL
508 <https://arxiv.org/abs/2301.12597>.
509
- 510 Liang, Z., Yuan, Y., Gu, S., Chen, B., Hang, T., Cheng, M.,
511 Li, J., and Zheng, L. Aesthetic post-training diffusion
512 models from generic preferences with step-by-step pref-
513 erence optimization. *arXiv preprint arXiv:2406.04314*,
514 2024.
515
- 516 Lipman, Y., Chen, R. T., Ben-Hamu, H., Nickel, M., and
517 Le, M. Flow matching for generative modeling. *arXiv*
518 *preprint arXiv:2210.02747*, 2022.
519
- 520 Liu, D., Li, S., Liu, Y., Li, Z., Wang, K., Li, X., Qin, Q.,
521 Liu, Y., Xin, Y., Li, Z., et al. Lumina-video: Efficient and
522 flexible video generation with multi-scale next-dit. *arXiv*
523 *preprint arXiv:2502.06782*, 2025a.
524
- 525 Liu, J., Liu, G., Liang, J., Li, Y., Liu, J., Wang, X., Wan,
526 P., Zhang, D., and Ouyang, W. Flow-grpo: Training
527 flow matching models via online rl. *arXiv preprint*
528 *arXiv:2505.05470*, 2025b.
529
- 530 Liu, X., Gong, C., and Liu, Q. Flow straight and fast:
531 Learning to generate and transfer data with rectified flow.
532 *arXiv preprint arXiv:2209.03003*, 2022.
533
- 534 Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin,
535 S., and Guo, B. Swin transformer: Hierarchical vision
536 transformer using shifted windows. In *Proceedings of the*
537 *IEEE/CVF international conference on computer vision*,
538 pp. 10012–10022, 2021.
539
- 540 Mohamed, S., Rosca, M., Figurnov, M., and Mnih, A.
541 Monte carlo gradient estimation in machine learning.
542 *Journal of Machine Learning Research*, 21(132):1–62,
543 2020.
544
- 545 Mori, S., Nishida, H., and Yamada, H. *Optical character*
546 *recognition*. John Wiley & Sons, Inc., 1999.
547
- 548 Peng, X., Kumar, A., Zhang, G., Levine, S., and Regres-
549 sion, A.-W. Simple and scalable off-policy reinforcement
learning. *arXiv preprint arXiv:1910.00177*.
- Podell, D., English, Z., Lacey, K., Blattmann, A., Dockhorn,
T., Müller, J., Penna, J., and Rombach, R. Sdxl: Im-
proving latent diffusion models for high-resolution image
synthesis, 2023. URL <https://arxiv.org/abs/2307.01952>.
- Prabhudesai, M., Goyal, A., Pathak, D., and Fragkiadaki,
K. Aligning text-to-image diffusion models with reward
backpropagation. 2023.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G.,
Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark,
J., Krueger, G., and Sutskever, I. Learning transferable
visual models from natural language supervision, 2021.
URL <https://arxiv.org/abs/2103.00020>.
- Rafailov, R., Sharma, A., Mitchell, E., Manning, C. D., Er-
mon, S., and Finn, C. Direct preference optimization:
Your language model is secretly a reward model. *Ad-
vances in Neural Information Processing Systems*, 36:
53728–53741, 2023.
- Rafailov, R., Hejna, J., Park, R., and Finn, C. From r to q^* :
Your language model is secretly a q-function, 2024. URL
<https://arxiv.org/abs/2404.12358>.
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and
Ommer, B. High-resolution image synthesis with latent
diffusion models. In *Proceedings of the IEEE/CVF con-
ference on computer vision and pattern recognition*, pp.
10684–10695, 2022.
- Ruiz, N., Li, Y., Jampani, V., Pritch, Y., Rubinstein, M.,
and Aberman, K. Dreambooth: Fine tuning text-to-image
diffusion models for subject-driven generation. In *Pro-
ceedings of the IEEE/CVF conference on computer vision*
and pattern recognition, pp. 22500–22510, 2023.
- Schuhmann, C. and Beaumont, R. Aesthetic-predictor:
A linear estimator on top of clip to predict the aes-
thetic quality of pictures. [https://github.com/
LAION-AI/aesthetic-predictor](https://github.com/LAION-AI/aesthetic-predictor), June 2021.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and
Klimov, O. Proximal policy optimization algorithms.
arXiv preprint arXiv:1707.06347, 2017.
- Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Bi, X., Zhang,
H., Zhang, M., Li, Y., Wu, Y., et al. Deepseekmath: Push-
ing the limits of mathematical reasoning in open language
models. *arXiv preprint arXiv:2402.03300*, 2024.
- Skreta, M., Akhound-Sadegh, T., Ohanesian, V., Bondesan,
R., Aspuru-Guzik, A., Doucet, A., Brekelmans, R., Tong,
A., and Neklyudov, K. Feynman-kac correctors in diffu-
sion: Annealing, guidance, and product of experts. *arXiv*
preprint arXiv:2503.02819, 2025.

- 550 Song, J., Meng, C., and Ermon, S. Denoising diffusion im-
551 plicit models. *arXiv preprint arXiv:2010.02502*, 2020a.
552
- 553 Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Er-
554 mon, S., and Poole, B. Score-based generative modeling
555 through stochastic differential equations. *arXiv preprint*
556 *arXiv:2011.13456*, 2020b.
557
- 558 Sun, X., Xiao, R., Mo, J., Wu, B., Yu, Q., and Wang, B.
559 F5r-tts: Improving flow-matching based text-to-speech
560 with group relative policy optimization. *arXiv preprint*
561 *arXiv:2504.02407*, 2025.
562
- 563 Tong, A., Fatras, K., Malkin, N., Huguët, G., Zhang, Y.,
564 Rector-Brooks, J., Wolf, G., and Bengio, Y. Improving
565 and generalizing flow-based generative models with mini-
566 batch optimal transport. *arXiv preprint arXiv:2302.00482*,
567 2023.
568
- 569 Vojnovic, M. and Yun, S.-Y. What is the alignment objective
570 of grpo? *arXiv preprint arXiv:2502.18548*, 2025.
571
- 572 Wallace, B., Dang, M., Rafailov, R., Zhou, L., Lou, A.,
573 Purushwalkam, S., Ermon, S., Xiong, C., Joty, S., and
574 Naik, N. Diffusion model alignment using direct prefer-
575 ence optimization. In *Proceedings of the IEEE/CVF*
576 *Conference on Computer Vision and Pattern Recognition*,
577 pp. 8228–8238, 2024.
- 578 Wan, T., Wang, A., Ai, B., Wen, B., Mao, C., Xie, C.-W.,
579 Chen, D., Yu, F., Zhao, H., Yang, J., et al. Wan: Open
580 and advanced large-scale video generative models. *arXiv*
581 *preprint arXiv:2503.20314*, 2025.
582
- 583 Williams, R. J. Simple statistical gradient-following algo-
584 rithms for connectionist reinforcement learning. *Machine*
585 *learning*, 8(3):229–256, 1992.
586
- 587 Xu, J., Liu, X., Wu, Y., Tong, Y., Li, Q., Ding, M., Tang,
588 J., and Dong, Y. Imagereward: Learning and evaluat-
589 ing human preferences for text-to-image generation. *Ad-*
590 *vances in Neural Information Processing Systems*, 36:
591 15903–15935, 2023.
592
- 593 Ye, H., Zhang, J., Liu, S., Han, X., and Yang, W. Ip-adapter:
594 Text compatible image prompt adapter for text-to-image
595 diffusion models. *arXiv preprint arXiv:2308.06721*,
596 2023.
597
- 598 Yujian, L. and Bo, L. A normalized levenshtein distance
599 metric. *IEEE transactions on pattern analysis and ma-*
600 *chine intelligence*, 29(6):1091–1095, 2007.
- 601 Zheng, Q., Le, M., Shaul, N., Lipman, Y., Grover, A., and
602 Chen, R. T. Guided flows for generative modeling and
603 decision making. *arXiv preprint arXiv:2311.13443*, 2023.
604
- Zhu, H., Xiao, T., and Honavar, V. G. Dspo: Direct score
preference optimization for diffusion model alignment.
In *The Thirteenth International Conference on Learning*
Representations, 2025.

A. The Use of Large Language Models (LLMs)

Large Language Models (LLMs) were used solely for language polishing and editorial refinement of this manuscript, including grammar correction and clarity improvements. All research content, methodology, analysis, and conclusions are the original work of the authors.

B. Model Specification

Table 8 lists the base models, RL-finetuned models, reward models, and their corresponding links.

C. Theoretical Derivations

This appendix provides the formal derivations discussed in the main paper.

C.1. Proof of the Optimal Policy for KL-Regularized RL

We aim to find the policy π^* that solves the optimization problem defined in Equation 7:

$$\pi^* = \arg \max_{\pi} \left(\mathbb{E}_{\mathbf{x} \sim \pi(\mathbf{x})} [R(\mathbf{x})] - \beta D_{\text{KL}}(\pi(\mathbf{x}) \parallel \pi_{\text{ref}}(\mathbf{x})) \right) \quad (18)$$

subject to the constraint that $\pi(\mathbf{x})$ is a valid probability distribution, i.e., $\int \pi(\mathbf{x}) d\mathbf{x} = 1$.

First, we expand the objective functional $J(\pi)$:

$$\begin{aligned} J(\pi) &= \int \pi(\mathbf{x}) R(\mathbf{x}) d\mathbf{x} - \beta \int \pi(\mathbf{x}) \log \frac{\pi(\mathbf{x})}{\pi_{\text{ref}}(\mathbf{x})} d\mathbf{x} \\ &= \int (\pi(\mathbf{x}) R(\mathbf{x}) - \beta \pi(\mathbf{x}) \log \pi(\mathbf{x}) + \beta \pi(\mathbf{x}) \log \pi_{\text{ref}}(\mathbf{x})) d\mathbf{x} \end{aligned} \quad (19)$$

This is a constrained optimization problem that can be solved using the calculus of variations with a Lagrange multiplier, λ , for the probability distribution constraint. The Lagrangian is:

$$\mathcal{L}(\pi, \lambda) = J(\pi) + \lambda \left(\int \pi(\mathbf{x}) d\mathbf{x} - 1 \right) \quad (20)$$

To find the optimal policy π^* , we take the functional derivative of \mathcal{L} with respect to $\pi(\mathbf{x})$ and set it to zero.

$$\begin{aligned} \frac{\delta \mathcal{L}}{\delta \pi(\mathbf{x})} &= \frac{\partial}{\partial \pi(\mathbf{x})} [\pi R - \beta \pi \log \pi + \beta \pi \log \pi_{\text{ref}} + \lambda \pi] = 0 \\ &= R(\mathbf{x}) - \beta (\log \pi(\mathbf{x}) + 1) + \beta \log \pi_{\text{ref}}(\mathbf{x}) + \lambda = 0 \end{aligned} \quad (21)$$

Now, we solve for $\log \pi(\mathbf{x})$:

$$\begin{aligned} \beta \log \pi(\mathbf{x}) &= R(\mathbf{x}) + \beta \log \pi_{\text{ref}}(\mathbf{x}) + \lambda - \beta \\ \log \pi(\mathbf{x}) &= \frac{1}{\beta} R(\mathbf{x}) + \log \pi_{\text{ref}}(\mathbf{x}) + \frac{\lambda - \beta}{\beta} \end{aligned} \quad (22)$$

Exponentiating both sides gives the form of the optimal policy $\pi^*(\mathbf{x})$:

$$\begin{aligned} \pi^*(\mathbf{x}) &= \exp \left(\frac{1}{\beta} R(\mathbf{x}) + \log \pi_{\text{ref}}(\mathbf{x}) + \frac{\lambda - \beta}{\beta} \right) \\ &= \pi_{\text{ref}}(\mathbf{x}) \exp \left(\frac{1}{\beta} R(\mathbf{x}) \right) \exp \left(\frac{\lambda - \beta}{\beta} \right) \end{aligned} \quad (23)$$

The term $\exp \left(\frac{\lambda - \beta}{\beta} \right)$ is a constant that does not depend on \mathbf{x} . This constant serves as the normalization factor to ensure that $\int \pi^*(\mathbf{x}) d\mathbf{x} = 1$. Let us denote this normalization constant as $1/Z(\beta)$. Therefore, the optimal distribution is:

$$\pi^*(\mathbf{x}) = \frac{1}{Z(\beta)} \pi_{\text{ref}}(\mathbf{x}) \exp \left(\frac{1}{\beta} R(\mathbf{x}) \right) \quad (24)$$

Table 8. Models and Their Corresponding Links.

Model/Reward Function	Link
SD3.5-M (Esser et al., 2024)	https://huggingface.co/stabilityai/stable-diffusion-3.5-medium
SD1.5 (Rombach et al., 2022)	https://huggingface.co/stable-diffusion-v1-5/stable-diffusion-v1-5
SD1.4 (Rombach et al., 2022)	https://huggingface.co/CompVis/stable-diffusion-v1-4
SDXL-base (Podell et al., 2023)	https://huggingface.co/stabilityai/stable-diffusion-xl-base-1.0
SD3.5M-FlowGRPO-PickScore (Liu et al., 2025b)	https://huggingface.co/jieliu/SD3.5M-FlowGRPO-PickScore
SD3.5M-FlowGRPO-Text (Liu et al., 2025b)	https://huggingface.co/jieliu/SD3.5M-FlowGRPO-Text
SD3.5M-FlowGRPO-GenEval (Liu et al., 2025b)	https://huggingface.co/jieliu/SD3.5M-FlowGRPO-GenEval
dpo-sd1.5-text2image-v1 (Lee et al., 2023)	https://huggingface.co/mhdang/dpo-sd1.5-text2image-v1
SPO-SDXL-4k-p-10ep (Liang et al., 2024)	https://huggingface.co/SPO-Diffusion-Models/SPO-SDXL_4k-p-10ep
dpo-sdxl-text2image-v1 (Lee et al., 2023)	https://huggingface.co/mhdang/dpo-sdxl-text2image-v1
ddpo-compressibility (Black et al., 2023)	https://huggingface.co/kvablack/ddpo-compressibility
ddpo-incompressibility (Black et al., 2023)	https://huggingface.co/kvablack/ddpo-incompressibility
Aesthetic Score (Schuhmann & Beaumont, 2021)	https://github.com/LAION-AI/aesthetic-predictor
ImageReward (Xu et al., 2023)	https://huggingface.co/THUDM/ImageReward
PickScore (Kirstain et al., 2023)	https://huggingface.co/yuvalkirstain/PickScore_v1
clip-vit-large-patch14 (Radford et al., 2021)	https://huggingface.co/openai/clip-vit-large-patch14
stable-diffusion-inpainting (Podell et al., 2023)	https://aihub.caict.ac.cn/models/runwayml/stable-diffusion-inpainting
prefpaint (Bui et al., 2025)	https://huggingface.co/kd5678/prefpaint-v1.0
prefpaintreward (Bui et al., 2025)	https://huggingface.co/kd5678/prefpaintReward
IP-Adapter-Plus (Ye et al., 2023)	https://huggingface.co/h94/IP-Adapter
PatchDPO (Huang et al., 2025)	https://huggingface.co/hqhQAQ/PatchDPO

This is equivalent to the proportional relationship given in Equation 8:

$$\pi^*(\mathbf{x}) \propto \pi_{\text{ref}}(\mathbf{x}) \exp\left(\frac{1}{\beta} R(\mathbf{x})\right) \quad (25)$$

This completes the proof.

C.2. Equivalence of the DPO Objective

The Direct Preference Optimization (DPO) framework is derived by re-parameterizing the KL-regularized RL objective in terms of preferences, thereby avoiding the need to explicitly train a reward model. The derivation shows that optimizing the DPO loss is equivalent to optimizing the policy towards the same theoretical distribution π^* derived above.

The derivation proceeds as follows:

1. **Express Reward in terms of Policies:** We start with the optimal policy solution from the previous section and rearrange it to solve for the reward function $R(\mathbf{x})$:

$$\begin{aligned} \pi^*(\mathbf{x}) &= \frac{1}{Z(\beta)} \pi_{\text{ref}}(\mathbf{x}) \exp\left(\frac{1}{\beta} R(\mathbf{x})\right) \\ \implies R(\mathbf{x}) &= \beta \log\left(\frac{\pi^*(\mathbf{x})}{\pi_{\text{ref}}(\mathbf{x})}\right) + \beta \log Z(\beta) \end{aligned} \quad (26)$$

The term $\beta \log Z(\beta)$ is a constant with respect to \mathbf{x} .

2. **Model Human Preferences:** Human preferences are typically collected as pairs $(\mathbf{x}_w, \mathbf{x}_l)$, where \mathbf{x}_w is preferred over \mathbf{x}_l . The Bradley-Terry model maps reward scores to preference probabilities:

$$p(\mathbf{x}_w \succ \mathbf{x}_l) = \sigma(R(\mathbf{x}_w) - R(\mathbf{x}_l)) \quad (27)$$

where $\sigma(\cdot)$ is the sigmoid function.

3. **Combine Reward and Preference Models:** We substitute the policy-based expression for the reward into the Bradley-Terry model. The constant term $\beta \log Z(\beta)$ cancels out perfectly:

$$\begin{aligned} &R(\mathbf{x}_w) - R(\mathbf{x}_l) \\ &= \left(\beta \log \frac{\pi^*(\mathbf{x}_w)}{\pi_{\text{ref}}(\mathbf{x}_w)} + C\right) - \left(\beta \log \frac{\pi^*(\mathbf{x}_l)}{\pi_{\text{ref}}(\mathbf{x}_l)} + C\right) \\ &= \beta \left(\log \frac{\pi^*(\mathbf{x}_w)}{\pi_{\text{ref}}(\mathbf{x}_w)} - \log \frac{\pi^*(\mathbf{x}_l)}{\pi_{\text{ref}}(\mathbf{x}_l)}\right) \end{aligned} \quad (28)$$

Thus, the ground-truth preference probability can be expressed entirely in terms of the optimal policy π^* and the reference policy π_{ref} :

$$p(\mathbf{x}_w \succ \mathbf{x}_l) = \sigma\left(\beta \left(\log \frac{\pi^*(\mathbf{x}_w)}{\pi_{\text{ref}}(\mathbf{x}_w)} - \log \frac{\pi^*(\mathbf{x}_l)}{\pi_{\text{ref}}(\mathbf{x}_l)}\right)\right) \quad (29)$$

4. **Construct the DPO Loss:** DPO seeks to find a policy π_θ that maximizes the log-likelihood of the observed human preferences. This is equivalent to minimizing the negative log-likelihood loss. By replacing the theoretical optimal policy π^* with our parameterized model policy π_θ , we arrive at the DPO loss function:

$$\mathcal{L}_{\text{DPO}} = -\mathbb{E}_{(\mathbf{x}_w, \mathbf{x}_l) \sim \mathcal{D}} (\log p(\mathbf{x}_w \succ \mathbf{x}_l)) \quad (30)$$

By minimizing this loss, we are directly training the policy π_θ to satisfy the same mathematical relationship that defines the optimal RL policy π^* . Therefore, the policy obtained by successfully optimizing the DPO objective, π_{DPO}^* , converges to the same theoretical distribution as the one found by KL-regularized RL, where the reward function $R(\mathbf{x})$ is implicitly defined by the human preference dataset.

C.3. Derivation of Velocity-Score Relationship

This section provides a detailed derivation of the relationship between the velocity field $\mathbf{v}(\mathbf{x}, t)$ used in Flow Matching models and the score function $\mathbf{s}(\mathbf{x}, t)$ used in Denoising Diffusion Models, as stated in Equation 2.

The unifying perspective relies on a common reference path $(X_t)_{t \in [0,1]}$ that interpolates between an initial noise variable $X_1 \sim p_1 = \mathcal{N}(0, \mathbf{I})$ and a data sample $X_0 \sim p_{\text{data}}$. This path is defined by linear interpolation:

$$X_t = \beta_t X_1 + \alpha_t X_0 \quad (31)$$

where α_t and β_t are scalar functions of time, with $\alpha_0 = \beta_1 = 0$ and $\alpha_1 = \beta_0 = 1$.

In Denoising Diffusion Models, the model learns to predict the score function $\mathbf{s}(\mathbf{X}_t, t) = \nabla_{\mathbf{X}_t} \log p_t(\mathbf{X}_t)$. For the chosen linear interpolation path where $X_1 \sim \mathcal{N}(0, \mathbf{I})$, it's a known property that the score function is related to the conditional expectation of X_0 and X_1 given X_t . Specifically, the optimal denoised estimate of X_0 , denoted $\hat{X}_0(\mathbf{X}_t, t)$, and the optimal estimate of the noise X_1 , denoted $\hat{X}_1(\mathbf{X}_t, t)$, can be expressed in terms of X_t and its score:

$$\hat{X}_0(\mathbf{X}_t, t) = \mathbb{E}[X_0|X_t] = \frac{1}{\alpha_t}(\mathbf{X}_t + \beta_t^2 \mathbf{s}(\mathbf{X}_t, t)) \quad (32)$$

$$\hat{X}_1(\mathbf{X}_t, t) = \mathbb{E}[X_1|X_t] = -\beta_t \mathbf{s}(\mathbf{X}_t, t) \quad (33)$$

These relationships hold under the assumption that the conditional distribution $p(X_t|X_0)$ is a Gaussian $X_t \sim \mathcal{N}(\alpha_t X_0, \beta_t^2 \mathbf{I})$, which is implied by the path definition with $X_1 \sim \mathcal{N}(0, \mathbf{I})$.

For Flow Matching models, the objective is to learn a velocity field $\mathbf{v}(\mathbf{X}_t, t)$ that describes the deterministic trajectory of samples via an ordinary differential equation $d\mathbf{X}_t = \mathbf{v}(\mathbf{X}_t, t)dt$. This velocity field matches the time derivative of the reference flow, $\frac{d}{dt} X_t$. Differentiating Equation 31 with respect to time t :

$$\mathbf{v}(\mathbf{X}_t, t) = \frac{d}{dt} X_t = \dot{\beta}_t X_1 + \dot{\alpha}_t X_0 \quad (34)$$

To express the velocity field in terms of the current state \mathbf{X}_t and the score function $\mathbf{s}(\mathbf{X}_t, t)$, we substitute the expressions for $\hat{X}_0(\mathbf{X}_t, t)$ (Equation 32) and $\hat{X}_1(\mathbf{X}_t, t)$ (Equation 33) into Equation 34:

$$\begin{aligned} \mathbf{v}(\mathbf{X}_t, t) &= \dot{\beta}_t(-\beta_t \mathbf{s}(\mathbf{X}_t, t)) + \dot{\alpha}_t \left(\frac{1}{\alpha_t}(\mathbf{X}_t + \beta_t^2 \mathbf{s}(\mathbf{X}_t, t)) \right) \\ &= -\dot{\beta}_t \beta_t \mathbf{s}(\mathbf{X}_t, t) + \frac{\dot{\alpha}_t}{\alpha_t} \mathbf{X}_t + \frac{\dot{\alpha}_t}{\alpha_t} \beta_t^2 \mathbf{s}(\mathbf{X}_t, t) \end{aligned}$$

Rearranging the terms by grouping \mathbf{X}_t and $\mathbf{s}(\mathbf{X}_t, t)$:

$$\begin{aligned} \mathbf{v}(\mathbf{X}_t, t) &= \left(\frac{\dot{\alpha}_t}{\alpha_t} \right) \mathbf{X}_t + \left(\frac{\dot{\alpha}_t}{\alpha_t} \beta_t^2 - \dot{\beta}_t \beta_t \right) \mathbf{s}(\mathbf{X}_t, t) \\ &= \left(\frac{\dot{\alpha}_t}{\alpha_t} \right) \mathbf{X}_t + \beta_t \left(\frac{\dot{\alpha}_t}{\alpha_t} \beta_t - \dot{\beta}_t \right) \mathbf{s}(\mathbf{X}_t, t) \end{aligned} \quad (35)$$

This derivation confirms Equation 2 from the main paper, establishing the precise mathematical connection between the velocity field learned by Flow Matching and the score function predicted by Denoising Diffusion Models under the common linear interpolation path.

D. Derivation of the Implicit Time-Dependent Reward

To formalize this, we first need to establish that for any given generative model policy π_θ (represented by its distribution $p_{\theta,t}$), we can define a corresponding reward function for which π_θ is the optimal policy. This concept is well-established in inverse reinforcement learning for discrete MDPs, such as those used for aligning LLMs (Sun et al., 2025; Rafailov et al., 2023). We can extend this framework to diffusion models by considering the generation process as a continuous-time MDP (Black et al., 2023; Rafailov et al., 2024).

In this diffusion MDP, the state at time t is the noisy sample \mathbf{x}_t , and the policy $\pi(\cdot|\mathbf{x}_t)$ determines the transition to the next state \mathbf{x}_{t-dt} . Recent theoretical work has shown a bijection between reward functions and optimal Q-functions (and thus optimal policies) in such MDPs. Specifically, following (Rafailov et al., 2024), for a given reference policy π_{ref} and a temperature parameter β , the optimal policy π^* for a reward function $r(s_t, a_t)$ satisfies the relationship:

$$\beta \log \frac{\pi^*(a_t|s_t)}{\pi_{\text{ref}}(a_t|s_t)} = r(s_t, a_t) + \Phi(s_{t-dt}) - \Phi(s_t) \quad (36)$$

where Φ is a potential function, corresponding to the optimal value function V^* . This means that any policy π_θ can be viewed as the optimal policy for an implicitly defined reward function, equivalent to the log-policy ratio up to a potential-based shaping term.

By adapting this principle to the continuous state-space of diffusion models, we can define an instantaneous, time-dependent reward function $R_t(\mathbf{x}_t)$ directly in terms of the model’s probability density. The policy $\pi_\theta(\cdot|\mathbf{x}_t)$ is governed by the underlying score function $\mathbf{s}_\theta(\mathbf{x}_t, t)$, which itself is the gradient of the log-density $\log p_{\theta,t}(\mathbf{x}_t)$. We can therefore define the implicit reward by relating the marginal densities of the RL-aligned model ($p_{\theta,t}$) and the reference model ($p_{\text{ref},t}$):

$$R_t(\mathbf{x}_t) \triangleq \beta \log \frac{p_{\theta,t}(\mathbf{x}_t)}{p_{\text{ref},t}(\mathbf{x}_t)} \quad (37)$$

Here, $p_{\theta,t}$ is the marginal probability distribution of the noisy image \mathbf{x}_t under the RL-aligned model, and $p_{\text{ref},t}$ is the corresponding distribution for the pre-trained reference model. The parameter β is the same KL-regularization coefficient from the original RL objective (Equation 7). This equation defines the reward that the RL-aligned model π_θ is implicitly optimizing for at every point (\mathbf{x}_t, t) in the generation process, relative to the reference model.

E. One-dimensional case study of model convergence

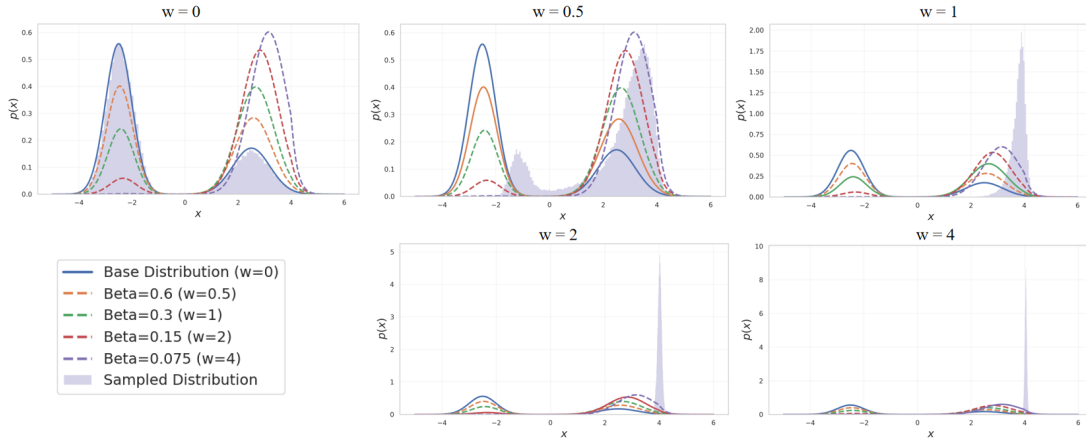


Figure 4. Small-scale demonstration supporting the theoretical justification of RLG. Each subplot shows the sampled distribution under a different RLG weight w , while the curves represent the corresponding theoretically predicted RL-fine-tuned distributions. Here, β denotes the KL regularization coefficient.

We further investigate how RLG behaves under different RL algorithms when the RL policy does *not* converge to the ideal reward–reweighted distribution, and why RLG can still improve performance in this realistic regime. To this end, we construct a small-scale one-dimensional case study. Our setting uses a flow-matching model defined on the real line with a pretrained base Gaussian mixture distribution

$$p_{\text{base}}(x) \triangleq 0.7 \mathcal{N}(-2.5, 0.25) + 0.3 \mathcal{N}(2.5, 0.49).$$

To bias the model toward a specific region, we introduce a reward function that assigns higher reward to samples near $x = 4$; concretely, we use $r(x) = -0.1|x - 4|$ (up to an additive constant), so that points closer to $x = 4$ receive larger reward. We then fine-tune the pretrained flow with either a vanilla policy-gradient method or the GRPO algorithm, using a KL coefficient $\beta = 0.3$, batch size 64, and learning rate 1×10^{-5} , with the pretrained flow as the reference policy.

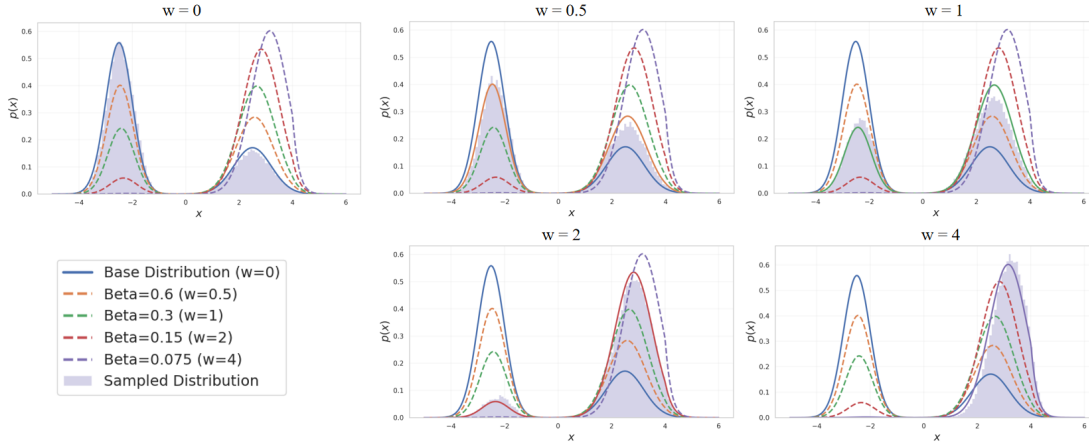


Figure 5. Small-scale demonstration supporting the theoretical justification of RLG. Each subplot shows the sampled distribution under a different RLG weight w , while the curves represent the corresponding theoretically predicted RL-fine-tuned distributions. Here, β denotes the KL regularization coefficient.

Figures 4 and 5 show the sampled distributions under different RLG weights w , together with the corresponding theoretical reward-reweighted target $p_{\text{rl}}(x) \propto p_{\text{base}}(x) \exp(r(x)/\beta)$. For both the policy-gradient and GRPO experts, increasing w from 0.5 to 1, 2, and 4 consistently shifts probability mass toward the high-reward region around $x = 4$: the left-hand mode is gradually suppressed, and the right-hand mode becomes sharper near the reward peak. In the GRPO case, the standalone RL policy already produces a very sharp and imperfect approximation of the target distribution, yet applying RLG on top of this expert still moves the overall sampler closer to the high-reward region. This toy example illustrates that even when the RL policy does not exactly realize the ideal distribution (as is typical for practical algorithms such as GRPO), RLG with a moderate weight w can reliably steer samples toward higher-reward regions and improve alignment over both the base model and the raw RL-fine-tuned model.

F. RLG vs. Equivalent KL Training Across Scales

We designed this experiment to verify that inference-time RLG achieves stronger alignment than retraining with an equivalent KL coefficient. We trained Flow-GRPO models on the PickScore human-preference dataset using SD3.5-M as the base model with four GPUs. Training used 512×512 resolution, 10 sampling steps (40 for evaluation), and no classifier-free guidance (guidance scale 1.0). The batch size was 9 per GPU with 18 images per prompt and gradient accumulation; we used one inner epoch, a small clipping range of 10^{-5} , EMA, and the same stochastic settings across runs (global std enabled, noise level 0.8, CPS SDE with a short window). Our base KL coefficient is $\beta = 0.001$, and we trained additional models by varying β , each for 1200 steps.

We then applied RLG to the base model and compared its performance against these independently trained KL variants. Figure 6 plots the Aesthetic Score as a function of the RLG scale for the base model (blue) and overlays the performance of KL-trained models (orange) mapped to their equivalent RLG scales. The results show that RLG consistently outperforms the model trained with the corresponding KL across the scale range.

Two trends are evident. First, increasing the RLG scale on the fixed base model yields a smooth, monotonic improvement. Second, the KL-trained variants lie consistently below the RLG curve at their matched scales, even though they require separate training runs. This gap suggests that the inference-time interpolation in RLG is not only more flexible but also more sample-efficient in practice, delivering stronger alignment without re-optimizing the policy for each β .

G. Additional Ablations and Analysis

In this section, we provide further experimental validations conducted to assess the statistical significance, computational efficiency, and multi-objective capabilities of Reinforcement Learning Guidance (RLG).

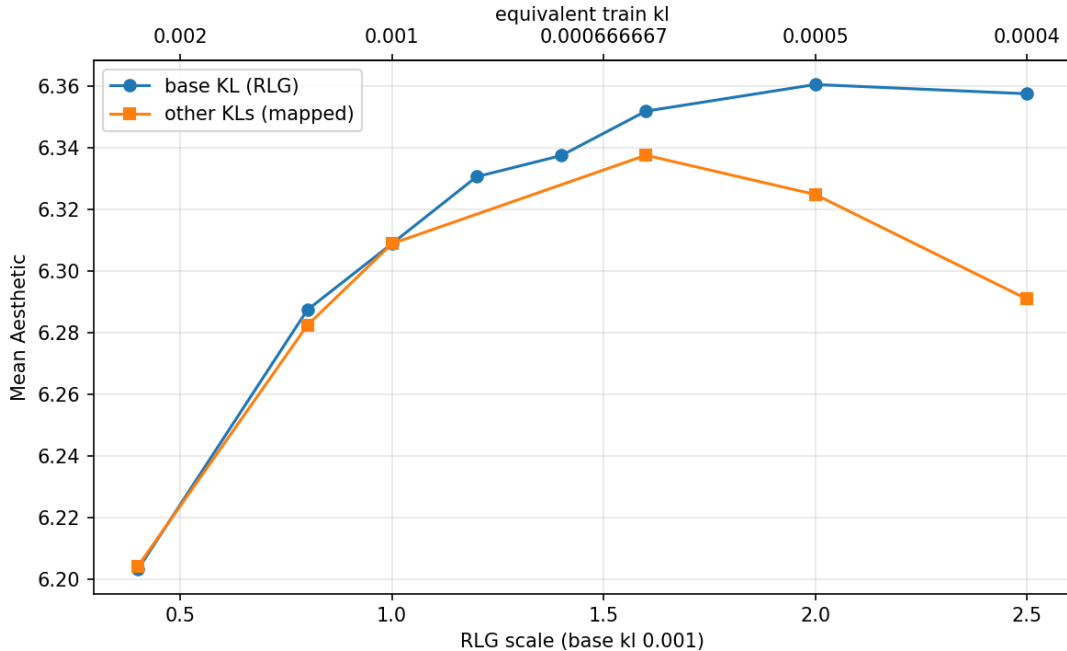


Figure 6. Aesthetic Score versus RLG scale for the base model ($\beta = 0.001$), with the top axis indicating the equivalent training KL. The blue curve is RLG applied to the base model; the orange curve corresponds to independently trained Flow-GRPO models with different KL coefficients, mapped to the equivalent RLG scale. RLG consistently yields higher scores than the corresponding KL-trained models across the scale range.

G.1. Statistical Significance and Multi-Seed Evaluation

To ensure the robustness of our improvements, we evaluated the Flow-GRPO and SPO models on the DrawBench aesthetics task using five random seeds. We report the Mean \pm Standard Error of the Mean (SEM). As shown in Table 9, RLG yields statistically significant improvements ($p < 10^{-5}$ via one-sided t-test) in both Aesthetic Score and PickScore compared to the RL-only baseline, confirming that the gains are not artifacts of seed selection.

Table 9. Multi-seed evaluation (5 seeds) on DrawBench. RLG consistently improves preference metrics with negligible variance.

Model	Guidance	Aesthetic \uparrow	PickScore \uparrow	ImageReward \uparrow
Flow-GRPO	RL-only ($w = 1.0$)	6.44 \pm 0.002	23.28 \pm 0.004	1.40 \pm 0.002
Flow-GRPO	RLG ($w = 2.2$)	6.64 \pm 0.002	23.58 \pm 0.003	1.39 \pm 0.003
SPO	RL-only ($w = 1.0$)	6.42 \pm 0.001	23.69 \pm 0.007	1.13 \pm 0.005
SPO	RLG ($w = 1.4$)	6.49 \pm 0.003	23.73 \pm 0.005	1.15 \pm 0.004

G.2. Interaction between RLG and CFG Scales

To verify that RLG provides additive value beyond simply correcting the Classifier-Free Guidance (CFG) scale, we performed a grid search over CFG and RLG scales on the DrawBench task. Table 10 demonstrates that for any fixed CFG scale, increasing the RLG weight (w_{RL}) improves performance. Furthermore, a two-way ANOVA analysis ($F_{RLG} = 2960.5$, $p < 10^{-27}$) confirms that RLG contributes a strong, independent main effect to the generation quality, rather than merely compensating for suboptimal CFG.

G.3. Compute Efficiency vs. Test-Time Scaling

RLG requires two forward passes per step. We investigated whether the same compute budget would be better spent simply increasing the number of diffusion steps for the single RL-tuned model. As shown in Table 11, increasing the RL-only

Table 10. PickScore performance across a grid of CFG and RLG scales. The combination of CFG=4.5 and RLG=2.4 yields the global optimum.

CFG Scale	RLG Scale (w_{RL})				
	1.4	1.8	2.2	2.4	3.0
2.0	23.14	23.27	23.31	23.32	23.32
3.0	23.42	23.50	23.54	23.55	23.54
4.0	23.48	23.56	23.58	23.59	23.57
4.5 (Default)	23.48	23.56	23.57	23.59	23.58
5.0	23.47	23.55	23.57	23.58	23.57

model from 20 to 60 steps yields no performance gain. Conversely, RLG at 20 steps (computationally equivalent to ~ 40 steps of a single model) provides a substantial boost in Aesthetic Score (+0.19) and PickScore (+0.29), making it a more effective test-time scaling strategy.

Table 11. Comparison of RLG against increasing diffusion steps for the RL-only baseline.

Method	Steps	Effective Cost	Aesthetic \uparrow	PickScore \uparrow
RL-only	20	1 \times	6.45	23.29
RL-only	40	2 \times	6.41	23.30
RL-only	60	3 \times	6.40	23.29
RLG ($w = 2.2$)	20	2\times	6.64	23.58

H. Experimental Details

H.1. Generation Hyper-parameters

This section provides a detailed overview of the parameters and settings used for the text-to-image generation experiments discussed in the main paper.

To ensure consistency across our evaluations, several parameters were standardized across all models. All images were generated at a resolution of 512×512 pixels, and text prompts were processed using each model’s native text encoder. For the generative process, we uniformly set the number of sampling steps to 20 for all experiments to balance computational cost and output quality.

The Classifier-Free Guidance (CFG) scale, which controls the adherence to the text prompt, was set to the generally recommended value for each model to ensure optimal baseline performance. The specific CFG scales used were:

- **Stable Diffusion 1.4:** A CFG scale of 7.5 was used.
- **Stable Diffusion 1.5:** A CFG scale of 7.5 was used.
- **Stable Diffusion inpainting:** A CFG scale of 7.5 was used.
- **Stable Diffusion XL:** A CFG scale of 5.0 was used.
- **Stable Diffusion 3.5:** A CFG scale of 4.5 was used.

These CFG scales were held constant across all experiments for a given model, allowing for a direct assessment of the impact of our RLG guidance scale (w). All other parameters were kept at the default settings of the corresponding Diffusers pipeline.

H.2. Human Preference Evaluation Metric Details

To quantitatively evaluate the performance of our method, we established automated reward models. These models are designed to assess different aspects of image quality and text-to-image alignment, providing a comprehensive evaluation

framework.

- **Aesthetic Score:** This metric provides a general assessment of an image’s visual appeal. It uses a pre-trained CLIP model (`clip-vit-large-patch14`) (Radford et al., 2021) to extract an image embedding. This embedding is then processed by a multilayer perceptron (MLP) head, loaded with the weights provided in (Liu et al., 2025b), which maps the features to a single scalar score, typically on a 1–10 scale. The corresponding links are listed in Table 8. A higher score indicates higher predicted aesthetic quality according to the human judgments used to train the MLP.
- **ImageReward:** Developed by (Xu et al., 2023), this is a sophisticated reward model designed to evaluate both the semantic alignment of an image to its text prompt and its overall visual fidelity. Built upon the BLIP-2 architecture (Li et al., 2023), it was fine-tuned on a large-scale dataset of human preference feedback, enabling it to serve as a robust, general-purpose proxy for human judgment in text-to-image generation tasks.
- **PickScore:** Introduced by (Kirstain et al., 2023), this reward model is specifically trained to predict human preferences from direct pairwise comparisons. It is derived from the extensive Pick-a-Pic dataset, which contains a large number of human choices between two images generated from the same prompt. We use the v1 version, which leverages a powerful CLIP-ViT-H-14 model (Radford et al., 2021). Its strong correlation with human preference makes it particularly relevant to our work.

I. GenEval Benchmark Details

The **GenEval** benchmark provides an automated, object-focused framework for evaluating the compositional capabilities of text-to-image models (Ghosh et al., 2023). Unlike holistic metrics that measure overall image quality or text alignment, GenEval offers a fine-grained analysis of a model’s ability to adhere to specific compositional instructions within a prompt.

The official test set comprises 553 prompts. These prompts are designed to probe several distinct aspects of compositional generation:

- **Single Object:** Tests the model’s fundamental ability to render a single specified object.
- **Two Objects:** Assesses the capacity to generate two distinct objects in the same image, testing for co-occurrence.
- **Counting:** Evaluates whether the model can generate a precise number of a given object.
- **Colors:** Measures if an object can be rendered with a specific, designated color.
- **Position:** Tests spatial reasoning by requiring two objects to be placed in a specified relative position (e.g., “a cat to the left of a dog”).
- **Attribute Binding:** The most complex task, which requires binding specific attributes (like color) to specific objects (e.g., generating “a red cube and a blue sphere” without swapping attributes).

The evaluation protocol is fully automated, using a sophisticated object detection model to parse the generated images. For our experiments, we adhered to the official methodology, which employs a **Mask2Former** (Cheng et al., 2022) model with a Swin-S transformer backbone (Liu et al., 2021). This detector identifies objects and verifies their properties and spatial arrangements as dictated by the input prompt.

J. Image Compressibility Experimental Details

This section provides a detailed description of the experimental setup for the image compressibility task.

Models and Task Definition. The goal of this experiment was to verify that Reinforcement Learning Guidance (RLG) can effectively control a low-level, non-semantic property of generated images: their compressibility. We used the standard Stable Diffusion v1.4 model as our base reference (\mathbf{v}_{ref}). For the aligned models (\mathbf{v}_θ), we utilized two sets of weights from the official DDPO implementation (Black et al., 2023):

- **Low Compressibility Model:** We use ddp-compressibility. This model was fine-tuned to generate images that are less compressible, resulting in larger file sizes when saved in JPEG format. This typically corresponds to images with higher texture detail and complexity.
- **High Compressibility Model:** We use ddp-incompressibility. This model was fine-tuned to prefer images that are more compressible, resulting in smaller JPEG file sizes. This often corresponds to images with smoother regions and less high-frequency detail.

Dataset and Prompts. Following the DDPO study (Black et al., 2023), our evaluation prompts were based on animal classes from the ImageNet dataset. We used 45 distinct animal classes. For each class, we generated 4 images, resulting in a total of 180 images per RLG scale setting. The prompt template used was: “{class_name}”.

The 45 animal classes used are: ant, bat, bear, bee, beetle, bird, butterfly, camel, cat, chicken, cow, deer, dog, dolphin, duck, fish, fly, fox, frog, goat, goose, gorilla, hedgehog, horse, kangaroo, lion, lizard, llama, monkey, mouse, pig, rabbit, raccoon, rat, shark, sheep, snake, spider, squirrel, tiger, turkey, turtle, whale, wolf, and zebra.

Evaluation Metric. We evaluated performance using the **Compression Ratio**. For a given prompt, let \mathbf{x}_{base} be the image generated by the base SD1.4 model, and let $\mathbf{x}_{\text{RLG}}(w)$ be the image generated using RLG with a guidance scale of w . Let $S_{\text{jpeg}}(\cdot)$ denote the file size of an image after being compressed and saved in JPEG format. The Compression Ratio is defined as:

$$\text{Compression Ratio}(w) = \frac{S_{\text{jpeg}}(\mathbf{x}_{\text{RLG}}(w))}{S_{\text{jpeg}}(\mathbf{x}_{\text{base}})}$$

The final score reported in Table 7 is the mean of this ratio, averaged across all 180 generated images for each guidance scale. A ratio of 1.0 indicates no change in compressibility compared to the base model.

K. Image Inpainting Experimental Details

This section provides a comprehensive overview of the experimental setup, models, and evaluation protocol used for the image inpainting task discussed in the main paper. Our methodology closely follows that of the PrefPaint (Bui et al., 2025) study to ensure a fair and direct comparison.

K.1. Models and Task Definition

The experiment focuses on conditional image generation for image inpainting, the task of filling in masked (missing) regions of an image in a semantically and visually plausible manner.

- **Base Model (\mathbf{v}_{ref}):** We use the standard Stable Diffusion inpainting model as our baseline. This model, accessible on diffusers as `runwayml/stable-diffusion-inpainting`, is widely used and serves as the un-aligned reference point for our experiments.
- **RL-Finetuned Model (\mathbf{v}_{θ}):** For the human-preference-aligned expert model, we employ **PrefPaint** (Bui et al., 2025). This model is a direct descendant of the base model, which has been further fine-tuned using reinforcement learning. The training process for PrefPaint leveraged a large-scale dataset of over 51,000 human preference judgments on inpainted images, making it an expert policy specialized in generating completions that align with human aesthetic and contextual expectations.

K.2. Evaluation Dataset and Protocol

All quantitative results were generated using the dataset provided by the authors of PrefPaint (Bui et al., 2025). The test set was constructed using the following procedure:

1. **Image Sourcing:** A diverse set of high-resolution images was initially sourced from established datasets such as ADE20K and ImageNet. All images were resized to a standard 512×512 pixel resolution.
2. **Mask Generation:** To simulate realistic inpainting and outpainting scenarios, two distinct masking strategies were employed to create the incomplete images:

- **Warping Holes (for Inpainting):** This method creates complex, non-rectangular masks inside the image. It simulates the disocclusion that occurs from a slight change in camera viewpoint. A depth map is first estimated for the source image, and then a new virtual camera view is generated with small shifts. The newly visible (disoccluded) regions form the mask that the model must fill. This tests the model’s ability to reason about 3D geometry and handle irregular shapes.
- **Boundary Masks (for Outpainting):** This strategy tests the model’s ability to extend a scene beyond its original borders. Masks are created at the edges of the image using two different cropping techniques:
 - *Square Cropping:* A central square region, covering 75% to 85% of the image area, is preserved, masking the outer frame.
 - *Rectangular Cropping:* The full height of the image is preserved, while a central vertical slice, comprising 60% to 65% of the original width, is kept, masking the left and right sides.

K.3. Evaluation Metrics

The quality of the generated inpainted images was assessed using the following automated metric:

- **Preference Reward:** We use the specialized reward model developed and released as part of the PrefPaint study (Bui et al., 2025). This model was trained on their custom dataset of nearly 51,000 human preference annotations. Unlike a general-purpose aesthetics model, it is specifically tailored to judge the quality of image inpainting, considering factors like structural rationality, local texture coherence, and overall aesthetic feeling. The reward scores reported in our table are the normalized values from this model, averaged over the official test set, as done in the original paper.

L. Personalized Image Generation Experimental Details

This appendix provides a detailed overview of the experimental setup for evaluating Reinforcement Learning Guidance (RLG) on the task of personalized image generation, as presented in the main paper.

L.1. Task and Model Background

Task Definition. Personalized image generation aims to synthesize novel images of a specific subject provided through one or more reference images. The model is given a reference image containing the subject (e.g., a specific pet dog) and a text prompt (e.g., ”a photo of [V] sleeping on a couch,” where [V] is a placeholder for the subject). The primary goal is to generate an image that not only matches the prompt’s description but also maintains high fidelity to the unique appearance and characteristics of the subject in the reference image.

Model Selection. Our experiment is designed to test whether RLG can amplify the effects of a fine-grained, RL-based alignment process. We therefore select models based on the work of PatchDPO (Huang et al., 2025).

- **Base Model (v_{ref}):** We use the publicly available, pre-trained **IP-Adapter-Plus** (Ye et al., 2023) model built on SDXL (Podell et al., 2023). IP-Adapter is a powerful method for subject-driven generation that injects image features into the cross-attention layers of a diffusion model. We use it as our baseline because it represents a strong, general-purpose personalization model before any preference-based fine-tuning.
- **RL-aligned Model (v_{θ}):** We use the model fine-tuned from IP-Adapter-Plus using the **PatchDPO** algorithm. PatchDPO is a form of preference optimization that operates at a sub-image or ”patch” level. During its training, generated images are compared against the reference image. Patches from the generated image that are consistent with the reference subject receive a positive reward, while inconsistent patches are penalized. This process, analogous to reinforcement learning with fine-grained rewards, tunes the model to be highly specialized in preserving subject fidelity.

L.2. Benchmark and Evaluation Metrics

Benchmark Dataset. All evaluations are conducted on **DreamBench** (Ruiz et al., 2023), the standard benchmark for personalized image generation. DreamBench consists of 30 unique subjects, each with a set of reference images and 80 corresponding text prompts. This benchmark is designed to test a model’s ability to generate the subject in various contexts, poses, and interactions.

Evaluation Metrics. To quantitatively measure the fidelity of the generated images to the reference subject, we employ two standard, complementary metrics. For each prompt in DreamBench, we generate an image and compare it to the ground-truth reference images of the subject.

- **CLIP-I (Image Similarity):** This metric, introduced by the DreamBooth authors, measures the semantic similarity between the generated and reference images. It works by encoding both images into high-dimensional feature vectors using a pre-trained CLIP ViT-L/14 image encoder. The final score is the average cosine similarity between the embedding of the generated image and the embeddings of the reference images. A higher CLIP-I score indicates that the generated image is semantically and stylistically closer to the reference subject from the perspective of the CLIP model.
- **DINO (Structural Similarity):** This metric uses features extracted from a self-supervised ViT-S/16 DINO (Caron et al., 2021) model. DINO is trained without labels and learns to capture rich information about object structure, texture, and shape. The metric is calculated as the average cosine similarity between the DINO features of the generated and reference images. It is particularly effective at measuring the preservation of fine-grained details and the structural integrity of the subject, making it an excellent indicator of subject fidelity.

M. Detailed Human Preference Alignment Results

This section provides the complete quantitative results for the human preference alignment experiments, complementing the summary presented in Table 1 of the main paper. For each model, we present two tables: one detailing the absolute mean scores for Aesthetic Score, ImageReward, and PickScore across various RLG guidance scales (w_{RL}), and another showing the pairwise win rates against the base model ($w_{RL} = 0.0$) and the standard RL-finetuned model ($w_{RL} = 1.0$).

Table 12. Mean scores for the **SD3.5-M** model series on human preference metrics across various RLG scales (w_{RL}). The scale $w_{RL} = 0.0$ corresponds to the original **SD3.5-M** base model, while $w_{RL} = 1.0$ represents the model after GRPO finetuning, named **SD3.5M-FlowGRPO-PickScore**.

w_{RL}	Aesthetic Score (\uparrow)	ImageReward (\uparrow)	PickScore (\uparrow)
0.0	5.97	0.99	21.75
1.0	6.45	1.40	23.29
1.2	6.48	1.41	23.36
1.4	6.54	1.40	23.48
1.6	6.57	1.40	23.53
1.8	6.60	1.41	23.56
2.0	6.62	1.40	23.57
2.2	6.64	1.39	23.58
2.4	6.66	1.39	23.58
2.6	6.68	1.37	23.59
2.8	6.68	1.36	23.56

Table 13. Win rates (%) for **SD3.5-M** model series at various RLG scales (w_{RL}) compared against the base ($w_{RL} = 0.0$) and GRPO ($w_{RL} = 1.0$) models.

w_{RL}	Win Rate vs. Base ($w_{RL} = 0.0$)			Win Rate vs. GRPO ($w_{RL} = 1.0$)		
	Aesthetic	ImageReward	PickScore	Aesthetic	ImageReward	PickScore
1.0	88.67	82.71	97.61	-	-	-
1.2	89.60	83.59	98.14	57.96	53.42	60.25
1.4	91.31	82.71	97.80	69.29	54.44	74.95
1.8	92.19	80.76	97.71	75.24	54.25	76.27
2.2	92.72	78.91	97.07	77.39	53.56	73.68
2.4	92.87	77.54	96.78	79.10	51.66	72.80

Table 14. Mean scores for **SDXL-base** model series on human preference metrics at various RLG scales (w_{RL}). The scale $w_{\text{RL}} = 0.0$ corresponds to the original **SDXL-base** base model, while $w_{\text{RL}} = 1.0$ represents the model after SPO finetuning, named **SPO-SDXL.4k-p.10ep**.

w_{RL}	Aesthetic Score (\uparrow)	ImageReward (\uparrow)	PickScore (\uparrow)
0.0	6.10	0.72	21.66
1.0	6.42	1.12	22.69
1.2	6.45	1.13	22.71
1.4	6.48	1.14	22.71

w_{RL}	Win Rate vs. Base ($w_{\text{RL}} = 0.0$)			Win Rate vs. SPO ($w_{\text{RL}} = 1.0$)		
	Aesthetic	ImageReward	PickScore	Aesthetic	ImageReward	PickScore
1.0	82.13	81.98	92.38	-	-	-
1.2	83.06	81.98	92.19	59.81	54.10	54.64
1.4	83.64	81.20	91.46	62.99	54.15	54.35

Table 15. Win rates (%) for **SDXL-base** model series at various RLG scales (w_{RL}) compared against the base ($w_{\text{RL}} = 0.0$) and SPO ($w_{\text{RL}} = 1.0$) models.

Table 16. Mean scores for **SD1.5** model series on human preference metrics at various RLG scales (w_{RL}). The scale $w_{\text{RL}} = 0.0$ corresponds to the original **SD1.5** base model, while $w_{\text{RL}} = 1.0$ represents the model after DPO finetuning, named **dpo-sd1.5-text2image-v1**.

w_{RL}	Aesthetic Score (\uparrow)	ImageReward (\uparrow)	PickScore (\uparrow)
0.0 (Base)	5.51	-0.02	20.03
1.0 (DPO)	5.61	0.20	20.39
1.2	5.62	0.22	20.42
1.4	5.62	0.25	20.46
1.6	5.64	0.26	20.51
1.8	5.63	0.29	20.51
2.0	5.64	0.31	20.54
2.2	5.64	0.31	20.55
2.4	5.64	0.32	20.56

Table 17. Win rates (%) for **SD1.5** at various RLG scales (w_{RL}) compared against the base ($w_{\text{RL}} = 0.0$) and standard DPO ($w_{\text{RL}} = 1.0$) models.

w_{RL}	Win Rate vs. Base ($w_{\text{RL}} = 0.0$)			Win Rate vs. DPO ($w_{\text{RL}} = 1.0$)		
	Aesthetic	ImageReward	PickScore	Aesthetic	ImageReward	PickScore
1.0	61.52	63.33	72.66	-	-	-
1.4	64.11	63.82	75.34	53.56	53.47	57.86
1.8	63.92	66.36	76.27	55.71	56.69	61.43
2.2	64.55	66.46	76.22	56.40	56.64	61.72
2.4	64.21	66.31	76.61	56.25	57.08	61.23

1320	N. Selected Images Generated
1321	
1322	N.1. Aesthetic Images Generated
1323	
1324	N.2. OCR Images Generated
1325	
1326	N.3. Compressibility and Incompressibility Images Generated
1327	
1328	N.4. GenEval Images Generated
1329	
1330	
1331	
1332	
1333	
1334	
1335	
1336	
1337	
1338	
1339	
1340	
1341	
1342	
1343	
1344	
1345	
1346	
1347	
1348	
1349	
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	

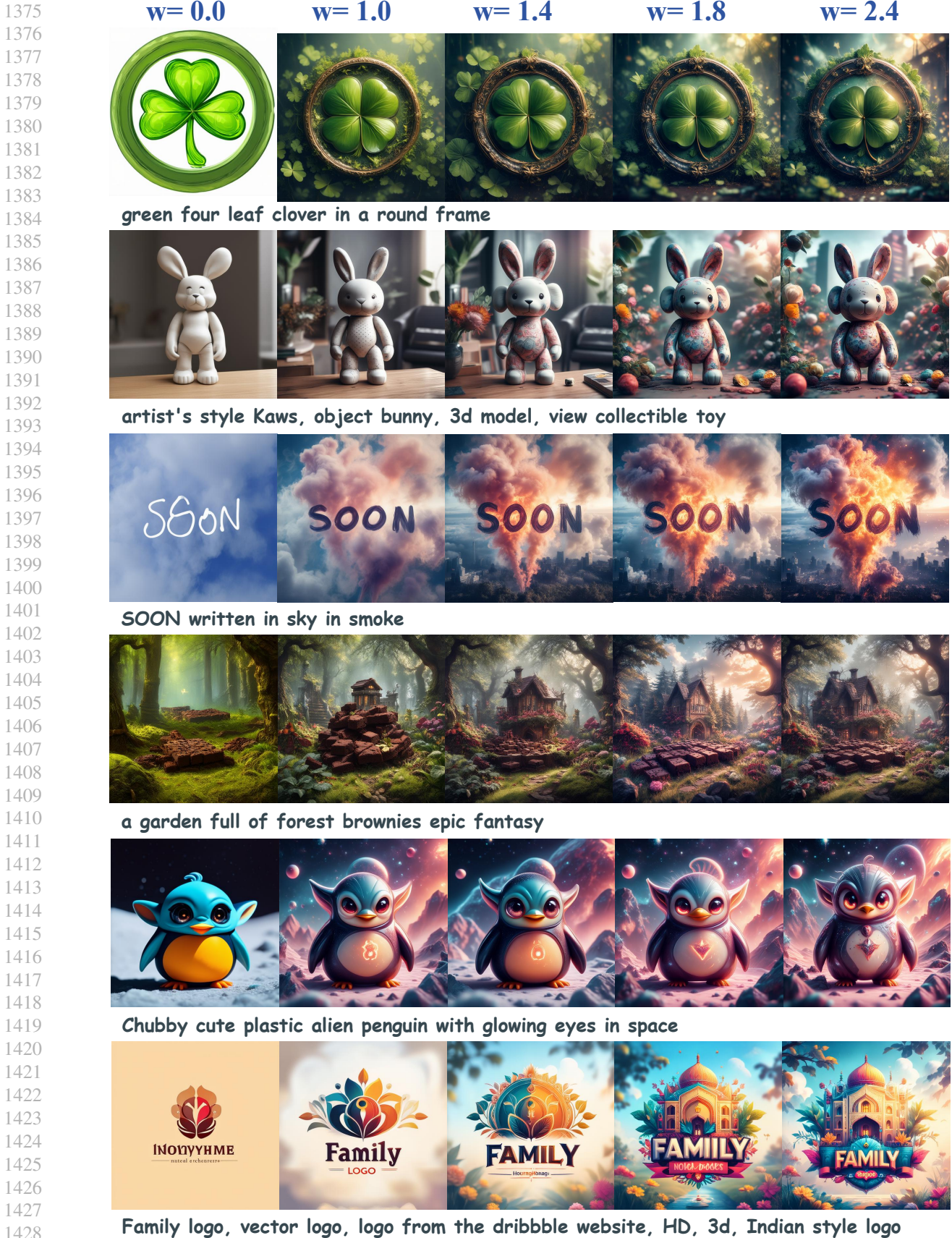


Figure 7. Selected qualitative results for the human preference task. Images are generated from SD3.5 trained with GRPO, with different RLK scales.

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
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



Illustration of a School girl and a giant steampunk robot



A painting, a beautiful portrait of a girl, oil painting style



Realistic Black and white portrait of ... a 19 year old girl , ...



anime girl eating pizza, hd, 4k, anime



A tiny mouse holding a sword



Chinese face joyful, fine art, HD, 8K

Figure 8. Selected qualitative results for the human preference task. Images are generated from SD1.5 trained with DPO, with different RLG scales.

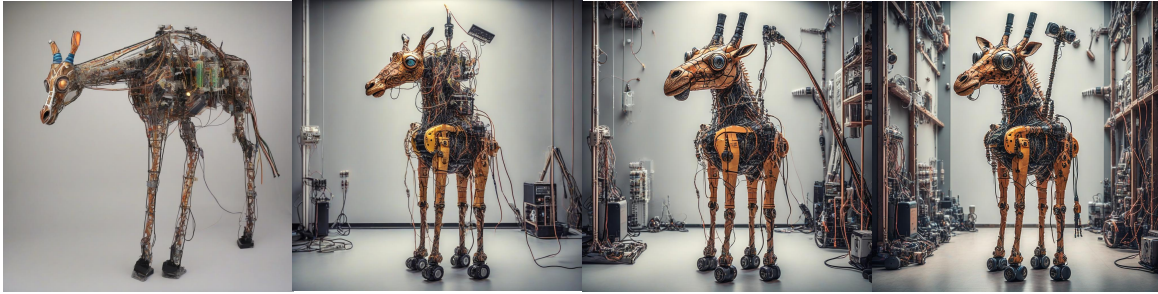
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539

w= 0.0

w= 1.0

w= 1.2

w= 1.4



mechanical giraffe, electronics, ... led instead of eyes



Marilyn Manson sticking tongue out wearing sunglasses holding a sign that says Famousa



Egirl with orange hair, gorgeous, high-quality, beautiful



The background of the cover should be the campus scenery of the primary school ...



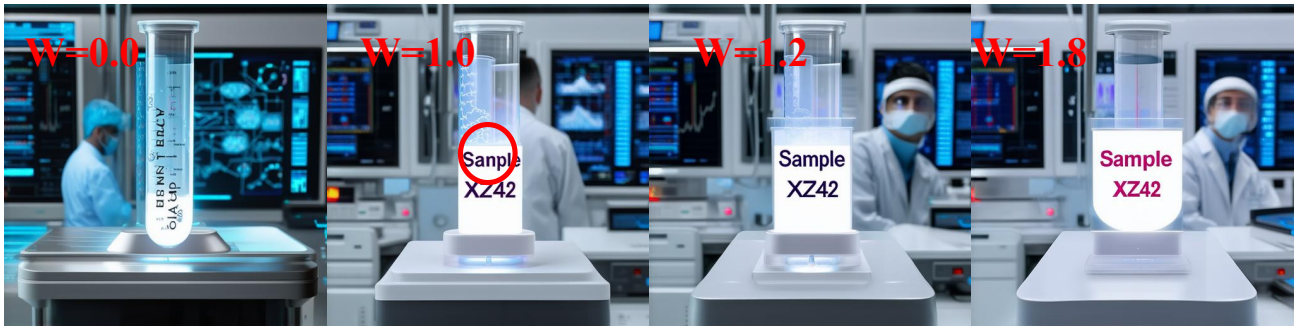
a minimalistic style logo for a game ... studies market and information design

Figure 9. Selected qualitative results for the human preference task. Images are generated from SDXL trained with SPO, with different RLG scales.

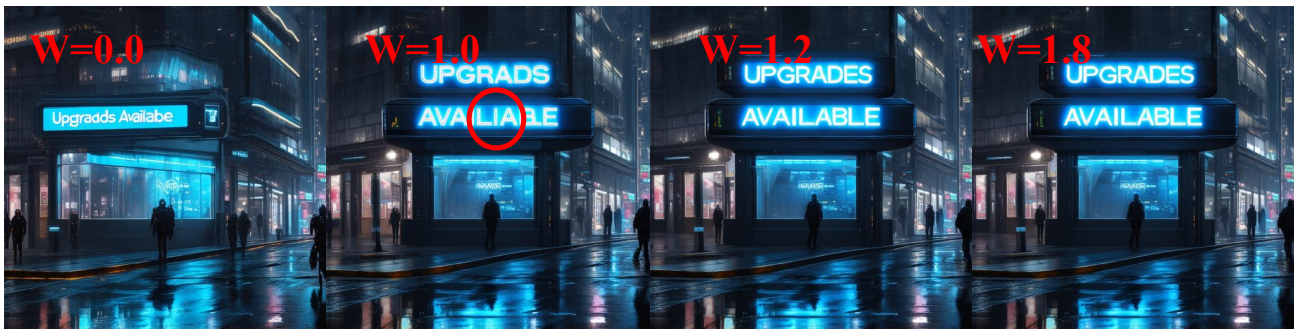
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594



" Bandit Who Stole the D "



" Sample XZ42 "



" Upgrades Available "

Figure 10. Selected qualitative results for the visual text rendering task. As can be seen, the standard RL-finetuned model ($w = 1.0$) still produces some errors in the generated text. By applying RLG with a higher guidance scale ($w > 1.0$), the model correctly renders the specified text without any loss in image quality. This illustrates how RLG effectively enhances the model's ability to adhere to precise instructions.

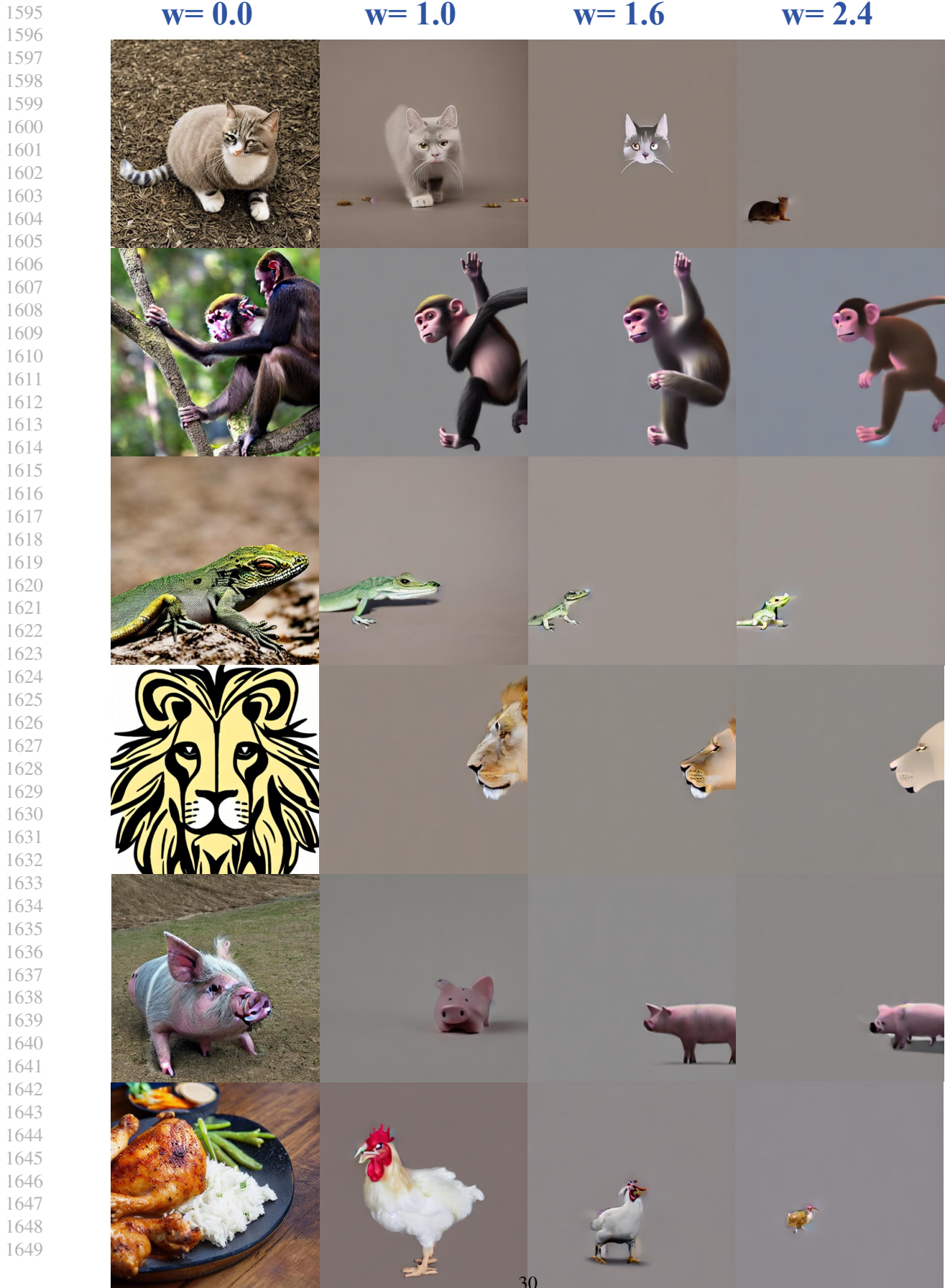


Figure 11. Selected qualitative results for the image compressibility task.

1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
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



Figure 12. Selected qualitative results for the image incompressibility task.

1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
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

w= 0.0

w= 1.0

w= 1.2

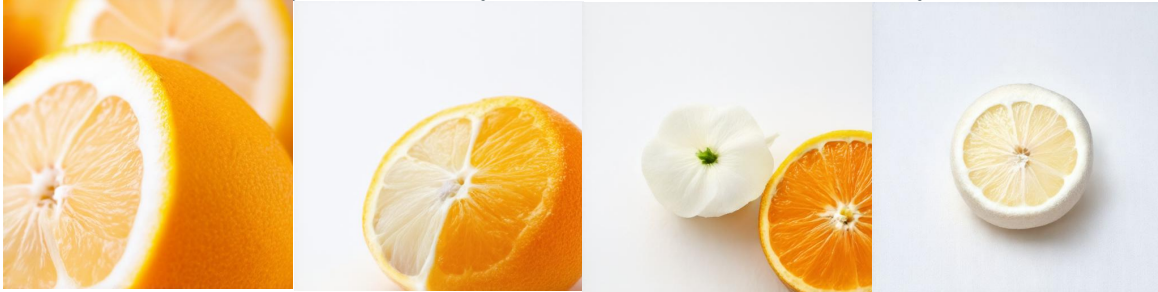
w= 1.8



a photo of a blue pizza and a yellow baseball glove



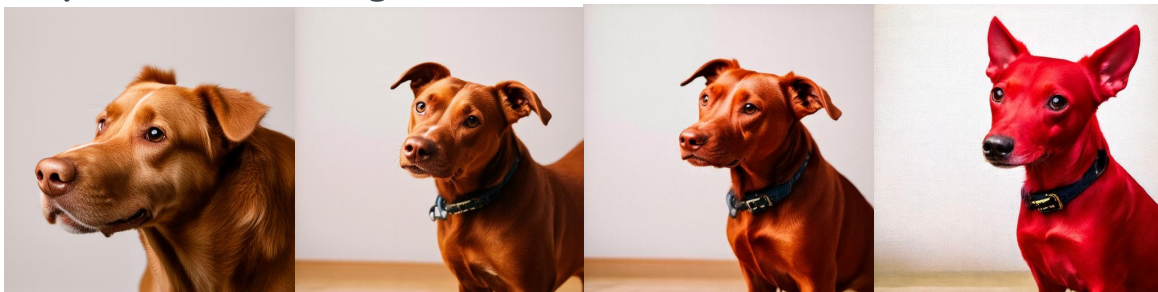
a photo of a yellow bicycle and a red motorcycle



a photo of a white orange



a photo of a red giraffe



a photo of a red dog

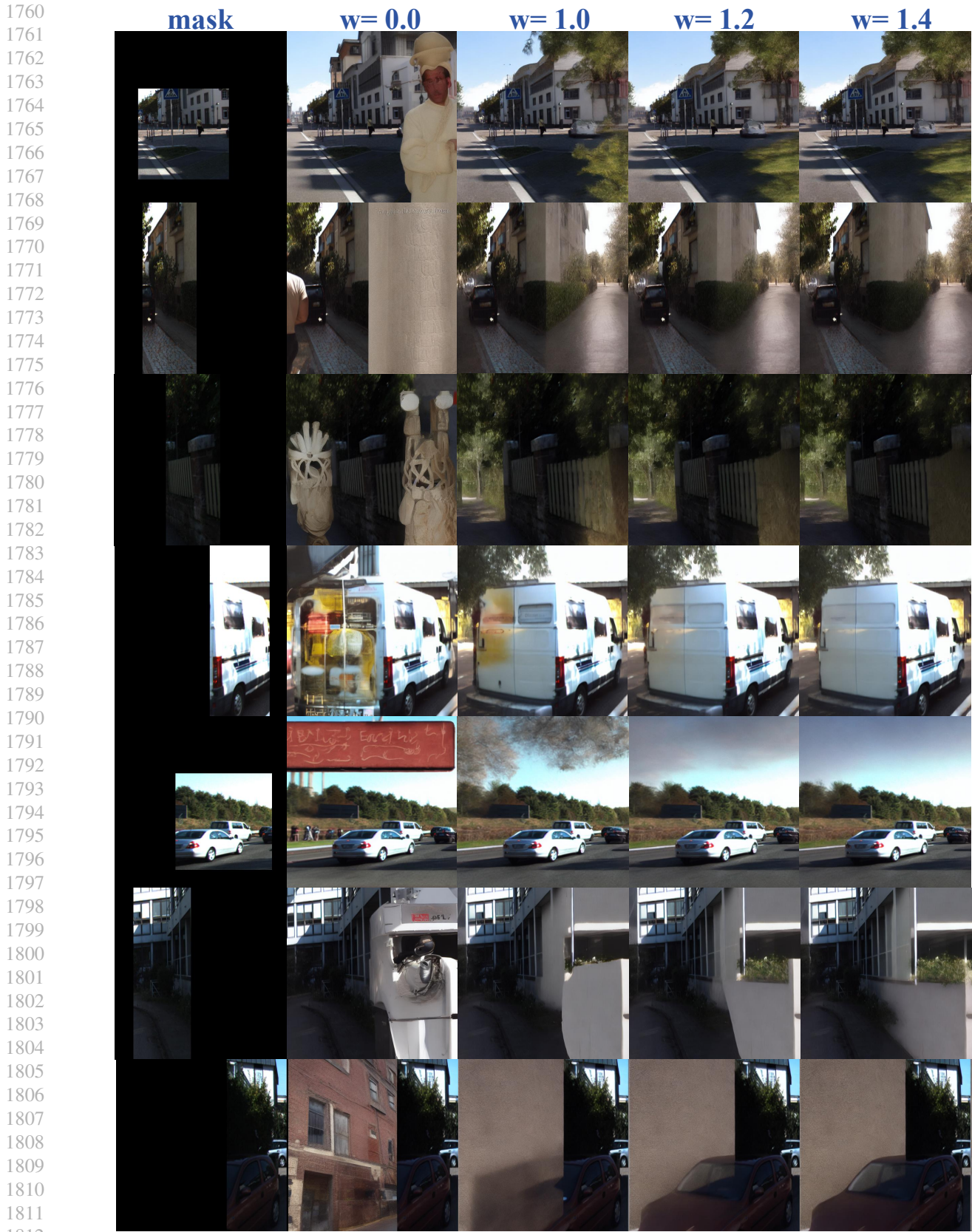


Figure 14. Selected qualitative results for the image inpainting task.