# HyperAlign: Hypernetwork for Efficient Test-Time Alignment of Diffusion Models

Xin Xie[1], Jiaxian Guo[2], Dong Gong[1†]

[1]University of New South Wales (UNSW Sydney), [2]Google Research

{xin.xie3, dong.gong}@unsw.edu.au, jeffguo@google.com

Figure 1. Sample images generated by our method based on FLUX backbone. The generated images not only achieve a high alignment with text prompt and human preferences, but also exhibit visually attractive and stunning aesthetics.

## Abstract

*Diffusion models achieve state-of-the-art performance but often fail to generate outputs that align with human preferences and intentions, resulting in images with poor aesthetic quality and semantic inconsistencies. Existing alignment methods present a difficult trade-off: fine-tuning approaches suffer from loss of diversity with reward over-optimization, while test-time scaling methods introduce significant computational overhead and tend to under-optimize. To address these limitations, we propose HyperAlign, a novel framework that trains a hypernetwork for efficient and effective test-time alignment. Instead of modifying latent states, HyperAlign dynamically generates low-rank adaptation weights to modulate the diffusion model's generation operators. This allows the denoising trajectory to be adaptively adjusted based on input latents, timesteps and prompts for reward-conditioned alignment. We introduce multiple variants of HyperAlign that differ in how frequently the hypernetwork is applied, balancing between performance and efficiency. Furthermore, we optimize the hypernetwork using a reward score objective regularized with preference data to reduce reward hacking. We evaluate HyperAlign on multiple extended generative paradigms, including Stable Diffusion and FLUX. It significantly outperforms existing fine-tuning and test-time scaling baselines in enhancing semantic consistency and visual appeal. The project page: hyperalign.github.io.*

## 1. Introduction

Diffusion models learn score function [46] to gradually transform a random noise into a structured output, offering state-of-the-art performance in Text-to-Image (T2I) generation [26, 36, 40]. However, these models are typically trained on a large set of pre-collected datasets (*e.g.*, web images) that may not accurately represent target conditional distribution aligned with human intention and preferences. Consequently, the generated images often misrepresent users' textual instructions and fail to reflect their aesthetic preferences. Despite classifier-based or free advances [9, 21], models only improve prompt controllability but still struggle to reflect fine-grained human preferences.

---

[†]D. Gong is the corresponding author.

1

These challenges highlight the necessity of diffusion model *alignment* [31] to bridge the gap between the generated images and human preferences, enhancing semantic consistency with textual prompts and visual appeals.

Alignment of diffusion models is generally approached through fine-tuning and test-time scaling. Fine-tuning alignment approaches, including Reinforcement Learning (RL) [27, 28, 30, 52, 59] and direct backpropagation [37, 44], optimize target rewards based on explicit reward signals or implicit feedback from preference datasets [25, 54]. While training-based alignment methods effectively reshape the distribution to close to a desired target distribution, their overall performance remains constrained. Since the generation requirements and inputs to the model (*i.e.*, users' input prompts and the sampled initial random noise) vary across use cases, fine-tuning through a set of model parameters may not account for every combination of the desires in Fig. 2. As a result, it often suffers from the reward over-optimization problem (*e.g.*, reward hacking), leading to a severe loss of diversity and creativity.

On the other hand, test-time scaling methods perform input-specific computing for the alignment goal during inference, through gradient-based approaches [24, 49, 56, 63] or sampling-based approaches [33, 39]. This allows the alignment process to be tailored to each specific requirement, *i.e.*, executing necessary computation to dynamically adjust the denoising trajectory and align the generated outputs with query-specific objectives. However, these test-time scaling methods incur the additional computational overhead introduced by gradient calculation and repeated forward sampling. Meanwhile, they also suffer from the reward under-optimization, failing to effectively optimize target rewards since the externally injected test-time prior is isolated from the diffusion models' training dynamics.

To address these limitations, we propose to train a hypernetwork to achieve efficient and effective test-time alignment of diffusion models, termed as **HyperAlign**. Starting from the sampled noise, T2I diffusion models generate results along a denoising trajectory, reflecting the input prompts. However, the results often exhibit poor aesthetic quality misaligned with human preferences and semantic inconsistency with input prompts. Leveraging the score-based nature of diffusion processes, we formulate the alignment task as aligning the trajectories. Unlike directly modifying the latent states through gradient-based guidance, we aim to adjust step-wise generation operators for reward-conditioned alignment, namely modifying the network weights of given generative models. We design a hypernetwork that inputs latents, timesteps and user prompts and generates test-time dynamic modulation parameters to adapt the generation trajectory accordingly. Considering the prohibitive cost of generating full model parameters, the hypernetwork is designed to produce low-rank adaptation
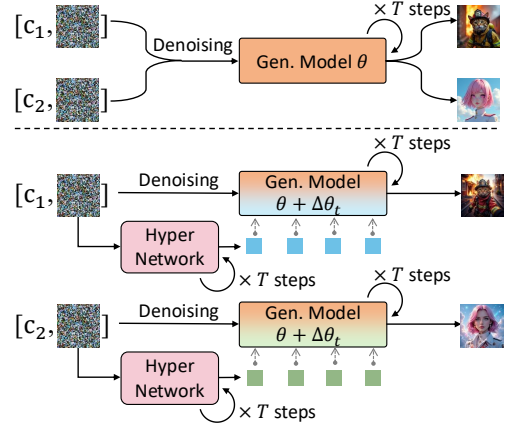


Figure 2. Task-specific test-time alignment of HyperAlign. Compared to the original generative model, HyperAlign adapts the model's behavior to each combination of prompt and temporal states, producing aligned and visually appealing results

weights (LoRA). To enhance efficiency, we introduce three weights generation strategies: step-wise generation for fine-grained alignment, generation at the starting point for minimal computation, and piece-wise generation that updates adapters only at key timesteps to balance performance and efficiency. To optimize the hypernetwork, we use reward score as the training objective. To mitigate reward hacking, we regularize generated trajectories with preference data, preventing the model from overfitting to proxy scores while maintaining fidelity to genuine human preferences. We implement the proposed method on different generative paradigms (diffusion models [40] and rectified flows [26]).

The main contributions are summarized as:

- We propose HyperAlign, a hypernetwork that adaptively adjusts denoising operations for efficient and effective test-time alignment of diffusion models, ensuring that generated images better reflect user-intended semantics in text and appealing visual quality.
- We design different strategies for adaptive weight generation, enabling efficient and flexible alignment. Apart from reward score as training objective, we introduce a preference regularization term to prevent reward hacking.
- We evaluate the performance of the proposed method with different generative models, *e.g.*, SD V1.5 and FLUX. HyperAlign outperforms different baseline models and other state-of-the-art fune-tuning and test-time scaling methods significantly on different metrics, demonstrating effectiveness and superiority.

## 2. Related Work

### 2.1. Fine-tuning Diffusion Model Alignment

Diffusion models [26, 36, 40] exhibit remarkable generative performance, yet suffer from misalignment with hu-

man expectations. Early works [7, 10, 37, 55, 57] directly learn preferences from reward models, but are constrained by unstable long-trajectory gradients. Therefore, SRPO [44] employs a noise prior to predict clear data and yields accurate reward gradient for each step. Alternatively, DDPO [5] and DPOK [14] integrate RL to optimize the score function [46] through policy gradient updates. Subsequently, D3PO [60] and Diffusion-DPO [52] first introduce offline Direct Preference Optimization (DPO), modeling human preferences from win–lose paired data. SPO [30] and LPO [64] extend step-wise preference alignment by training timestep-aware reward models. Diffusion-KTO [28] adopt human utility maximization to reduce reliance on offline paired data. Recently, Flow-GRPO [32] and Dance-GRPO [59] pioneer the integration of group-wise policy optimization paradigm [43] for improved diffusion model alignment. TempFlow-GRPO [19] designs temporal-aware credit assignment for intermediate-step advantage estimation and Pref-GRPO [53] replaces pointwise score maximization objective with pairwise preference fitting. To improve the training efficiency, MixGRPO [27] adopts a mixed ODE–SDE paradigm and BranchGRPO [29] restructures rollout process into a branching tree. Despite substantial advances in aligning diffusion models, discrepancies with human preferences and considerable computational burdens continue to pose challenges.

## 2.2. Test-time Computing for Diffusion Models

The goal of test-time scaling is to spend additional compute during inference to obtain more desirable generations. One naive scaling law in diffusion model sampling is to increase the number of denoising steps [3, 58], enabling marginal improvements. Beyond this, there are two mainstream test-time techniques: One is sampling-based strategies, relying the reward models to evaluate multiple noise candidates and select more favorable denoising trajectory, such as Best-of-N search [33], evolutionary search [18], $\varepsilon$-greedy search [39], etc. The other one is gradient-based approaches, built on the score-based formulation of diffusion models [46]. These methods use the differentiable reward functions to iteratively refine noise [12, 16, 49], prompt embeddings [8] or latents [4, 24, 51, 56, 61, 63] through gradient descent. However, these test-time scaling methods suffer from inaccurate guidance and limited practicality. Optimizing a single image on DiT-based models takes over minutes.

## 2.3. Hypernetworks

Ha *et al.* [17] proposed hypernetworks that predict weights of a primary network, showing notable success in language modeling [6, 23]. For vision tasks, hypernetworks are applied across various domains, including segmentation [35], image editing [2], continue learning [50], 3D modeling [48], personalization [20, 41] and initial noise prediction for diffusion model [1, 13], among others.

# 3. Problem Setup: Diffusion Model Alignment

## 3.1. Preliminary on Score-based Generative Models

Diffusion models [22] capture a data distribution by learning to reverse a gradual noising process of applied to clean data. Given a data distribution $p_{\text{data}}(\mathbf{x})$, the forward process of a diffusion model [22, 45, 46] progressively perturbs a clean sample $\mathbf{x}_0 \sim p_{\text{data}}(\mathbf{x})$ with Gaussian noise toward a Gaussian noise, following a stochastic differential equation (SDE) under certain conditions:

$$\mathrm{d}\mathbf{x}_t = \mathbf{f}(\mathbf{x}_t)\,\mathrm{d}t + g_t\,\mathrm{d}\mathbf{w}, \tag{1}$$

where $t \in [0, T]$, $\mathbf{w}$ is a standard Wiener process, $\mathbf{f}(\mathbf{x}_t)$ and $g_t$ represent the drift and diffusion coefficients, respectively [46].

By running the above process backwards starting from $\mathbf{x}_T \sim \mathcal{N}(0, \mathbf{I})$, we obtain a data generation process through the reverse SDE:

$$\mathrm{d}\mathbf{x}_t = \left[\mathbf{f}(\mathbf{x}_t) - g_t^2 \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)\right]\mathrm{d}t + g_t\,\mathrm{d}\mathbf{w}, \tag{2}$$

where $p_t(\mathbf{x}_t)$ denotes the marginal distribution of $\mathbf{x}_t$ at time $t$. The score function $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)$ can be estimated by training a model $\mathbf{s}_\theta(\mathbf{x}_t, t)$ [46, 47]:

$$\min_\theta \mathbb{E}_{t,\mathbf{x}_0,\mathbf{x}_t} \left\{ \lambda(t) \left\| \mathbf{s}_\theta(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{x}_0) \right\|_2^2 \right\}, \tag{3}$$

where $\lambda(t)$ is the weight function, $\mathbf{x}_0 \sim p_{\text{data}}(\mathbf{x})$, $p_t(\mathbf{x}_t|\mathbf{x}_0)$ is a transition density in Gaussian, and $\mathbf{x}_t \sim p_t(\mathbf{x}_t|\mathbf{x}_0)$. The approximated $\mathbf{s}_\theta(\cdot)$ defines a learned distribution $p_\theta(\cdot)$.

The score-based model unifies the formulations of diffusion models [22, 45] and flow matching models [11, 32], where the sample trajectories of $\mathbf{x}_t$ are generated through a stochastic or ordinary differential equation (SDE or ODE) [46]. For clarity and simplicity, we focus on diffusion models in the following presentation without loss of generality. Under this unified formulation, we can naturally generalize our analyses and approach to both diffusion and flow-matching models. More details in supplementary material.

## 3.2. Aligning Diffusion Model with Reward

**Conditional diffusion models and score functions.** We consider conditional diffusion models that learn a distribution $p_\theta(\mathbf{x}|\mathbf{c})$ with $\mathbf{c}$ denotes the conditioning variable. It is trained to generate samples through a reverse diffusion process via denoising a sampled noise $\mathbf{x}_T$ under the control conditioning on $\mathbf{c}$. For image generation, $\mathbf{c}$ is the input prompts indicating user's instruction for the generated contents. We resort to discrete score-based model with the variance-preserving setting [22, 45] for better discussion, and its sampling formula is:

$$\mathbf{x}_{t-1} = (1 + \frac{1}{2}\beta_t)\mathbf{x}_t + \beta_t \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|c) + \sqrt{\beta_t}\,\boldsymbol{\epsilon}, \tag{4}$$
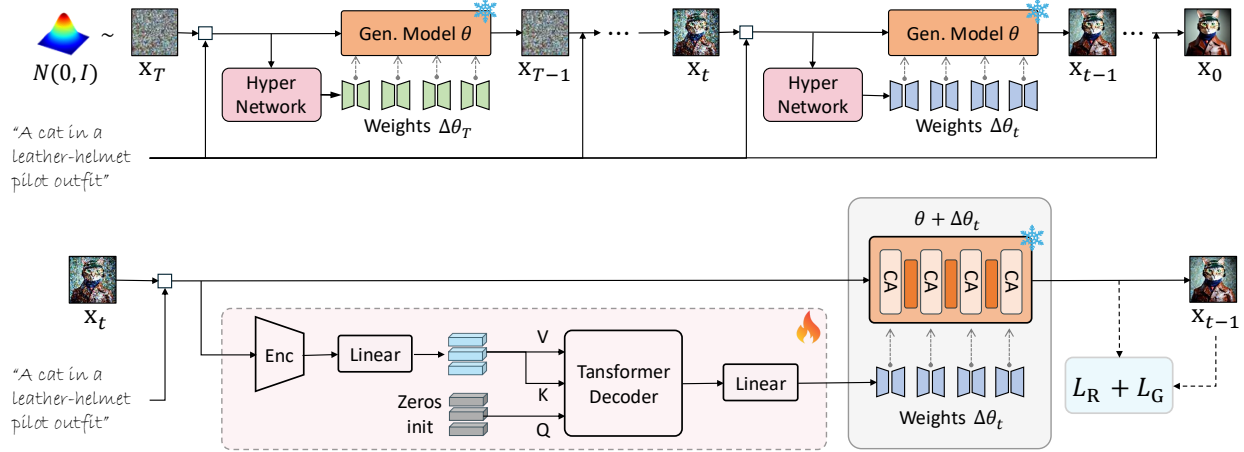
Figure 3. The framework of HyperAlign. Given a user prompt, the hypernetwork produces step-wise modulation weights $\Delta\theta_t$ that are injected into the generative model to steer the denoising trajectory (top). During training (bottom), the hypernetwork is optimized using the reward loss and the preference-regularization loss, enabling it to produce input-specific adjustments.

where $\mathbf{x}_t \sim \mathcal{N}(\sqrt{\bar{\alpha}_t}\mathbf{x}_0, \sigma_t^2\mathbf{I})$, $\bar{\alpha}_t = \prod_{i=1}^{t}(1 - \beta_i)$, $\sigma_t = \sqrt{1 - \bar{\alpha}_t}$, and $\beta_t$ is a linearly increasing noise scheduler. This iterative denoising process forms a trajectory $\{\mathbf{x}_t\}_{t=T}^{0}$ in the latent space, gradually transforming the noise $\mathbf{x}_T$ into a clean sample $\mathbf{x}_0$ reflecting the input prompt $\mathbf{c}$.

**Diffusion model alignment with reward.** While existing T2I models demonstrate strong generative capabilities, the results frequently fall short of user expectations, showing poor visual appeal and semantic inconsistency with the input prompts. This limitation arises because the score functions are learned from large-scale uncurated datasets, which diverge from distribution of human preferences. To bridge this gap, diffusion model alignment is introduced to enhance the consistency between the generated images and human user preferences.

Relying on human preference data [25, 34, 54], we can obtain a reward model $R(\mathbf{x})$ that captures human preference, *e.g.*, aesthetic preference [42]. Through connecting with the condition $\mathbf{c}$, the reward model can be formulated as $R(\mathbf{x}, \mathbf{c})$, which can be assumed to partially capture the consistency between $\mathbf{x}$ and $\mathbf{c}$ together with the visual aesthetic preference. It can be explicitly learned from preference data or implicitly modeled directly using data. Given a learned $p_\theta(\mathbf{x}|\mathbf{c})$ and a reward model, diffusion model alignment can be formulated as solving for a new distribution:

$$p_{\theta,R}(\mathbf{x}|\mathbf{c}) = \frac{1}{\mathcal{Z}}p_\theta(\mathbf{x}|\mathbf{c})\exp(\frac{R(\mathbf{x}, \mathbf{c})}{\gamma}), \quad (5)$$

where $\gamma$ is the KL regularization coefficient controlling the balance between reward maximization and consistency with the base model. Prevalent training-based alignment methods optimize the target rewards through RL [27, 28, 30, 52, 59] and direct backpropagation [37, 44]. Although effective, these approaches often incur substantial computational overhead and risk over-optimization, leading to degraded

generation diversity. In contrast, test-time scaling methods achieve alignment goal by using guidance to modify the temporal states. Since the generative distribution is manifested as the trajectory of $\mathbf{x}_t$ in the sampling process, test-time alignment can be regarded as steering this trajectory to better match the desired conditional distribution $p_{\theta,R}(\mathbf{x}|\mathbf{c})$.

## 4. Methodology: HyperAlign

In this work, we aim to learn a hypernetwork for efficient and effective test-time alignment of diffusion models, termed as HyperAlign.

### 4.1. Test-time Alignment with Diffusion Guidance

As discuss in Sec. 3.2, test-time diffusion alignment methods adjust the generative trajectory to better satisfy alignment objectives. Existing test-time computing strategies can be broadly categorized into noise sampling-based and gradient-based diffusion guidance. Noise sampling methods attempt to identify favorable noise candidates based on reward feedback. However, exploring the vast high-dimensional noise space is computationally expensive and hard to converge, leading to inefficiency and under-optimized outcomes. In contrast, gradient-based diffusion guidance directly compute the gradient from specific objectives and uses them to steer the denoising trajectory by modifying the temporal states.

To effectively align the diffusion model through directly injecting the guidance from reward, we aim to train a hypernetwork that generates prompt-specific and state-aware adjustments at each denoising step. This design maintains computational efficiency by amortizing the costly test-time optimization into a compact and learnable modeling process during finetuning.

Before introducing the proposed method, we first ana-
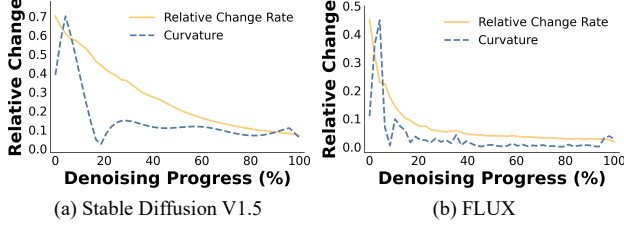
(a) Stable Diffusion V1.5    (b) FLUX

Figure 4. The prompt-invariant temporal dynamics of one-step predicted data. Average over 1000 prompts.

lyze diffusion guidance approaches that achieve alignment by leveraging generative gradients to steer the denoising trajectory. Based on Bayes' rule, we can derive an approximate expression of $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|c) \approx \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) + \nabla_{\mathbf{x}_t} R(\mathbf{x}_{0|t}, c)$, where the first term corresponds to the unconditional score and does not require extra optimization. Thus, we focus on the second term, which injects reward gradient into the denoising process:

$$\nabla_{\mathbf{x}_t} R(\mathbf{x}_{0|t}, c) = \frac{1}{\sqrt{\bar{\alpha}_t}} \cdot \frac{\partial R}{\partial \mathbf{x}_{0|t}} \cdot \left( \mathbf{I} - \sqrt{1 - \bar{\alpha}_t} \cdot \frac{\partial \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)}{\partial \mathbf{x}_t} \right), \quad (6)$$

where the reward function is actually applied on the decoded image domain through the decoder. For simplicity of discussion, we omit the decoder notation. By substituting Eq. (6) into Eq. (4), we observe that the guidance-based methods achieve alignment by injecting reward-aware diffusion dynamics into $\mathbf{x}_{t-1}$, which essentially changes the transition path from $\mathbf{x}_t$ to $\mathbf{x}_{t-1}$.

### 4.2. HyperNetwork for Test-time Alignment

As discussed in Sec. 4.1, gradient-guidance methods perform test-time alignment by directly modifying temporal states using reward-derived scores to adjust the denoising trajectory. However, backpropagating gradients from the reward model to the generator incurs substantial computational overhead, slows inference, and remains disconnected from the generator's training process.

To mitigate these issues while retaining task-specific modeling benefits, we train a hypernetwork that efficiently steers the generation trajectory according to the task, input, and current generation states. Its test-time alignment capability is learned during training by injecting reward-based guidance into the hypernetwork. Different from fine-tuning based alignments accommodate all combinations of user intentions using a fixed set of parameters, our methods is prompt-specific and state-aware, dynamically generating adaptive modulation parameters at each denoising step to align the generation trajectory.

**Hypernetwork as a dynamic LoRA predictor.** We aim to learn a hypernetwork that takes $\mathbf{x}_t$ and $\mathbf{c}$ as input and outputs adjustments for each step of the generative process. A naive approach would be to learn an alignment score as a substitute for Eq. (6), but this requires a formulation akin to the original generative score and thus incurs high complexity. Instead, we design the hypernetwork to directly adjust the score $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t \mid c)$, corresponding to the network parameters $\theta$ in the original generative model, through producing a lightweight low-rank adapter (LoRA) for $\theta$.

We divide hypernetwork architecture into two main components: *perception encoder* and *transformer decoder*, as shown in Fig. 3. Concretely, the inputs temporal latent $\mathbf{x}_t$, timestep $t$ and prompt $\mathbf{c}$ are first passed into perception encoder, which consists of downsampling blocks from the pretrained U-Net of the generative models. The pretrained U-Net carries rich diffusion priors, making it a natural encoder to capture semantic representations across diverse input combinations. The encoded features are then projected through a linear layer and passed to a transformer decoder, where we use zero-initialized tokens to generate query (Q) and use the encoded features to generate the key (K) and the value (V). The transformer decoder integrates temporal and semantic information via cross-attention, and a subsequent linear layer maps the decoded features into LoRA weights:

$$\Delta\theta_t = h_\psi(\mathbf{x}_t, \mathbf{c}, t), \quad (7)$$

where $\psi$ denotes the parameters of hypernetwork $h_\psi$. Temporally, integrating the generated LoRA weights into the original model parameters yields a input-and-step-specific score function $\mathbf{s}_{\theta + \Delta\theta_t}$ (with an abuse of notation $+$), thereby modifying the underlying denoising trajectory.

**Efficient HyperAlign.** By default, the hypernetwork design in Eq. (7) can be applied at all generation steps adaptively starting from initial step $T$ (termed as HyperAlign-S). We further develop two variants for balancing inference efficiency. (1) HyperAlign-I is trained to only predict the LoRA weights once at the starting point, $\Delta\theta_T = h_\psi(\mathbf{x}_T, \mathbf{c}, T)$, and used for all steps. (2) A piece-wise variants, HyperAlign-P, produces new weights at several key timesteps, where all timesteps within the same segment share the same LoRA weights. We compute the the relative $\ell_1$ distance of one-step predicted latents, shown in Fig. 4, while a small value indicates that adjacent latents are similar to each other. The observations support that similar latent states can be grouped into a single segment and share same LoRA weights, aligning with the diffusion behavior across different denoising stages. We compute the curvature rate to identify $M$ keypoints that exhibit greater influence on the trajectory. The hypernetwork is trained to regenerate LoRA weights only at these keysteps to adaptively modulate the diffusion process with less computations than HyperAlign-S, balancing between efficiency and performance.

### 4.3. HyperAlign Training

To optimize the hypernetwork, we can use the reward score as training objective. By maximizing the reward signals,

5

| FLUX | BoN | ε-greedy | FreeDoM | DyMO | DanceGRPO | MixGRPO | SRPO | HyperAlign |
|------|-----|----------|---------|------|-----------|---------|------|-----------|



A still of Doraemon from "Shaun the Sheep" by Aardman Animation

A 3D rendering of anime schoolgirls with a sad expression underwater, surrounded by dramatic lighting

A kitten with a panda coloring eating bamboo

A photo of a cow left of a stop sign
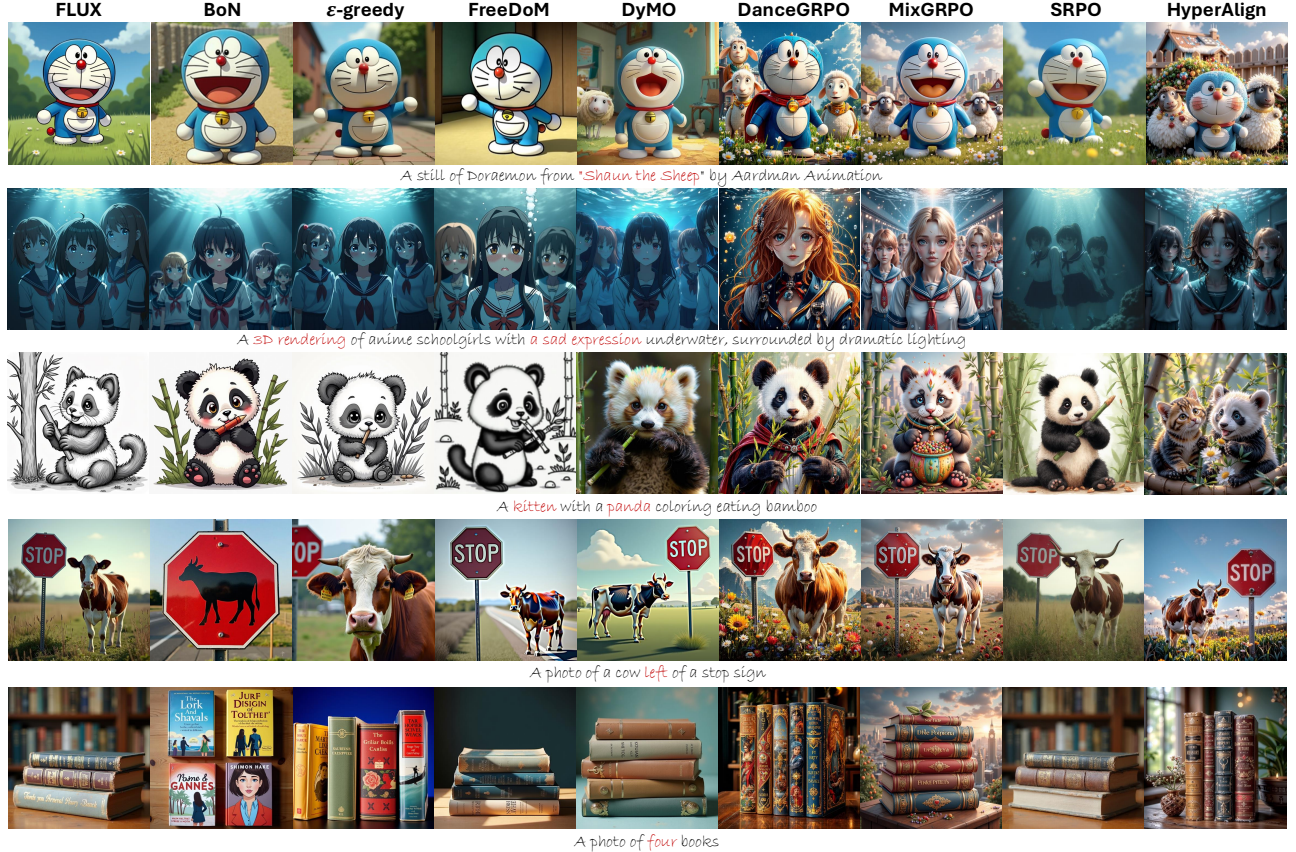
A photo of four books

Figure 5. Qualitative comparison examples based on FLUX backbones.

the model is encouraged to generate intermediate predictions with higher conditional likelihood, thereby aligning the latent trajectory with the true conditional distribution:

$$\mathcal{L}_{\mathrm{R}} = -\mathbb{E}_{p(\mathbf{x}_t)}\left[R(\mathbf{x}_{0|t}, \mathbf{c})\right]. \tag{8}$$

**Regularization on reward optimization.** While maximizing reward objective drives the model to produce high-reward, condition-aligned latent states, it also exposes two key challenges: (1) inaccurate reward signals due to the blurriness of one-step predictions in early denoising stages, and (2) the risk of over-optimization, where aggressive reward maximization leads to reward hacking or degraded visual fidelity. To mitigate these issues, we incorporate a regularization loss to constrain the alignment process and preserve generation quality:

$$\mathcal{L}_{\mathrm{G}} = \mathbb{E}_{\mathbf{x}_0, \mathbf{x}_t|\mathbf{x}_0}\left[\eta_t \big\|\nabla_{\mathbf{x}_t}\log p_t(\mathbf{x}_t|\mathbf{x}_0) - \nabla_{\mathbf{x}_t}\log q(\mathbf{x}_t|\mathbf{x}_0)\big\|_2^2\right], \tag{9}$$

where $\eta_t$ denotes the hyperparameter, $\mathbf{x}_0$ sampled from preferred data $q(\mathbf{x}_0)$ and $\mathbf{x}_t \sim q(\mathbf{x}_t|\mathbf{x}_0)$. We encourages the learned denoising conditional score to match the score in preferred data, regularizing reward hacking.

The final learning objective for the hypernetwork optimization can be described as follows:

$$\psi^* = \arg\min_{\psi}\left\{\mathcal{L}_{\mathrm{R}} + \mathcal{L}_{\mathrm{G}}\right\}. \tag{10}$$

Our method is not limited to diffusion models, As mentioned, HyperAlign is not limited to diffusion models and is also compatible with flow-matching models (*e.g.*, FLUX in experiments). More details are in supplementary material.

## 5. Experiments

In this section, we conduct comprehensive experimental evaluations to verify the effectiveness and efficiency of our method. We flexibly apply it across various generative paradigms and validate its performance through comparisons with existing state-of-the-art approaches. Moreover, ablation studies are performed to substantiate the contribution of each component in our designs.

### 5.1. Experimental Setting

**Implementation Details.** We employ SD V1.5 [40] and FLUX [26] as base models, paired HPSv2 [54] as the reward model. The preferred data used for regularization loss originates from Pick-a-Pic [25] and HPD [54]. All our experiments uses four NVIDIA H100 GPUs.
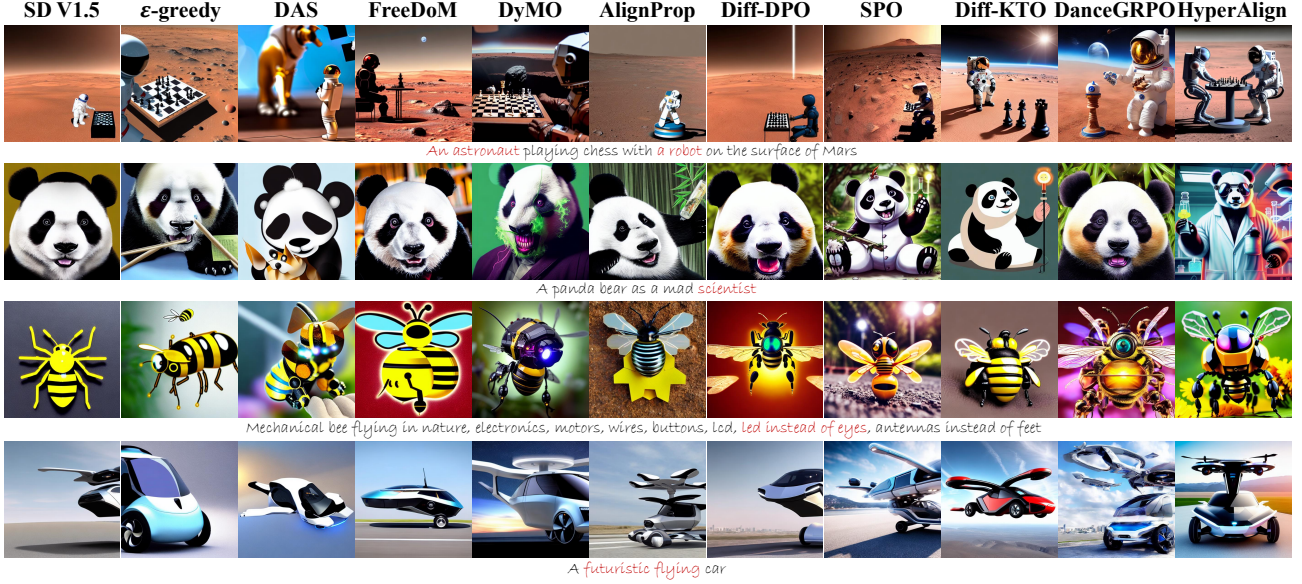
6

Figure 6. Qualitative comparison based on SD V1.5 backbones.

**Datasets and Metrics.** We evaluate our method on four datasets: 1K prompts from Pick-a-Pic [25], 2K from GenEval [15], 500 from HPD [54], and 1K from Partiprompt [62]. We choose six AI feedback models to assess the image quality: PickScore [25] and ImageReward (IR) [57] for general human preference, HPSv2 [54], CLIP [38] and GenEval Scorer [15] for prompt alignment, Aesthetic Predictor [42] for visual appeal. All test images are produced with 50 denoising steps for fair comparison, where the CFG scale is set to 7.5 for SD V1.5-based and 3.5 for FLUX-based methods. *For all metrics, higher values indicate better performance.*

### 5.2. Comparison with Existing Methods

For comprehensive assessment, we compare our method with training-based and test-time scaling models. The former covers RL methods with direct reward backpropagation [37, 44] and different policy optimization paradigms including DPO [30, 52], KTO [28] and GRPO [59]. The latter comprises reward gradient-based guidance [24, 56, 63] and noise candidate search strategies [33, 39]. For noise sampling methods, we follow the original configuration by setting the number of noise candidates to 20 for BoN [33], and using 20 local search iterations with 4 noise candidates for $\varepsilon$-greedy [39].

#### 5.2.1. Quantitative Analysis

To objectively evaluate the performance of our method, we conduct a quantitative comparison on the Pick-a-Pic dataset. For fair comparison, all alignment methods only use the HPSv2 scorer [54] as the reward model. The results are organized in Tab. 2 for SD V1.5-based backbones and Tab. 1 for FLUX-based backbones, respectively. It is observed that

Table 1. Comparison of AI feedback on FLUX-based methods.

| Method | Aes | Pick | IR | CLIP | HPS | Time |
|---|---|---|---|---|---|---|
| FLUX | 6.182 | 22.26 | 1.023 | 0.2651 | 0.3072 | 15s |
| BoN | 6.016 | 22.26 | 1.124 | 0.2668 | 0.3114 | 300s |
| $\varepsilon$-greedy | 6.166 | <u>22.30</u> | 1.077 | 0.2662 | 0.3115 | 1100s |
| FreeDoM | 5.933 | 21.67 | 0.842 | <u>0.2672</u> | 0.2970 | 240s |
| DyMO | 6.165 | 21.74 | 1.062 | **0.2714** | 0.3004 | 300s |
| DanceGRPO | 6.706 | 22.08 | 1.135 | 0.2341 | 0.3527 | 15s |
| MixGRPO | 6.760 | 22.25 | 1.240 | 0.2499 | 0.3460 | 15s |
| SRPO | 6.061 | 22.16 | 0.833 | 0.2624 | 0.2917 | 15s |
| HyperAlign-I | 6.733 | 22.18 | 1.242 | 0.2511 | 0.3529 | 16s |
| HyperAlign-P | <u>6.769</u> | 22.13 | **1.280** | 0.2506 | <u>0.3530</u> | 17s |
| HyperAlign-S | **6.853** | **22.37** | <u>1.251</u> | 0.2602 | **0.3611** | 20s |

our method effectively achieve alignment and outperform the previous methods by adjusting the generation trajectory step by step. The other two variants of our method also keep competitive performance with faster inference. By contrast, test-time methods suffer from under-optimization in preference alignment. However, DyMO [56], benefiting from its semantic consistency objective, retains relatively high text–image alignment reflected by the CLIP score. The fine-tuning alignment methods produce suboptimal results due to due to the lack of input-specific adaptability. More metric comparison results across various benchmarks are provided in the supplementary material.

#### 5.2.2. Qualitative Evaluation

We provide a visual comparison of the generated images in Fig. 6 for SD V1.5-based backbones and Fig. 5 for FLUX-based backbones. It is evident that consistently produces images with coherent layouts, semantically rich content aligned with the prompts, and superior visual and aes-
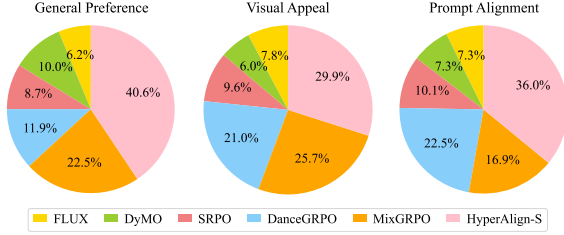
Figure 7. User study results.

Table 2. Comparison of AI feedback on SD V1.5-based methods.

| Method | Aes | Pick | IR | CLIP | HPS | Time |
|---|---|---|---|---|---|---|
| SD v1.5 | 5.443 | 20.66 | 0.1111 | 0.2745 | 0.2500 | 3s |
| $\varepsilon$-greedy | 5.509 | 21.02 | 0.4906 | 0.2821 | 0.2737 | 250s |
| DAS | 5.231 | 19.43 | -0.4524 | 0.2420 | 0.2387 | 25s |
| FreeDoM | 5.690 | 21.16 | 0.4419 | 0.2822 | 0.2779 | 120s |
| DyMO | 5.723 | 21.89 | 0.6944 | **0.2901** | 0.2944 | 162s |
| AlignProp | 5.409 | 20.58 | 0.1345 | 0.2750 | 0.2501 | 3s |
| Diffusion-DPO | 5.519 | 21.01 | 0.3201 | 0.2795 | 0.2601 | 3s |
| SPO | 5.664 | 21.15 | 0.2978 | 0.2599 | 0.2741 | 3s |
| Diffusion-KTO | 5.660 | 21.16 | 0.6970 | 0.2809 | 0.2824 | 3s |
| DanceGRPO | 5.711 | 21.15 | 0.6886 | 0.2748 | 0.2925 | 3s |
| HyperAlign-I | 5.791 | 21.94 | 0.5947 | 0.2870 | 0.2877 | 3s |
| HyperAlign-P | **5.878** | 21.89 | 0.6979 | 0.2823 | 0.2885 | 4s |
| HyperAlign-S | 5.824 | **22.01** | **0.7731** | 0.2851 | **0.2957** | 5s |

Table 3. Ablation study results.

| Method | Aes | Pick | IR | CLIP | HPS |
|---|---|---|---|---|---|
| HPSv2+Pick-a-Pic | 5.824 | 22.01 | 0.7731 | 0.2851 | 0.2957 |
| HPSv2+HPD | 5.852 | 21.24 | 0.6627 | 0.2765 | 0.2833 |
| PickScore+HPD | 5.871 | 21.45 | 0.6317 | 0.2734 | 0.2837 |
| PickScore+Pick-a-Pic | 5.720 | 21.73 | 0.6832 | 0.2735 | 0.2842 |
| Only Pick-a-Pic | 5.607 | 20.88 | 0.5615 | 0.2745 | 0.2743 |
| Only HPD | 5.413 | 20.37 | 0.3314 | 0.2734 | 0.2551 |
| Only HPSv2 | 5.702 | 20.43 | 0.7629 | 0.2496 | 0.3128 |
| Only PickScore | 6.000 | 21.45 | 0.6450 | 0.2498 | 0.2796 |

### 5.2.4. User Study

We conduct a subjective user study on FLUX-based backbones by randomly sampling 100 unique prompts from the HPD benchmark [54] and generating the corresponding images using our method and several state-of-the-art baselines. A total of 100 participants are invited to evaluate each comparison group by selecting the most favorable image across three criteria: Q1 General Preference (Which image do you prefer given the prompt?), Q2 Visual Appeal (Which image is more visually appealing?), Q3 Prompt Alignment (Which image better fits the text description?). Fig. 7 shows the approval percentage of each method in three aspects, which demonstrates our method outperforms the previous preference learning models on human feedback.

### 5.3. Ablation Study

To better understand the contributions of each component in our framework, we conduct a series of ablation studies on SD V1.5 [40] under the HyperAlign-S configuration unless otherwise specified. The results are summarized in Tab. 3.

**Effect of preference data for regularization loss $\mathcal{L}_\mathbf{G}$.** Our default configuration adopts HPSv2 as the reward model and Pick-a-Pic as the preference dataset for regularization. When replacing Pick-a-Pic with HPD while keeping HPSv2 fixed, our method still achieves strong performance, demonstrating the robustness and effectiveness of our method.

**Effect of reward–regularization configurations.** Beyond HPSv2, we combine PickScore and different preference datasets to optimize the hypernetwork. All combinations lead to consistently solid outcomes, verifying that HyperAlign can adapt to different reward and regularization sources. Our default choice, HPSv2 leans toward text–image alignment while Pick-a-Pic dataset favors visual appeal, provides balanced supervision that yields stronger overall improvements across metrics.

**Effect of reward loss $\mathcal{L}_\mathbf{R}$.** We further examine the influence of the reward loss by supervised fine-tuning using only preference data (Pick-a-Pic and HPD) and optimization using only reward signals (HPSv2 and PickScore). Results show that supervised fine-tuning with preference data alone yields marginal gains. Reward-only optimization boosts

thetic quality. In contrast, test-time alignment methods generate image with unstable effects and noticeable artifacts. Although training-based approaches achieve higher proxy scores, they tend to be over-optimized. The generated results lack practical usability, exhibiting anthropomorphic distortions and excessively saturated color tones. More visual results are provided in supplementary material.

### 5.2.3. Inference Efficiency

We report the average inference time for generating a single image in Tab. 2 for SD V1.5-based backbones and Tab. 1 for FLUX-based backbones. Our method achieves superior performance while requiring only a few seconds per image. When adopting the two-stage weight generation strategies, the inference efficiency can be further improved without sacrificing performance. In contrast, test-time scaling methods incur substantial computational overhead due to gradient computation or repeated sampling during model forwarding. Although such cost may be tolerable for small-scale models, it becomes prohibitive for large-scale backbones, where optimizing a single image can take several minutes, making the approach impractical. Both test-time alignment and our HyperAlign adjust the generative trajectory, however, the time cost of generating and loading adaptive weights in our method is nearly negligible, further demonstrating its efficiency and practicality.

most preference scores but severely degrades CLIP, indicating clear reward over-optimization.

## 6. Conclusion

We propose HyperAlign, a hypernetwork-based framework for efficient and effective test-time alignment of generative models. HyperAlign dynamically generates low-rank modulation weights across denoising steps, enabling trajectory-level alignment guided by reward signals. Its variants provide flexible trade-offs between computational efficiency and alignment precision. Extensive experiments on both diffusion and rectified flow backbones show that Hyper-Align delivers superior semantic consistency and aesthetic quality compared to existing fine-tuning and test-time alignment approaches. In the future, we aim to further enhance dynamic adaptation while developing more lightweight hypernetwork designs to improve efficiency and scalability.

## References

[1] Donghoon Ahn, Jiwon Kang, Sanghyun Lee, Jaewon Min, Minjae Kim, Wooseok Jang, Hyoungwon Cho, Sayak Paul, SeonHwa Kim, Eunju Cha, et al. A noise is worth diffusion guidance. *arXiv preprint arXiv:2412.03895*, 2024. 3

[2] Yuval Alaluf, Omer Tov, Ron Mokady, Rinon Gal, and Amit Bermano. Hyperstyle: Stylegan inversion with hypernetworks for real image editing. In *Proceedings of the IEEE/CVF conference on computer Vision and pattern recognition*, pages 18511–18521, 2022. 3

[3] Lichen Bai, Shitong Shao, Zikai Zhou, Zipeng Qi, Zhiqiang Xu, Haoyi Xiong, and Zeke Xie. Zigzag diffusion sampling: Diffusion models can self-improve via self-reflection. *arXiv preprint arXiv:2412.10891*, 2024. 3

[4] Arpit Bansal, Hong-Min Chu, Avi Schwarzschild, Soumyadip Sengupta, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Universal guidance for diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 843–852, 2023. 3

[5] Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion models with reinforcement learning. *arXiv preprint arXiv:2305.13301*, 2023. 3

[6] Rujikorn Charakorn, Edoardo Cetin, Yujin Tang, and Robert Tjarko Lange. Text-to-lora: Instant transformer adaption. *arXiv preprint arXiv:2506.06105*, 2025. 3

[7] Kevin Clark, Paul Vicol, Kevin Swersky, and David J Fleet. Directly fine-tuning diffusion models on differentiable rewards. *arXiv preprint arXiv:2309.17400*, 2023. 3

[8] Niklas Deckers, Julia Peters, and Martin Potthast. Manipulating embeddings of stable diffusion prompts. *arXiv preprint arXiv:2308.12059*, 2023. 3

[9] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 1

[10] Hanze Dong, Wei Xiong, Deepanshu Goyal, Yihan Zhang, Winnie Chow, Rui Pan, Shizhe Diao, Jipeng Zhang, Kashun Shum, and Tong Zhang. Raft: Reward ranked finetuning for generative foundation model alignment. *arXiv preprint arXiv:2304.06767*, 2023. 3

[11] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first International Conference on Machine Learning*, 2024. 3

[12] Luca Eyring, Shyamgopal Karthik, Karsten Roth, Alexey Dosovitskiy, and Zeynep Akata. Reno: Enhancing one-step text-to-image models through reward-based noise optimization. *arXiv preprint arXiv:2406.04312*, 2024. 3

[13] Luca Eyring, Shyamgopal Karthik, Alexey Dosovitskiy, Nataniel Ruiz, and Zeynep Akata. Noise hypernetworks: Amortizing test-time compute in diffusion models. *arXiv preprint arXiv:2508.09968*, 2025. 3

[14] Ying Fan, Olivia Watkins, Yuqing Du, Hao Liu, Moonkyung Ryu, Craig Boutilier, Pieter Abbeel, Mohammad Ghavamzadeh, Kangwook Lee, and Kimin Lee. Reinforcement learning for fine-tuning text-to-image diffusion models. *Advances in Neural Information Processing Systems*, 36, 2024. 3

[15] Dhruba Ghosh, Hannaneh Hajishirzi, and Ludwig Schmidt. Geneval: An object-focused framework for evaluating text-to-image alignment. *Advances in Neural Information Processing Systems*, 36:52132–52152, 2023. 7, 2

[16] Xiefan Guo, Jinlin Liu, Miaomiao Cui, Jiankai Li, Hongyu Yang, and Di Huang. Initno: Boosting text-to-image diffusion models via initial noise optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9380–9389, 2024. 3

[17] David Ha, Andrew Dai, and Quoc V Le. Hypernetworks. *arXiv preprint arXiv:1609.09106*, 2016. 3

[18] Haoran He, Jiajun Liang, Xintao Wang, Pengfei Wan, Di Zhang, Kun Gai, and Ling Pan. Scaling image and video generation via test-time evolutionary search. *arXiv preprint arXiv:2505.17618*, 2025. 3

[19] Xiaoxuan He, Siming Fu, Yuke Zhao, Wanli Li, Jian Yang, Dacheng Yin, Fengyun Rao, and Bo Zhang. Tempflow-grpo: When timing matters for grpo in flow models. *arXiv preprint arXiv:2508.04324*, 2025. 3

[20] Eric Hedlin, Munawar Hayat, Fatih Porikli, Kwang Moo Yi, and Shweta Mahajan. Hypernet fields: Efficiently training hypernetworks without ground truth by learning weight trajectories. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 22129–22138, 2025. 3

[21] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 1

[22] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 3

[23] Hamish Ivison, Akshita Bhagia, Yizhong Wang, Hannaneh Hajishirzi, and Matthew E Peters. Hint: hypernetwork instruction tuning for efficient zero-and few-shot generalisation. In *Proceedings of the 61st annual meeting of the Association for Computational Linguistics (volume 1: long papers)*, pages 11272–11288, 2023. 3

[24] Sunwoo Kim, Minkyu Kim, and Dongmin Park. Test-time alignment of diffusion models without reward over-optimization. *arXiv preprint arXiv:2501.05803*, 2025. 2, 3, 7

[25] Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-a-pic: An open dataset of user preferences for text-to-image generation. *Advances in Neural Information Processing Systems*, 36: 36652–36663, 2023. 2, 4, 6, 7

[26] Black Forest Labs. Flux.1-schnell, 2024. Accessed: 2024-08-17. 1, 2, 6

[27] Junzhe Li, Yutao Cui, Tao Huang, Yinping Ma, Chun Fan, Miles Yang, and Zhao Zhong. Mixgrpo: Unlocking flow-based grpo efficiency with mixed ode-sde. *arXiv preprint arXiv:2507.21802*, 2025. 2, 3, 4

[28] Shufan Li, Konstantinos Kallidromitis, Akash Gokul, Yusuke Kato, and Kazuki Kozuka. Aligning diffusion models by optimizing human utility. *arXiv preprint arXiv:2404.04465*, 2024. 2, 3, 4, 7

[29] Yuming Li, Yikai Wang, Yuying Zhu, Zhongyu Zhao, Ming Lu, Qi She, and Shanghang Zhang. Branchgrpo: Stable and efficient grpo with structured branching in diffusion models. *arXiv preprint arXiv:2509.06040*, 2025. 3

[30] Zhanhao Liang, Yuhui Yuan, Shuyang Gu, Bohan Chen, Tiankai Hang, Ji Li, and Liang Zheng. Step-aware preference optimization: Aligning preference with denoising performance at each step. *arXiv preprint arXiv:2406.04314*, 2024. 2, 3, 4, 7

[31] Buhua Liu, Shitong Shao, Bao Li, Lichen Bai, Haoyi Xiong, James Kwok, Sumi Helal, and Zeke Xie. Alignment of diffusion models: Fundamentals, challenges, and future. *arXiv preprint arXiv:2409.07253*, 2024. 2

[32] Jie Liu, Gongye Liu, Jiajun Liang, Yangguang Li, Jiaheng Liu, Xintao Wang, Pengfei Wan, Di Zhang, and Wanli Ouyang. Flow-grpo: Training flow matching models via online rl. *arXiv preprint arXiv:2505.05470*, 2025. 3, 1

[33] Nanye Ma, Shangyuan Tong, Haolin Jia, Hexiang Hu, Yu-Chuan Su, Mingda Zhang, Xuan Yang, Yandong Li, Tommi Jaakkola, Xuhui Jia, et al. Inference-time scaling for diffusion models beyond scaling denoising steps. *arXiv preprint arXiv:2501.09732*, 2025. 2, 3, 7

[34] Yuhang Ma, Xiaoshi Wu, Keqiang Sun, and Hongsheng Li. Hpsv3: Towards wide-spectrum human preference score, 2025. *URL https://arxiv. org/abs/2508.03789*. 4

[35] Yuval Nirkin, Lior Wolf, and Tal Hassner. Hyperseg: Patch-wise hypernetwork for real-time semantic segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4061–4070, 2021. 3

[36] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 1, 2

[37] Mihir Prabhudesai, Anirudh Goyal, Deepak Pathak, and Katerina Fragkiadaki. Aligning text-to-image diffusion models with reward backpropagation. *arXiv preprint arXiv:2310.03739*, 2023. 2, 3, 4, 7

[38] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PmLR, 2021. 7

[39] Vignav Ramesh and Morteza Mardani. Test-time scaling of diffusion models via noise trajectory search. *arXiv preprint arXiv:2506.03164*, 2025. 2, 3, 7

[40] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 1, 2, 6, 8

[41] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Wei Wei, Tingbo Hou, Yael Pritch, Neal Wadhwa, Michael Rubinstein, and Kfir Aberman. Hyperdreambooth: Hypernetworks for fast personalization of text-to-image models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6527–6536, 2024. 3

[42] Christoph Schuhmann. Laion-aesthetics. https://laion.ai/blog/laion-aesthetics/, 2022. Accessed: 2023 - 11- 10. 4, 7

[43] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024. 3

[44] Xiangwei Shen, Zhimin Li, Zhantao Yang, Shiyi Zhang, Yingfang Zhang, Donghao Li, Chunyu Wang, Qinglu Lu, and Yansong Tang. Directly aligning the full diffusion trajectory with fine-grained human preference. *arXiv preprint arXiv:2509.06942*, 2025. 2, 3, 4, 7

[45] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 3

[46] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020. 1, 3

[47] Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of score-based diffusion models. *Advances in neural information processing systems*, 34: 1415–1428, 2021. 3

[48] Przemyslaw Spurek, Artur Kasymov, Marcin Mazur, Diana Janik, Slawomir Konrad Tadeja, J Tabor, T Trzciński, et al. Hyperpocket: Generative point cloud completion. In *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 6848–6853. IEEE, 2022. 3

[49] Zhiwei Tang, Jiangweizhi Peng, Jiasheng Tang, Mingyi Hong, Fan Wang, and Tsung-Hui Chang. Tuning-free alignment of diffusion models with direct noise optimization. *arXiv preprint arXiv:2405.18881*, 2024. 2, 3

[50] Johannes Von Oswald, Christian Henning, Benjamin F Grewe, and João Sacramento. Continual learning with hypernetworks. *arXiv preprint arXiv:1906.00695*, 2019. 3

[51] Bram Wallace, Akash Gokul, Stefano Ermon, and Nikhil Naik. End-to-end diffusion latent optimization improves

classifier guidance. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7280–7290, 2023. 3

[52] Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou, Aaron Lou, Senthil Purushwalkam, Stefano Ermon, Caiming Xiong, Shafiq Joty, and Nikhil Naik. Diffusion model alignment using direct preference optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8228–8238, 2024. 2, 3, 4, 7

[53] Yibin Wang, Zhimin Li, Yuhang Zang, Yujie Zhou, Jiazi Bu, Chunyu Wang, Qinglin Lu, Cheng Jin, and Jiaqi Wang. Pref-grpo: Pairwise preference reward-based grpo for stable text-to-image reinforcement learning. *arXiv preprint arXiv:2508.20751*, 2025. 3

[54] Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis. *arXiv preprint arXiv:2306.09341*, 2023. 2, 4, 6, 7, 8

[55] Xiaoshi Wu, Yiming Hao, Manyuan Zhang, Keqiang Sun, Zhaoyang Huang, Guanglu Song, Yu Liu, and Hongsheng Li. Deep reward supervisions for tuning text-to-image diffusion models. *arXiv preprint arXiv:2405.00760*, 2024. 3

[56] Xin Xie and Dong Gong. Dymo: Training-free diffusion model alignment with dynamic multi-objective scheduling. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 13220–13230, 2025. 2, 3, 7

[57] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation. *Advances in Neural Information Processing Systems*, 36, 2024. 3, 7

[58] Yilun Xu, Mingyang Deng, Xiang Cheng, Yonglong Tian, Ziming Liu, and Tommi Jaakkola. Restart sampling for improving generative processes. *Advances in Neural Information Processing Systems*, 36:76806–76838, 2023. 3

[59] Zeyue Xue, Jie Wu, Yu Gao, Fangyuan Kong, Lingting Zhu, Mengzhao Chen, Zhiheng Liu, Wei Liu, Qiushan Guo, Weilin Huang, et al. Dancegrpo: Unleashing grpo on visual generation. *arXiv preprint arXiv:2505.07818*, 2025. 2, 3, 4, 7, 1

[60] Kai Yang, Jian Tao, Jiafei Lyu, Chunjiang Ge, Jiaxin Chen, Weihan Shen, Xiaolong Zhu, and Xiu Li. Using human feedback to fine-tune diffusion models without any reward model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8941–8951, 2024. 3

[61] Haotian Ye, Haowei Lin, Jiaqi Han, Minkai Xu, Sheng Liu, Yitao Liang, Jianzhu Ma, James Zou, and Stefano Ermon. Tfg: Unified training-free guidance for diffusion models. *arXiv preprint arXiv:2409.15761*, 2024. 3

[62] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2(3):5, 2022. 7

[63] Jiwen Yu, Yinhuai Wang, Chen Zhao, Bernard Ghanem, and Jian Zhang. Freedom: Training-free energy-guided conditional diffusion model. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 23174–23184, 2023. 2, 3, 7

[64] Tao Zhang, Cheng Da, Kun Ding, Huan Yang, Kun Jin, Yan Li, Tingting Gao, Di Zhang, Shiming Xiang, and Chunhong Pan. Diffusion model as a noise-aware latent reward model for step-level preference optimization. *arXiv preprint arXiv:2502.01051*, 2025. 3

# HyperAlign: Hypernetwork for Efficient Test-Time Alignment of Diffusion Models

## Supplementary Material

## 7. More Details of HyperAlign with Flow-Matching Models

As discussed in the main paper and demonstrated in the experiments (*e.g.*, experiments with FLUX), our HyperAlign method can be applied to both diffusion and flow-matching models, although the main paper primarily presents the formulation using diffusion models.

In Sec. 3.1 and Sec. 4.1, we discussed conditional generation under test-time diffusion guidance, where the denoising trajectory is adjusted by directly modifying the temporal states. This paradigm is also compatible with flow-matching models [26]. While the connections between diffusion and flow-matching models have been established and unified formulations have been presented in prior work [59], in this section, we provide more details for our method with flow-matching models.

**Conditional flow-matching models and score functions.** Different from reverse SDE in Eq. (2), a deterministic reverse probability flow ODE [46] takes the following form:

$$\mathrm{d}\mathbf{x}_t = \left[\mathbf{f}(\mathbf{x}_t) - \frac{1}{2}g_t^2\nabla_{\mathbf{x}_t}\log p_t(\mathbf{x}_t)\right]\mathrm{d}t. \quad (11)$$

For flow matching, the score $\nabla_{\mathbf{x}_t}\log p_t(\mathbf{x}_t)$ is implicitly linked to the velocity field $v_t$. Specifically, we define $\mathbf{x}_0 \sim p_{data}(\mathbf{x})$ and $\mathbf{x}_1 \sim \mathcal{N}(0,\mathbf{I})$, then the forward process can be formulated as a linear interpolation:

$$\mathbf{x}_t = \alpha_t\mathbf{x}_0 + \beta_t\mathbf{x}_1, \quad (12)$$

where $\alpha_t = 1-t$, $\beta_t = t$ and $t \in [0,1]$. Under this construction, we have the distribution $\mathbf{x}_t \sim \mathcal{N}(\alpha_t\mathbf{x}_0, \beta_t^2\mathbf{I})$, yielding the marginal score:

$$\begin{aligned}
\nabla_{\mathbf{x}_t}\log p_t(\mathbf{x}_t) &= \mathbb{E}\left[\nabla_{\mathbf{x}_t}\log p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)|\mathbf{x}_t\right] \\
&= -\frac{1}{\beta_t}\mathbb{E}[\mathbf{x}_1|\mathbf{x}_t].
\end{aligned} \quad (13)$$

For the velocity field $v_t(\mathbf{x}_t)$, we derive:

$$\begin{aligned}
v_t(\mathbf{x}_t) &= \mathbb{E}\left[\dot{\alpha}_t\mathbf{x}_0 + \dot{\beta}_t\mathbf{x}_1 \mid \mathbf{x}_t\right] \\
&= \dot{\alpha}_t\mathbb{E}[\mathbf{x}_0|\mathbf{x}_t] + \dot{\beta}_t\mathbb{E}[\mathbf{x}_1|\mathbf{x}_t] \\
&= \dot{\alpha}_t\mathbb{E}\left[\frac{\mathbf{x}_t - \beta_t\mathbf{x}_1}{\alpha_t} \mid \mathbf{x}_t\right] + \dot{\beta}_t\mathbb{E}[\mathbf{x}_1|\mathbf{x}_t] \\
&= \frac{\dot{\alpha}_t}{\alpha_t}\mathbf{x}_t - \frac{\dot{\alpha}_t\beta_t}{\alpha_t}\mathbb{E}[\mathbf{x}_1|\mathbf{x}_t] + \dot{\beta}_t\mathbb{E}[\mathbf{x}_1|\mathbf{x}_t] \\
&= \frac{\dot{\alpha}_t}{\alpha_t}\mathbf{x}_t - \left(\beta_t\dot{\beta}_t - \frac{\dot{\alpha}_t\beta_t^2}{\alpha_t}\right)\nabla_{\mathbf{x}_t}\log p_t(\mathbf{x}_t).
\end{aligned} \quad (14)$$

In particular, we resort to rectified flow [32] setting for better discussion. Substituting $\alpha_t$, $\beta_t$ and Eq. (14) into the sampling process $v_t = \frac{\mathrm{d}\mathbf{x}_t}{\mathrm{d}t}$ and applying discretization yields the update rule:

$$\mathbf{x}_{t+\Delta t} = (1 - \frac{\Delta t}{1-t})\mathbf{x}_t - \frac{t\,\Delta t}{1-t}\nabla_{\mathbf{x}_t}\log p_t(\mathbf{x}_t|\mathbf{c}), \quad (15)$$

where $p_t(\mathbf{x}_t|\mathbf{c})$ represents that flow-matching models learn a distribution with conditioning variable $\mathbf{c}$. For image generation, $\mathbf{c}$ is the input prompts indicating user's instruction for the generated contents. Similar to Eq. (4), this iterative denoising process also forms a trajectory $\{\mathbf{x}_t\}_{t=1}^0$ in the latent space, gradually transforming the noise $\mathbf{x}_1$ into a clean sample $\mathbf{x}_0$.

**Test-time alignment with reward-based guidance.** As mentioned in Sec. 4.1, test-time diffusion alignment methods adjust the generative trajectory to better satisfy alignment objectives. Specifically, gradient-based diffusion guidance directly compute the gradient from reward signals and uses them to steer the denoising trajectory by modifying the temporal states. Similarly, based on Bayes' rule, the score $\nabla_{\mathbf{x}_t}\log p_t(\mathbf{x}_t|\mathbf{c})$ in Eq. (14) can be divided into the unconditional score $\nabla_{\mathbf{x}_t}\log p_t(\mathbf{x}_t)$ and the correction gradient $\nabla_{\mathbf{x}_t}R(\mathbf{x}_{0|t},\mathbf{c})$. Since the first term is independent of the condition $\mathbf{c}$, we focus on the second term, which injects reward gradient into the velocity field:

$$\begin{aligned}
\nabla_{\mathbf{x}_t}R(\mathbf{x}_{0|t},\mathbf{c}) &= \nabla_{\mathbf{x}_t}R(\mathbf{x}_t - t\cdot v_\theta(\mathbf{x}_t,t)) \\
&= \frac{\partial R}{\partial\mathbf{x}_{0|t}}\cdot\left(\mathbf{I} - t\cdot\frac{\partial v_\theta(\mathbf{x}_t,t)}{\partial\mathbf{x}_t}\right),
\end{aligned} \quad (16)$$

where the reward function is actually applied on the decoded image domain through the decoder. For simplicity of discussion, we omit the decoder notation. By substituting Eq. (16) into Eq. (15), flow-matching models can achieve alignment by injecting reward-aware dynamics into the latent states of the next timesteps. Essentially, modifying the intermediate states between two timesteps corresponds to adjusting the sampling trajectory, which shows that our proposed HyperAlign naturally extends to flow-matching models as well.

## 8. Additional Experimental Details and Results

### 8.1. Additional Qualitative Results

We provide more visual results for qualitative evaluation.
**More results on visual comparison.** We provide additional visual comparison results as shown in Fig. 10 for SD

Table 4. GenEval Benchmark [15] evaluation based on SD V1.5.

| Method | Overall | Single object | Two object | Counting | Colors | Position | Color attribution |
|---|---|---|---|---|---|---|---|
| SD v1.5 | 0.39 | 0.95 | 0.39 | 0.30 | 0.73 | 0.04 | 0.05 |
| $\varepsilon$-greedy | 0.48 | 0.98 | 0.66 | 0.38 | 0.72 | 0.05 | 0.05 |
| DAS | 0.34 | 0.89 | 0.31 | 0.22 | 0.63 | 0.04 | 0.04 |
| FreeDoM | 0.45 | 0.95 | 0.49 | 0.44 | 0.79 | 0.06 | 0.07 |
| DyMO | 0.50 | 0.98 | 0.60 | 0.55 | 0.79 | 0.05 | 0.14 |
| AlignProp | 0.39 | 0.95 | 0.37 | 0.31 | 0.71 | 0.04 | 0.06 |
| Diffusion-DPO | 0.41 | 0.97 | 0.39 | 0.37 | 0.75 | 0.04 | 0.06 |
| SPO | 0.40 | 0.96 | 0.36 | 0.34 | 0.73 | 0.05 | 0.06 |
| Diffusion-KTO | 0.42 | 0.98 | 0.43 | 0.35 | 0.77 | 0.05 | 0.07 |
| DanceGRPO | 0.41 | 0.95 | 0.46 | 0.32 | 0.72 | 0.07 | 0.07 |
| HyperAlign | 0.52 | 0.98 | 0.66 | 0.40 | 0.82 | 0.09 | 0.24 |

Table 5. Comparison results on GenEval [15].

| Method | Overall | Single object | Two object | Counting | Colors | Position | Color attribution |
|---|---|---|---|---|---|---|---|
| FLUX | 0.63 | 0.97 | 0.77 | 0.68 | 0.78 | 0.21 | 0.44 |
| DanceGRPO | 0.68 | 0.98 | 0.86 | 0.72 | 0.78 | 0.22 | 0.46 |
| HyperAlign | 0.70 | 0.98 | 0.88 | 0.70 | 0.88 | 0.20 | 0.54 |

V1.5-based backbones and Fig. 11 for FLUX-based backbones, respectively. Compared to the baseline and existing models, our approach generates high-quality images more closely aligned with contextual semantics and better cater to human preferences. Moreover, all three variants of our LoRA weight generation strategy achieve strong performance, further demonstrating the effectiveness and flexibility of the proposed framework.

**Qualitative results on ablation studies.** In Tab. 3, We report ablation results examining the effects of different reward models and preference datasets. The experimental metrics show that our method remains effective and robust across diverse reward–preference configurations. In Fig. 13, we visualize the ablation study results. It is observed that the visual qualities of the generated image by our method (HPSv2-based and PickScore-based) are consistent with the numerical results, enhancing both aesthetic appeal and semantic correctness. Compared with SD V1.5 [40], supervised fine-tuning solely on preference datasets yields only marginal gains. Additionally, we observe that although reward-only optimization attains higher metric scores, it leads to over-optimized and visually saturated samples, which further demonstrates the effectiveness of our proposed method.

**Qualitative comparison on diversity.** Although Dance-GRPO [59] and MixGRPO [27] achieve high quantitative scores in Tab. 1, we further examine their generation behavior by sampling multiple outputs from the same prompt under different random initial noises, as shown in Fig. 12. We observe that both methods significantly reduce the inherent diversity of the FLUX backbone [26], producing im-

ages that collapse toward a single style or even a single identity, which is an indication of over-optimization. In contrast, our HyperAlign framework achieves strong preference alignment while preserving the model's native diversity, maintaining varied yet semantically faithful outputs across different noise initializations.

### 8.2. Additional Quantitative Results

We conduct quantitative evaluation on GenEval benchmark [15] and show comparisons in Tab. 4 for SD V1.5-based backbones. The results show that our method performs very well and shows superiority in many aspects, *e.g.*, overall, attribute binding and object synthesis. To further evaluate the ability to capture high-level semantics, we incorporate the CLIP score into the training objective following DanceGRPO [59]. The main quantitative results on the GenEval benchmark for FLUX-based backbones are presented in Tab. 5. Corresponding qualitative comparisons are provided in Fig. 14 and Fig. 15. Compared with HPS-only optimization, jointly optimizing with both HPS and CLIP objectives yields noticeably better semantic consistency.

### 8.3. Additional Analyses

**Dynamics of LoRA Weight.** We further analyze the intermediate dynamics of the hypernetwork-based trajectory alignment using Stable Diffusion v1.5 [40]. Specifically, we examine the LoRAs generated at different steps by HyperAlign-S. Fig. 8 illustrates how the generated LoRA evolves across the diffusion process relative to the LoRA at the initial step. We compute both the cosine similarity and the $\ell_1$ relative change of the LoRA at each step with respect
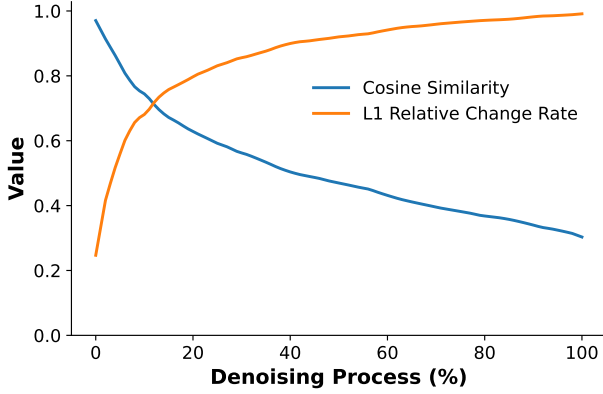
Figure 8. Visualization of LoRA weight variations at different timesteps relative to the initial state $T$. The cosine similarity and $\ell_1$ difference of the LoRA generated at each step relative to the LoRA at the initial step are calculated and demonstrated. Average over 1000 prompts.

to the initial step. The cosine similarity between the weights at each timestep and the initial state at $T = 1000$ steadily decreases, while the $\ell_1$ relative change rate consistently increases. This indicates that the LoRA weights progressively deviate from their initial configuration, highlighting the distinct functional roles played by different timesteps in the diffusion process.

**Variation of LoRA Weights.** To further examine the LoRAs generated for different prompts at different steps, we randomly sample 200 prompts and obtain their corresponding LoRAs across the diffusion process. We then perform PCA on the LoRA parameters and visualize the top two principal components in Fig. 9. The results show that HyperAlign produces distinct LoRAs for different inputs. The variances, reflected by the spread of points in each subplot of Fig. 9, are larger at the initial and early steps than at later steps. This aligns with our observation that the generation process requires more prompt-specific alignment during earlier stages.

### 8.4. Additional Details on Human Evaluation

We administer our user study using structured survey forms, where each prompt is presented as an independent section. Within each section, we present multiple images generated by different methods for the same prompt. Participants answer three questions: Q1 (Fig. 16), Q2 (Fig. 17) and Q3 (Fig. 18). Each question targets a different aspect of preference (overall preference, visual appeal, and prompt alignment), and participants are asked to select the most favorable image among the provided options. Participants are recruited via an online platform and remain fully anonymous. To ensure reliable evaluation, all participants are required to hold at least a bachelor's degree, and their privacy and identity are strictly protected throughout the study.

## 9. Ethical and Social Impacts

In this work, we propose HyperAlign, a hypernetwork-based test-time alignment framework for text-to-image diffusion or flow-matching models. While our method enhances alignment with human preferences and semantic consistency with input prompt, it also introduces ethical and social considerations that must be carefully addressed to ensure responsible AI deployment. Specifically, our method is built upon pre-trained backbones and preference datasets, which may inherit or amplify existing societal biases. To prevent reinforcing stereotypes, we recommend auditing reward signals and ensuring sufficient dataset diversity. Moreover, stronger semantic control by our methods can increase the risk of misuse, such as generating harmful, misleading, or identity-revealing content, raising significant privacy and safety concerns. To mitigate these issues, we advocate implementing safeguards, including content moderation and clear ethical usage guidelines. Overall, the benefits of HyperAlign substantially outweigh these potential concerns. The proposed framework lowers the barrier to high-quality test-time alignment and improves accessibility for diverse user groups. By balancing strong performance with ethical responsibility, HyperAlign aims to support fairness, transparency and safe deployment in real-world applications.
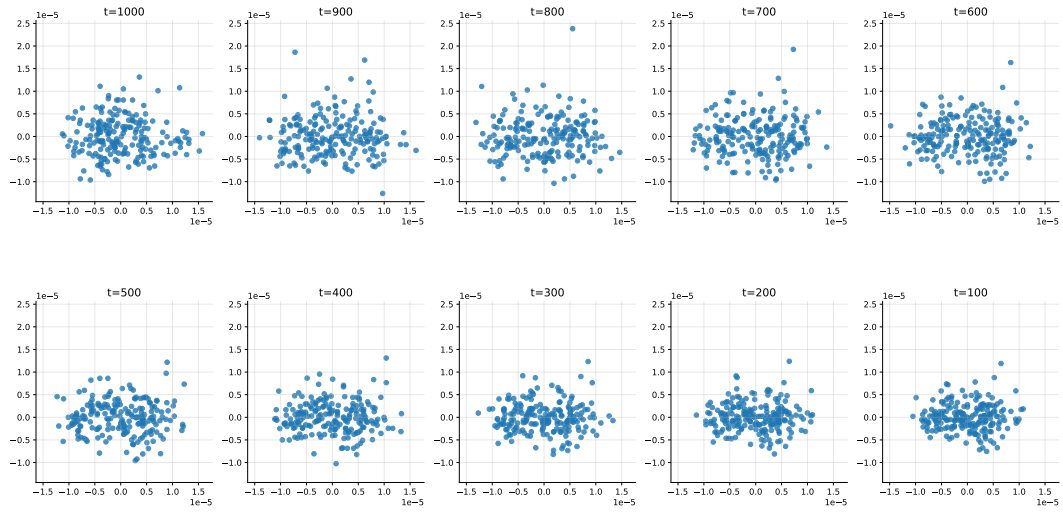
Figure 9. Visualization of the statistics of prompt-specific LoRA weights across different steps. The top two PCA components of the LoRAs generated for different prompts (200 examples) at each step are shown.
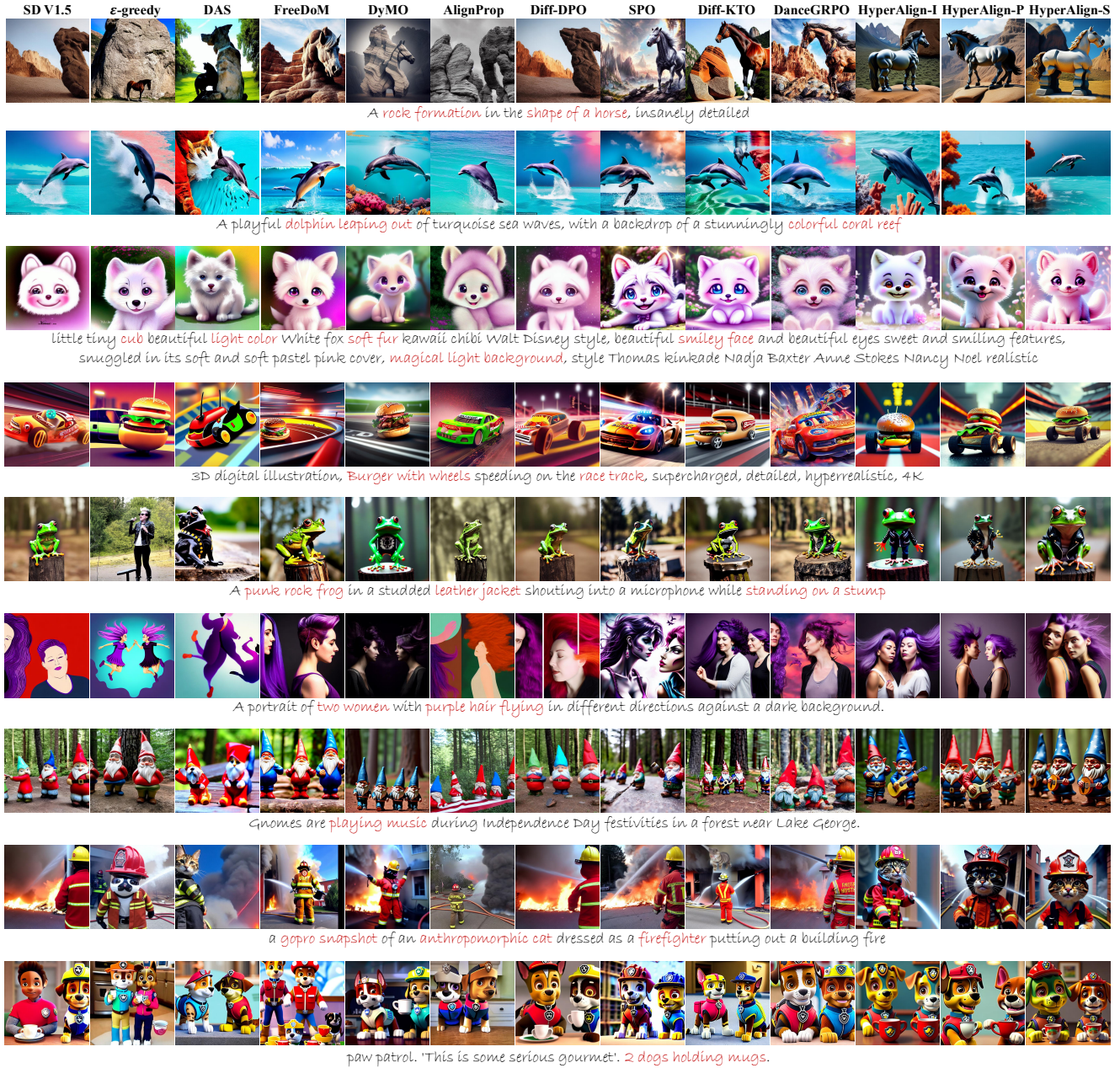
| SD V1.5 | ε-greedy | DAS | FreeDoM | DyMO | AlignProp | Diff-DPO | SPO | Diff-KTO | DanceGRPO | HyperAlign-I | HyperAlign-P | HyperAlign-S |

A rock formation in the shape of a horse, insanely detailed

A playful dolphin leaping out of turquoise sea waves, with a backdrop of a stunningly colorful coral reef

little tiny cub beautiful light color White fox soft fur kawaii chibi Walt Disney style, beautiful smiley face and beautiful eyes sweet and smiling features, snuggled in its soft and soft pastel pink cover, magical light background, style Thomas kinkade Nadja Baxter Anne Stokes Nancy Noel realistic

3D digital illustration, Burger with wheels speeding on the race track, supercharged, detailed, hyperrealistic, 4K

A punk rock frog in a studded leather jacket shouting into a microphone while standing on a stump

A portrait of two women with purple hair flying in different directions against a dark background.

Gnomes are playing music during Independence Day festivities in a forest near Lake George.

a gopro snapshot of an anthropomorphic cat dressed as a firefighter putting out a building fire

paw patrol. 'This is some serious gourmet'. 2 dogs holding mugs.

Figure 10. Qualitative comparison based on SD V1.5 backbones.

5

Figure 11. Qualitative comparison based on FLUX backbones.

The column headers from left to right are: FLUX, BoN, ε-greedy, FreeDoM, DyMO, DanceGRPO, MixGRPO, SRPO, HyperAlign-I, HyperAlign-P, HyperAlign-S.

The prompts (row captions) are:
- A cute little anthropomorphic Tropical fish knight wearing a cape and a crown in short, pale blue armor.
- Ralsei and Asriel from Deltarune eating pizza.
- A white polar bear cub wearing sunglasses sits in a meadow with flowers.
- Family assembling missile in living room.
- A ginger haired mouse mechanic in blue overalls in a cyberpunk scene with neon slums in the background.
- Mechanical bee flying in nature, electronics, motors, wires, buttons, lcd, led instead of eyes, antennas instead of feet
- A realistic furry anthro fox
- A tower of cheese
- At Song dynasty, a pretty woman in chinese was walking along the river

6

little tiny cub beautiful light color White fox soft fur kawaii chibi Walt Disney style, beautiful smiley face and beautiful eyes sweet and smiling features, snuggled in its soft and soft pastel pink cover, magical light background, style Thomas kinkade Nadja Baxter Anne Stokes Nancy Noel realistic



16-year-old teenager wearing a white bear-ear hat with a smirk on their face.

Figure 12. Diversity comparison based on FLUX backbones.

Figure 13. The visual results of ablation study.

| FLUX | HyperAlign (w HPSv2) | HyperAlign (w HPSv2&CLIP) |
|------|----------------------|---------------------------|

A photo of a toilet

A photo of a tie above a sink

A photo of four baseball gloves

A photo of a yellow bird and a black motorcycle

Figure 14. The comparison includes the baseline FLUX outputs, the results obtained through HPS-only optimization, and the versions further improved with joint HPS and CLIP objectives.

9

| FLUX | HyperAlign (w HPSv2) | HyperAlign (w HPSv2&CLIP) |

A photo of a laptop left of a cow

A group of mongooses scuttle about, set against a desert backdrop,
bathed in bright and warm earth tones

Neon rain alley at night: skyscrapers, a cyan motorcycle parked beside a magenta
umbrella, a warm yellow taxi streaking past; cinematic, volumetric fog, wet pavement

A slice of pizza floating through space with stars in the background

Figure 15. The comparison includes the baseline FLUX outputs, the results obtained through HPS-only optimization, and the versions further improved with joint HPS and CLIP objectives.

# Human preference investigation

Given a description, please choose the best one according to the question

## Test 1

Given a description prompt, select the most appropriate option based on the question.

1. Description: A 3D rendering of anime schoolgirls with a sad expression underwater, surrounded by dramatic lighting

Which image do you prefer given the prompt?

○ Image 1

○ Image 2
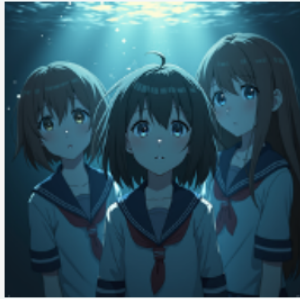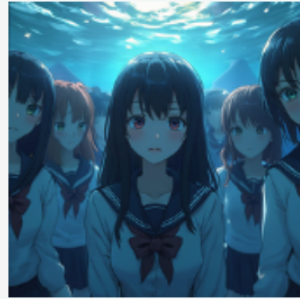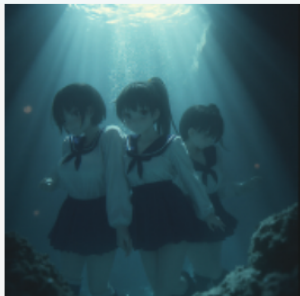
○ Image 3

○ Image 4

○ Image 5

○ Image 6

Figure 16. The screenshot of human preference investigation: Which image do you prefer given the prompt?
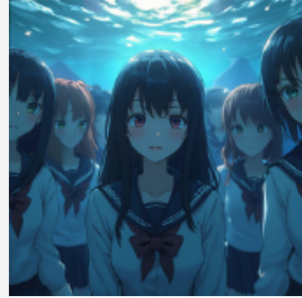
Figure 17. The screenshot of human preference investigation: Which image is more visually appealing?

3. Description: A 3D rendering of anime schoolgirls with a sad expression underwater, surrounded by dramatic lighting

Which image better fits the text description?



○ Image 1

○ Image 2

○ Image 3

○ Image 4

○ Image 5

○ Image 6

Figure 18. The screenshot of human preference investigation: Which image better fits the text description?