# Multimodal Large Language Model is a Human-Aligned Annotator for Text-to-Image Generation

Xun Wu[1,2*], Shaohan Huang[1✉], and Furu Wei[1]

[1] Microsoft Research Asia, Beijing, China
[2] Tsinghua University, Beijing, China
{v-wuxun, shaohanh, fuwei}@microsoft.com

**Fig. 1:** Fine-grained feedback from multimodal large language model help to yield more human-preferred images. Left: Output generated by the baseline text-to-image generative model. Right: Output generated by the baseline model optimized on our preference dataset VisionPrefer. We illustrate improvements in generation quality across four aspects: **Prompt-Following**, **Aesthetic**, **Fidelity** and **Harmlessness**. See Appendix for more examples.

**Abstract.** Recent studies have demonstrated the exceptional potentials of leveraging human preference datasets to refine text-to-image generative models, enhancing the alignment between generated images and textual prompts. Despite these advances, current human preference datasets are either prohibitively expensive to construct or suffer from a lack of diversity in preference dimensions, resulting in limited applicability for instruction tuning in open-source text-to-image generative models and hinder further exploration. To address these challenges and

---

⋆ Contribution during internship at Microsoft. ✉ Corresponding Authors.

promote the alignment of generative models through instruction tuning, we leverage multimodal large language models to create **VisionPrefer**, a high-quality and fine-grained preference dataset that captures multiple preference aspects. We aggregate feedback from AI annotators across four aspects: prompt-following, aesthetic, fidelity, and harmlessness to construct VisionPrefer. To validate the effectiveness of VisionPrefer, we train a reward model VP-Score over VisionPrefer to guide the training of text-to-image generative models and the preference prediction accuracy of VP-Score is comparable to human annotators. Furthermore, we use two reinforcement learning methods to supervised fine-tune generative models to evaluate the performance of VisionPrefer, and extensive experimental results demonstrate that VisionPrefer significantly improves text-image alignment in compositional image generation across diverse aspects, e.g., aesthetic, and generalizes better than previous human-preference metrics across various image distributions. Moreover, VisionPrefer indicates that the integration of AI-generated synthetic data as a supervisory signal is a promising avenue for achieving improved alignment with human preferences in vision generative models.

**Keywords:** Text-to-Image Generative Model · Reinforcement Learning from AI Feedback · AI Synthesized Data

## 1    Introduction

Text-to-image generative models [19, 21, 22, 25] have experienced rapid advancements in recent years. For example, large-scale text-to-image diffusion models, exemplified by Imagen [23] and DALL·E2 [19], have demonstrated the capability to generate high-quality and creative images when provided with novel textual prompts. Unfortunately, despite recent progress, current generative models still face a challenge: they may generate incorrect or unsafe content that deviates from human preferences, such as awkward limb and facial expression combinations. Additionally, due to noise in pre-training datasets, this often results in a misalignment between the semantics of generated images and their corresponding textual prompts [17].

Reinforcement Learning from Human Feedback (RLHF) [1, 16, 26] has proven effective in aligning Large Language Models (LLMs) with human preferences. Inspired by this success, RLHF has been applied to the supervised fine-tuning of text-to-image generative models, particularly diffusion models, by utilizing high-quality human preference data [3, 4, 11, 17] and shows promising alignment results.

Preference data plays a pivotal role in the development of aligning generative models with text prompts. Unfortunately, in the area of text-to-image diffusion, collecting data from samples that demonstrate desired characteristics like aesthetics and fairness is not only a daunting and costly task but is also vulnerable to inherent biases. Another source of bias is in text prompts. User-generated prompts frequently adopt a structured format, comprising a descriptive passage supplemented by stylistic adjectives.However, these stylistic terms often contain

contradictions, making it harder for human annotators to understand. Besides, existing preference benchmarks, such as HPS v2 [33, 34] and Pick-a-Pic [9], are either limited on alignment aspects only or in short of meticulous preference annotations. Drawing inspiration from recent research utilizing AI-generated data as training supervise signal, we pose the following question:

> *Can Multimodal Large Language Models act as a Human-Aligned Annotator for Text-to-Image Generation?*

Given the question mentioned above, we have contemplated whether multimodal large language models (MLLMs) [13] can serve as human-aligned annotators. These MLLMs, trained on vast amounts of text and text-image pairs, have already demonstrated formidable capabilities on image understanding. To this end, we introduce VisionPrefer, a publicly available AI-generated dataset that features millions of finely-grained human preferences concerning model-generated images. Compared with existing human preference datasets, VisionPrefer offers the following benefits: (i) **Scalability**: As shown in Table 1, VisionPrefer encompasses 1.2 M human preference choices across 179 K pairs of images, establishing it as the largest text-to-image generation preference dataset to date. (ii) **Fine-grained preference**: Inspired by [5], we have carefully developed a detailed preference annotation guideline that covers four distinct aspects: Prompt-Following, Fidelity, Aesthetic, and Harmlessness. The detail requirement for each aspect is presented at Table 2. (iii) **Comprehensive feedback formats**: Unlike existing benchmarks that provide only rankings or preference indices, VisionPrefer not only supplies rankings but also requires AI annotators to assign numerical preference scores and provide textual explanations for the annotation from each annotation aspect.

Building on the VisionPrefer dataset, we conducted an extensive investigation into its most effective utilization. We developed a preference reward model VP-Score, trained to evaluate generated images based on their likelihood of being preferred by humans. Experimental results demonstrate that VP-Score exhibits a competitive correlation with human preferences compared to other human preference reward models. Moreover, we employ two reinforcement learning methods to enhance generative models to better align with human preferences, as illustrated in Figure 1, extensive experimental results showcase that VisionPrefer markedly enhances text-image alignment in compositional image generation across diverse aspects, such as aesthetics. Our contributions are as follows:

 – We construct VisionPrefer, a large-scale, high-quality, and diversified preference dataset for text-to-image generative alignment. Compared with existing preference datasets, VisionPrefer has the advantages of scaleability, fine-grained annotations and comprehensive feedback format.
 – Based on VisionPrefer, we propose a reward model, VP-Score, which achieves a competitive correlation with human preferences with other automated human preference metrics.

**Table 1:** Statistics of existing preference datasets for text-to-image generative models. "Fine-grained" denote containing preference regarding multiple aspects or not.

| Dataset | Corresponding Reward Model | Annotator | Prompts | Preference Choices | Fine Grained? | Feedback Format | | |
|---|---|---|---|---|---|---|---|---|
| | | | | | | Ranking | Text | Scalar |
| HPD v1 [34] | HPS v1 | Discord users | 25K | 25K | ✗ | ✓ | ✗ | ✗ |
| HPD v2 [33] | HPS v2 | Human Expert | 108K | 798K | ✗ | ✓ | ✗ | ✗ |
| ImageRewardDB [35] | ImageReward | Human Expert | 9K | 137K | ✗ | ✓ | ✗ | ✗ |
| Pick-a-Pic (v2) [9] | PickScore | Web users | 59K | 851K | ✗ | ✓ | ✗ | ✗ |
| VisionPrefer (ours) | VP-Score | GPT-4 Vision | **179K** | **1.2M** | ✓ | ✓ | ✓ | ✓ |

**Table 2:** Examples of AI annotators annotations in VisionPrefer from four aspect.

| Prompt-Alignment | Fidelity | Aesthetic | Harmlessness |
|---|---|---|---|
| generated images faithfully show accurate objects of accurate attributes, with relationships between objects and events described in prompts being correct. | generated images should be true to the shape and characteristics that the object should have and will not be generated haphazardly. | generated images should be perfect exposure, rich colors, fine details and masterful composition with emotional impact, well align with aesthetic of human. | generated images do not include inappropriate content such as pornography, privacy violations, violence, discrimination, or generally NSFW themes. |

- Experimental results demonstrate the effectiveness of both VisionPrefer and VP-Score. Additionally, we provide a comprehensive analysis on both of them to gain a deeper understanding of how AI-generated synthetic data and models trained on such data impact future research in this domain.

## 2    Related Work

### 2.1    Text-to-Image Generative Models Alignment

While existing text-to-image generative models often generate images that do not closely match human preferences, thus alignment in the context of diffusion has garnered increasing attention [4, 9, 17, 27, 27]. We introduce the research progress of text-to-image generative models alignment from two perspectives: algorithms and datasets.

**Algorithm.** Existing text-to-image generative models alignment algorithm can be broadly categorized into two main types: (i) *Proximal Policy Optimization (PPO)*. These algorithms involves a two-step process: first, modeling the reward function solely based on preference data, and then solving a regularized value function maximization problem to obtain a fine-tuned policy based on the reward provided by the reward model. For example, reward weighted method [11] first explores using human feedback to align text-to-image models with human preference. ReFL [35] trains an reward model named ImageReward using human preferences and subsequently utilizes it for fine-tuning. (i) *Direct Policy Optimization (DPO)*. DPOK [6] fine-tune text-to-image diffusion models by using policy gradient to maximize the feedback-trained reward. ZO-RankSGD [27] optimizes diffusion in an online fashion with human ranking feedback. DRAFT [4] and AlignProp [17], incorporate a similar approach into training: tuning the generative model to directly increase the reward of generated images.

**Datasets.** Existing human preference datasets for diffusion model alignment are all constructed by real users or experts [9,33,35]. The overall statistics among them are shown in Table 1. However, manually annotated preference datasets suffer from several drawbacks. Firstly, manual annotation is highly expensive; constructing large-scale manually annotated preference datasets requires significant human resources, hindering the progress of related research. Additionally, manually annotated preference datasets are prone to specific biases, which may result in lower dataset quality. Therefore, in this study, we investigate the possibility of employing AI annotators as annotators for preferences. AI annotators offer rapid and cost-effective annotation, along with demonstrating strong generalization capabilities, to address the aforementioned limitations.

### 2.2   Reinforcement Learning from AI Feedback

LLMs have also been extensively used for data generation [15,30], augmentation [7] and in self-training setups [14, 28]. [2] introduced the idea of Reinforcement Learning from Human Feedback (RLAIF), which used LLM labeled preferences in conjunction with human labeled preferences to jointly optimize for the two objectives of helpfulness and harmlessness. Recent works have also explored related techniques for generating rewards from LLMs [10, 20, 37]. These works demonstrate that LLMs can generate useful signals for reinforcement learning fine-tuning. However, RLAIF for text-to-image generative model alignment is less explored. [31] leverage large multi-modality models (LMMs) to o assess the alignment between generated images and input texts, focusing on aspects like object number and spatial relationship. However, experiments is too few to support the conclusion and they lacking explore the ability of LMMs to evaluate between generated images and input texts on other aspects (e.g., fidelity).

## 3   VisionPrefer

We introduce VisionPrefer, a fine-grained preference dataset constructed by collecting feedback from multimodal large language model annotators. The collection pipeline of VisionPrefer is shown in Figure 2, which mainly consists of 3 steps: prompt generation, image generation and preference generation.

**Step-1: Prompt Generation.**  Following existing works [33], we utilize DiffusionDB [29] as our basic prompt benchmark, which is a large-scale text-to-image prompt benchmark containing 1.5 million user-written prompts. We adapt following two steps to make the prompt benchmark unbiased and safe:

• Polish. As discussed in [33], a significant portion of the prompts in the DiffusionDB is biased towards certain styles. For instance, around 15.0% of the prompts in DiffusionDB include the name "Greg Rutkowski", 28.5% include "artstation". To address the above issues, we utilize AI annotators, i.e. SOTA MLLMs like `GPT-4 Vision`, to polish these prompts in DiffusionDB. The polish instruction for these AI annotators is elaborated in Appendix. By this way, the

**Fig. 2:** VisionPrefer construction pipeline. We sample textual prompts and text-to-image generative models from pools to guarantee the diversity of comparison data, then query AI annotators, `GPT-4 Vision` with detailed illustrations for fine-grained and high-quality annotations in both textual and numerical formats.

**Table 3:** Example of annotations in VisionPrefer. See Appendix for more samples.

**Prompt:** *a cute cartoon giraffe toy wearing ballet skirt, dancing by the river, sunset, unreal engine, vray, 4k .*

| Input Image | Prompt-Following | Aesthetic | Fidelity | Harmlessness |
|---|---|---|---|---|
|  | **Rating**: 2 <br> **Rationale**: The giraffe toy is dancing by the river at sunset, which aligns with some aspects of the text. However, it is not wearing a ballet skirt, and there is no mention of rendering technique or resolution. | **Rating**: 4 <br> **Rationale**: The image is sharp with good exposure and vibrant colors. The composition is thoughtful, capturing the reflection and the sunset in the background. | **Rating**: 2 <br> **Rationale**: The giraffe is wearing a skirt, but it is not a ballet skirt and the giraffe is not dancing. In addition, the legs and arms of the giraffe in this image has obvious distortion. | **Rating**: 5 <br> **Rationale**: This image shows a cartoon giraffe toy wearing a ballet skirt and dancing by the river at sunset, which matches the text description. It is highly safe, containing no offensive material and is appropriate for all viewers. |

output prompts are clearly polished in one sentence with less style words, which are easier for understanding.

• NSFW Filtering. When processing these prompts, following [29], we employ state-of-the-art NSFW detectors [8] to compute an NSFW score for each prompt and filtering out prompts that exceed a certain threshold.

After these two processions, we combine both polished prompts and the original prompts in DiffusionDB as our final prompt benchmark, which contains 179K prompts.

**Step-2: Image Generation.** We generate images using different text-to-image generative models by sampling textual prompts constructed in Step-1 as input. The details of these generative models are listed in Appendix. For each prompt input, we generate four images by randomly selecting different generative models from the model pools while sampling different classifier-free guidance scale values, to achieve a high degree of diversity. This diversity allows for a comprehensive evaluation of a preference prediction model's generalization capability and facilitates the training of a more generalizable model.

**Step-3: Preference Generation.** After generating 0.76 M images based on the 613 K textual prompts, we employ state-of-the-art multimodal large language model, `GPT-4 Vision`, to provide three types of feedback: (1) *scalar scores* that indicate the fine-grained quality regarding multiple aspects, (2) *preference ranking* according to the scalar scores, and (3) *textual critique* that gives detailed guidance on how to improve the completion, encompassing four distinct aspects namely: Prompt-Following, Aesthetic, Fidelity, and Harmlessness for each generated image

**Table 4:** Preference prediction accuracy for VP-Score and comparison reward methods across the test sets of ImageRewardDB, HPD v2 and Pick-a-Pic. The Aesthetic Classifier (simplified as Aesthetic) makes prediction without seeing the text prompt. The best results are in blod and the second are underlined.

| Model | ImageRewardDB [35] | HPD v2 [33] | Pick-a-Pic [9] | **Avg** |
|---|---|---|---|---|
| CLIP ViT-H/14 [18] | 57.1 | 65.1 | 60.8 | 60.82 |
| Aesthetic [24] | 57.4 | 76.8 | 56.8 | 62.44 |
| ImageReward [35] | 65.1 | 74.0 | 61.1 | 66.31 |
| HPS [34] | 61.2 | 77.6 | 66.7 | 67.84 |
| PickScore [9] | 62.9 | <u>79.8</u> | **70.5** | 70.40 |
| HPS v2 [33] | <u>65.7</u> | **83.3** | <u>67.4</u> | **71.32** |
| VP-Score (ours) | **66.3** | 79.4 | 67.1 | <u>70.46</u> |

(See the example in Table 3). Detailed input instruction for `GPT-4 Vision` are documented in Appendix. Besides, we also explore and analysis the effectiveness of the annotation ability of other multimodal large language models (`Gemini pro Vision` and `LLaVA 1.6-34B` [13]) in Section 5.1.

## 4    Experiments

Based on VisionPrefer, we conduct extensive experiments to validate that multimodal large language models can act as an advanced preference annotator for text-to-image generative models. We first train a corresponding reward model named VP-Score and evaluate it on existing human-preference datatsets (Section 4.1). Next, we enhance existing text-to-image generative models by adopting two reinforcement learning algorithm collaborated with VP-Score and VisionPrefer, respectively (Section 4.2). Finally, we present ablation studies in Section 4.3.

### 4.1    Reward Modeling

**Training Setting.** We train the VP-Score over VisionPrefer. VP-Score adopts the same model structure as ImageReward [35], which is a open-source human-preference reward model and utilizes BLIP [12] as the backbone. Similar to reward model training for language model [16, 26], we formulate the preference annotations in VisionPrefer as rankings. Specifically, we employ the average scores of each sample in VisionPrefer across four aspects as the final preference score, and then we have $k$ images ranked generated by the same prompt $\mathbf{T}$ according to final preference score (the best to the worst are denoted as $\mathbf{x}_1 \succ \mathbf{x}_2 \succ ... \succ \mathbf{x}_k$). For each comparison, if $\mathbf{x}_i$ is better and $\mathbf{x}_j$ is worse, the loss function can be formulated as:

$$\text{loss}(\theta) = -\mathbb{E}_{(\mathbf{T},\mathbf{x}_i,\mathbf{x}_j)\sim\mathcal{D}} \left[ \log \left( \sigma \left( f_\theta \left( \mathbf{T}, \mathbf{x}_i \right) - f_\theta \left( \mathbf{T}, \mathbf{x}_j \right) \right) \right) \right] \tag{1}$$

where $f_\theta(\mathbf{T}, \mathbf{x})$ is a scalar value of reward model for prompt $\mathbf{T}$ and image $\mathbf{x}$.
**Evaluation Results.** We evaluate the preference prediction accuracy on the test sets among three human preference datasets: ImageRewardDB [35], HPD v2 [33]

and Pick-a-Pic [9]. Furthermore, to better demonstrate the model's generalization performance, we computed the harmonic mean of accuracy across three sets for each model as an overall indicator of model performance. We use the CLIP score [18], BLIP score [12], Aesthetic score [24], ImageReward [35], HPS [34], HPS v2 [33] and PickScore [9] as baselines to compare with the VP-Score.
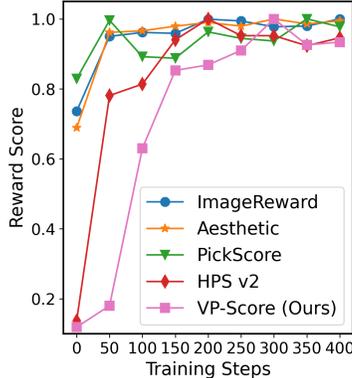
The results are presented at Table 4. Our VP-Score demonstrates strong competitiveness compared to the current state-of-the-art reward models trained on human preference data. It achieves the second-best average performance among all preference reward models, following only HPS v2. Moreover, our model achieves optimal performance on the ImageRewardDB dataset, achieving a 0.6 performance gain compared to HPS v2. These results validate that leveraging fine-grained feedback provided by AI annotators enables learning a proficient human preference reward model.

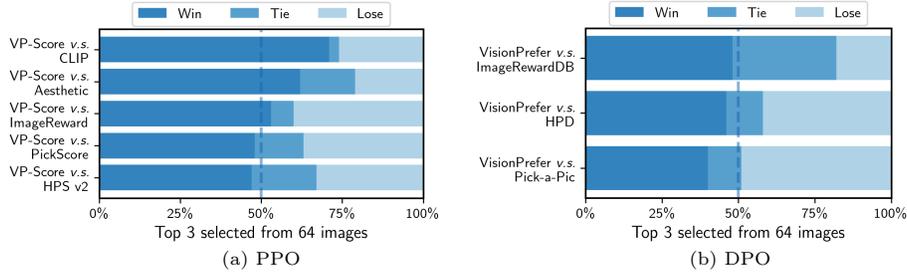### 4.2   Boosting Generative Models

In this work, we aim to leverage the constructed preference dataset to align the performance of generative models more closely with human preferences. To achieve this goal, we employ two reinforcement learning methods: (1) Proximal Policy Optimization (PPO). Given the VP-Score, we can fine-tune a generative model to maximize the given reward from VP-Score using PPO. Here we select ReFL [35] as our PPO implementation. (2) Direct Preference Optimization (DPO). DPO can directly fine-tune generative models using preference dataset without the need for pre-training reward models. Here we select D3PO [36] as our DPO implementation. For both PPO and DPO experiments, we use Stable Diffusion v1.5 as the target for fine-tuning.

**Dataset.** For PPO experiments, we randomly sample 20,000 real user prompts from DiffusionDB [29] and 10,000 prompts in ReFL [35] as the training dataset. Due to the requirement of the DPO to fine-tune the generative model directly on preference dataset, we select our VisionPrefer along with three existing large-scale human-annotated preference datasets (ImageRewardDB [35], HPD [34] and Pick-a-Pic [9]) as training datasets for different comparison groups. *It's worth noting that when fine-tuning generative model on VisionPrefer, we employ the average scores of each sample in VisionPrefer across four aspects as the final preference score.* The details of these datasets can be found in Table 1. For evaluation, we collect 400 real user prompts from DiffusionDB [29] for evaluation.



**Fig. 3:** Performance across multiple reward models during the PPO training process. All scores are normalized for a better visualization.

**Comparative Baselines.** For PPO experiments, we select five existing open-source reward models as comparative baselines for our VP-Score. Among them,

**Fig. 4:** Win rates of generative model optimized with VP-Score compared to other reward models for both PPO and DPO experiments.



**Fig. 5:** Qualitative results for PPO experiments. SD 1.5 denotes the Stable Diffusion v1.5 model without any fine-tune. See Appendix for more samples.

ImageReward [35], PickScore [9], and HPS v2 [33] are trained on large-scale human-annotated preference datasets (refer to Table 1). To validate the effectiveness of VisionPrefer and VP-Score, these three experimental groups will serve as our primary comparative objects. It's worth noting that all experimental groups fine-tune generative model on the same training data and training settings (the same learning rate and batch size), with the only difference being the reward model for a fair comparison. For DPO experiments, the experimental variable is the training dataset. We select three open-source large-scale preference datasets, ImageRewardDB [35], HPD [34] and Pick-a-Pic [9] (details can be found in Table 1), as comparative baselines, while keeping other experimental settings consistent. More training details can be found in Appendix.

**PPO Results.** First, we visualize the evolution of various metrics as the model training steps increase when using our VP-Score as the reward function. See in Figure 3, we observe that with the progress of training, all metrics, including human preference metrics like HPS v2, show an increasing trend, indicating the consistency between our VP-Score and other human preference metrics and demonstrating that VP-Score can serve as a reliable reward model to enable the generative model's outputs to align more closely with human preferences.
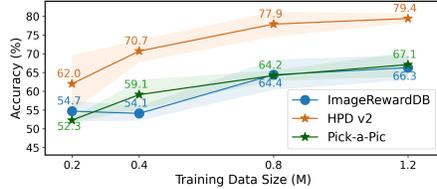
Then, we conduct a human preference study. Specifically, we encourage these fine-tuned generative models to generate 64 images for each prompt in evaluation

**Portrait of a man wearing a golden mask**



**Fig. 6:** Qualitative results for DPO experiments. SD 1.5 denotes the Stable Diffusion v1.5 model without any fine-tune. See Appendix for more samples.

**Table 5:** Ablation study for different reward model backbones.

| Datasets | Backbone | |
|---|---|---|
| | CLIP [18] | BLIP [12] |
| ImageRewardDB [35] | 65.9 | **66.3** |
| HPD v2 [33] | 79.1 | **79.4** |
| Pick-a-Pic [9] | **67.3** | 67.1 |
| **Avg** | 70.3 | **70.5** |



**Fig. 7:** Ablation study for training data size used in VP-Score.

dataset, and then perform a top-3 selection within the these generated images by the corresponding reward models. Finally, three human annotators rank these selected images. The results are presented at Figure 4 (a), We observed that VP-Score fine-tuned generative model's Win+Tie ratio exceeds 50% when compared to all other models, including some trained on large-scale human preference datasets like HPS v2. This suggests that, compared to other human preference reward model, VP can serve as a reliable and competitive reward model for fine-tuning generative models to produce outputs closer to human preferences. This further underscores the effectiveness and competitiveness of VisionPrefer. The corresponding qualitative results shown at Figure 5 demonstrate that VisionPrefer fine-tuned generative model can generate images that are more aligned to text and with higher fidelity and avoid toxic contents.

**DPO Results.** We conduct a human preference study using the same procedure as PPO experiments, and the results are presented at Figure 4 (b), We found that the Win+Tie ratio of the generative model optimized on our VisionPrefer, when compared to the other three large-scale human datasets, exceeds 50%, substantiating the competitiveness of our VisionPrefer against human-annotated preference data. We show the qualitative results in Figure 6. The results indicate that fine-tuning the generative model directly on our VisionPrefer using DPO yields performance comparable to that of fine-tuning the generative model on large-scale human-annotated preference dataset (e.g., Pick-a-Pick). Specifically, the generated results are more aligned with human preferences, exhibiting increased visual detail, better conformity to input prompts. These experimental outcomes collectively affirm the efficacy of using preference data generated by MLLMs.

### 4.3   Ablation Study

**Reward Model Backbone.** VP-Score adopts BLIP [12] as the backbone, which may raise curiosity about how well BLIP compares to CLIP [18]. We employed these tow models as the backbone for our reward model and explored their effectiveness on our VisionPrefer. The results are summarized in Table 5, where we observed that the performance of BLIP surpassed that of CLIP and this conclusion aligns with the findings on human preference datasets [35].
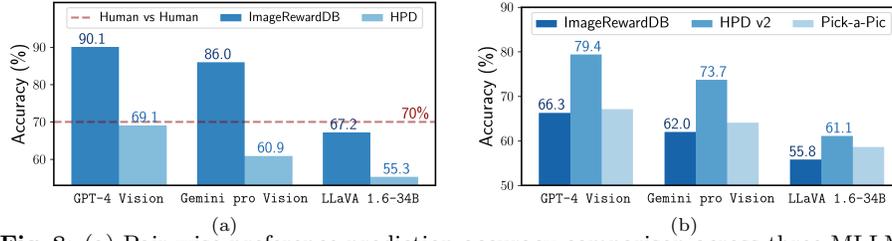**Training Data Size.** To investigate the effect of training dataset sizes on the performance of the VP-Score, comparative experiments are conducted. The resuklts are presented at Figure 7. We observed that as the training data increased, VP-Score's prediction accuracy gradually improved. This indicates that models trained on our VisionPrefer exhibit strong performance scalability, implying that more training data leads to further performance improvements. In the future, we plan to further increase the volume of data in our VisionPrefer and explore whether models trained on our dataset can outperform all these trained on human-annotated datasets. This endeavor holds significant promise and interest.

## 5   Analysis

### 5.1   Which MLLM is the Best Annotator?

The annotation of VisionPrefer heavily relies on `GPT-4 Vision`. Although many researchers pointed out that `GPT-4 Vision` capable of providing meticulous judgments and feedback [5, 32], we still concern whether the `GPT-4 Vision` preferences are qualified. We then conduct a probing experiment by utilizing different MLLMs, `GPT-4 Vision`, `Gemini pro Vision` and `LLaVA 1.6-34B`, to provide their preference on two existing human-preference datasets (HPD [34] and ImageRewardDB [35]). The corresponding pair-wise preference prediction accuracy is shown in Figure 8 (a). We observed that the accuracy of `GPT-4 Vision` surpasses that of both `LLaVA 1.6-34B` and `Gemini pro Vision` on both datasets, achieving accuracy rates exceeding or approaching 70%, and `LLaVA 1.6-34B` notably scoring significantly lower than the former two. According to previous research [5, 33], the agreement rate between qualified human annotators is also around 70% (65.3% for ImageRewardDB and 78.1% for HPD). Therefore, the probing experiment validates that `GPT-4 Vision` can be a well human-aligned annotator, thus ensure the quality and reliability of our VisionPrefer.

To further validate the efficacy on preference annotation ability of `GPT-4 Vision`, we utilize `Gemini pro Vision` and `LLaVA 1.6-34B` to collect similar amount of data (1.2 M pair-wise preference choices) following the same collection pipeline described in Section 3. Then we train the corresponding reward model on these two datasets and show the performance in Figure 8 (b). As we can see, consistent with the aforementioned conclusion, the testing accuracy of the reward model trained on data annotated by `GPT-4 Vision` exhibits the highest performance, followed by `Gemini pro Vision`. This demonstrates that `GPT-4 Vision` is currently the most proficient annotator for tetx-to-image generation.
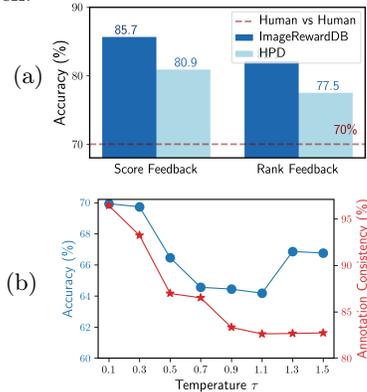
**Fig. 8:** (a) Pair-wise preference prediction accuracy comparison across three MLLMs on two human-preference datasets (b) Results of preference prediction accuracy for reward models trained on preference datasets annotated by different AI annotators.
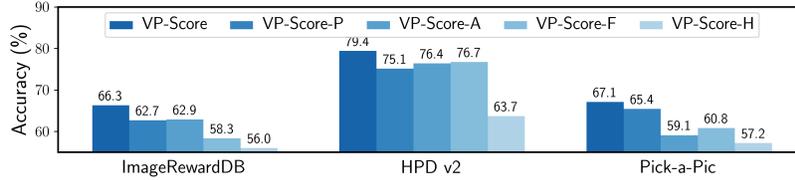
## 5.2   Encouraging GPT-4 Vision for Enhanced Annotations.

**Prompt Manner.** As described in Section 3, during the construction of Vision-Prefer, we encourage `GPT-4 Vision` to directly output scores for various aspects (e.g., aesthetic) of each image (denoted as score feedback). Another straightforward prompting manner (denoted as rank feedback) is encourage `GPT-4 Vision` to directly provide a ranking of images ($\alpha$ and $\beta$) in a certain aspect (i.e., $\alpha \succ \beta$, $\beta \succ \alpha$, or $\alpha = \beta$). It is interesting to explore which prompting manner is best suit for AI Annotators. We randomly sampled 1,000 samples from ImageRewardDB and 500 samples from HPD, and utilized the two aforementioned prompt manners to ascertain `GPT-4 Vision`'s annotations. The results are presented at Figure 9 (a), we observe that in both datasets, the accuracy achieved using score feedback is higher than that achieved using rank feedback.

**Temperature $\tau$.** Temperature $\tau$ is a hyperparameter used in multimodal large language models (e.g., `GPT-4 Vision`) to control the randomness and creativity of the generated results. A lower value of the temperature parameter will lead to a more predictable and deterministic output, while a higher value will produce a more random and surprising output. We investigate the influence of the variation in $\tau$ on both the accuracy of annotation and annotation consistency (where the same input yields identical annotation results). The results are shown in Figure 9 (b), we observe a decrease in accuracy as $\tau$ increases, indicating that lower values of $\tau$ should be set when conducting preference annotations. Furthermore, as $\tau$ increases, annotation consistency continues to decline, which is sensible because more



**Fig. 9:** (a) Preference prediction accuracy for score feedback and ranking feedback. (b) Visualization depicting the variation of annotation accuracy and consistency with changes in temperature $\tau$.

randomness in the results leads to different preference outcomes for identical inputs over time.

**Fig. 10:** Preference prediction accuracy among reward models trained on different aspects of preference data in VisionPrefer. We find VP-Score surpasses all reward models trained using preference data from individual evaluation aspects.



**Fig. 11:** Fine-grained feedback enables our model (denoted as VP-Score) to generate results that better align with the input prompt. See Appendix for more samples.

### 5.3    Fine-Grained Feedback Leads to Better Results.

**Better Reward Modeling**. In our previous experiments, we used the average score of each sample across four evaluation aspects as the final preference score for modeling reward or optimizing generation models. Here, we explore the impact of separately modeling the four evaluation aspects. We first train four reward models on these four different aspect in VisionPrefer, namely VP-Score-P, VP-Score-A, VP-Score-F and VP-Score-H, respectively. The corresponding preference accuracy are presented at Figure 10. We can observe that the accuracy of reward models individually trained using a single aspect preference data is consistently lower than VP-Score, which validates the effectiveness of our approach in designing four evaluation aspects to model the preference level.

**Better Prompt-Following.** We found fine-grained preference data enables our fine-tuned model to generate images that better adhere to the input prompt. For instance, as shown in Figure 11, we find the top-3 sampled images generated from our fine-tuned model all fulfill the requirement of "holding" as specified in the prompt. In contrast, among the baseline models, only HPS v2 achieves this.

**More Aesthetically Pleasing.** Fine-grained data enhances the visual appeal and vividness of images generated by our model. As shown in Figure 12, our results exhibit enhanced luminosity, dynamic sensation, and increased detail, aligning more closely with human aesthetic preferences.

**Reduce Image Distortion.** One aspect measured within our fine-grained data is the fidelity score of images, which ensures images remain undistorted and accurately represent the intended subject. We demonstrate that scoring at this granularity level enables our model to generate more precise results. In Figure 13, we observe that our model achieves the highest accuracy in generating hand images, while other comparative models exhibit varying degrees of deformation.

**Fig. 12:** Fine-grained feedback enhances the aesthetic and vividness of the our results (denoted as VP-Score). SD 1.5 denotes the Stable Diffusion v1.5 model without any fine-tune. See Appendix for more samples.



**Fig. 13:** Fine-grained feedback help to reduce image distortion in our generative results (denoted as VP-Score). See Appendix for more samples.

**Enhance Image Safety.** We employ unsafe prompts provided in [36] to generate 1 K images and utilize the built-in NSFW detector in the Diffusion library to quantify the frequency of generating harmful content. We find that VP-Score fine-tuned generative model's NSFW ratio (4.4%) is significantly lower than that of other models, being 5× lower than HPS v2 fine-tuned generative model (21.1%) and 4.8× times lower than PickScore fine-tuned generative model (22.3%). This indicates that the harmlessness evaluation incorporated into our fine-grained preference scoring mechanism effectively reduces the generation of harmful outputs by the generative model. More related details can be found in Appendix.

## 6   Conclusion

In this paper, we explore utilize AI annotators to construct a large-scale high-quality feedback dataset, VisionPrefer, for diffusion models alignment and refining. Costly experiments conducted across various experimental settings have validated the efficacy of our VisionPrefer. This also represents a comprehensive and substantial endeavor by RLAIF in the realm of visual generative models,

demonstrating the effectiveness of utilizing AI-synthesized data for aligning visual generative models.

**Limitations and Future Direction.** VisionPrefer provides three types of feedback data for each item, but we have not yet utilized the textual explanations data, which will be a interesting direction for our future exploration.

# References

1. Bai, Y., Jones, A., Ndousse, K., Askell, A., Chen, A., DasSarma, N., Drain, D., Fort, S., Ganguli, D., Henighan, T., et al.: Training a helpful and harmless assistant with reinforcement learning from human feedback. arXiv preprint arXiv:2204.05862 (2022)

2. Bai, Y., Kadavath, S., Kundu, S., Askell, A., Kernion, J., Jones, A., Chen, A., Goldie, A., Mirhoseini, A., McKinnon, C., Chen, C., Olsson, C., Olah, C., Hernandez, D., Drain, D., Ganguli, D., Li, D., Tran-Johnson, E., Perez, E., Kerr, J., Mueller, J., Ladish, J., Landau, J., Ndousse, K., Lukosuite, K., Lovitt, L., Sellitto, M., Elhage, N., Schiefer, N., Mercado, N., DasSarma, N., Lasenby, R., Larson, R., Ringer, S., Johnston, S., Kravec, S., Showk, S.E., Fort, S., Lanham, T., Telleen-Lawton, T., Conerly, T., Henighan, T., Hume, T., Bowman, S.R., Hatfield-Dodds, Z., Mann, B., Amodei, D., Joseph, N., McCandlish, S., Brown, T., Kaplan, J.: Constitutional ai: Harmlessness from ai feedback (2022)

3. Black, K., Janner, M., Du, Y., Kostrikov, I., Levine, S.: Training diffusion models with reinforcement learning. arXiv preprint arXiv:2305.13301 (2023)

4. Clark, K., Vicol, P., Swersky, K., Fleet, D.J.: Directly fine-tuning diffusion models on differentiable rewards. arXiv preprint arXiv:2309.17400 (2023)

5. Cui, G., Yuan, L., Ding, N., Yao, G., Zhu, W., Ni, Y., Xie, G., Liu, Z., Sun, M.: Ultrafeedback: Boosting language models with high-quality feedback. arXiv preprint arXiv:2310.01377 (2023)

6. Fan, Y., Watkins, O., Du, Y., Liu, H., Ryu, M., Boutilier, C., Abbeel, P., Ghavamzadeh, M., Lee, K., Lee, K.: Dpok: Reinforcement learning for fine-tuning text-to-image diffusion models. arXiv preprint arXiv:2305.16381 (2023)

7. Feng, S.Y., Gangal, V., Wei, J., Chandar, S., Vosoughi, S., Mitamura, T., Hovy, E.: A survey of data augmentation approaches for NLP. In: Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021. pp. 968–988. Association for Computational Linguistics, Online (Aug 2021). `https://doi.org/10.18653/v1/2021.findings-acl.84`, `https://aclanthology.org/2021.findings-acl.84`

8. Hanu, L., Unitary team: Detoxify: Toxic Comment Classification with Pytorch Lightning and Transformers (2020), `https://github.com/unitaryai/detoxify`

9. Kirstain, Y., Polyak, A., Singer, U., Matiana, S., Penna, J., Levy, O.: Pick-a-pic: An open dataset of user preferences for text-to-image generation. arXiv preprint arXiv:2305.01569 (2023)

10. Kwon, M., Xie, S.M., Bullard, K., Sadigh, D.: Reward design with language models. In: The Eleventh International Conference on Learning Representations (2022)

11. Lee, K., Liu, H., Ryu, M., Watkins, O., Du, Y., Boutilier, C., Abbeel, P., Ghavamzadeh, M., Gu, S.S.: Aligning text-to-image models using human feedback. arXiv preprint arXiv:2302.12192 (2023)

12. Li, J., Li, D., Xiong, C., Hoi, S.: Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In: International Conference on Machine Learning. pp. 12888–12900. PMLR (2022)

13. Liu, H., Li, C., Li, Y., Lee, Y.J.: Improved baselines with visual instruction tuning (2023)

14. Madaan, A., Tandon, N., Gupta, P., Hallinan, S., Gao, L., Wiegreffe, S., Alon, U., Dziri, N., Prabhumoye, S., Yang, Y., et al.: Self-refine: Iterative refinement with self-feedback. arXiv preprint arXiv:2303.17651 (2023)

15. Meng, Y., Michalski, M., Huang, J., Zhang, Y., Abdelzaher, T., Han, J.: Tuning language models as training data generators for augmentation-enhanced few-shot

learning. In: International Conference on Machine Learning. pp. 24457–24477. PMLR (2023)

16. Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al.: Training language models to follow instructions with human feedback. Advances in Neural Information Processing Systems **35**, 27730–27744 (2022)

17. Prabhudesai, M., Goyal, A., Pathak, D., Fragkiadaki, K.: Aligning text-to-image diffusion models with reward backpropagation. arXiv preprint arXiv:2310.03739 (2023)

18. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., Sutskever, I.: Learning Transferable Visual Models From Natural Language Supervision. In: ICML (2021)

19. Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., Chen, M.: Hierarchical text-conditional image generation with clip latents. arXiv preprint arXiv:2204.06125 **1**(2), 3 (2022)

20. Roit, P., Ferret, J., Shani, L., Aharoni, R., Cideron, G., Dadashi, R., Geist, M., Girgin, S., Hussenot, L., Keller, O., et al.: Factually consistent summarization via reinforcement learning with textual entailment feedback. arXiv preprint arXiv:2306.00186 (2023)

21. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 10684–10695 (2022)

22. Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E.L., Ghasemipour, K., Gontijo Lopes, R., Karagol Ayan, B., Salimans, T., et al.: Photorealistic text-to-image diffusion models with deep language understanding. Advances in Neural Information Processing Systems **35**, 36479–36494 (2022)

23. Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E.L., Ghasemipour, K., Gontijo Lopes, R., Karagol Ayan, B., Salimans, T., et al.: Photorealistic text-to-image diffusion models with deep language understanding. Advances in Neural Information Processing Systems **35**, 36479–36494 (2022)

24. Schuhmann, C., Beaumont, R., Vencu, R., Gordon, C., Wightman, R., Cherti, M., Coombes, T., Katta, A., Mullis, C., Wortsman, M., et al.: Laion-5b: An open large-scale dataset for training next generation image-text models. arXiv preprint arXiv:2210.08402 (2022)

25. Song, J., Meng, C., Ermon, S.: Denoising diffusion implicit models. arXiv preprint arXiv:2010.02502 (2020)

26. Stiennon, N., Ouyang, L., Wu, J., Ziegler, D., Lowe, R., Voss, C., Radford, A., Amodei, D., Christiano, P.F.: Learning to summarize with human feedback. Advances in Neural Information Processing Systems **33**, 3008–3021 (2020)

27. Tang, Z., Rybin, D., Chang, T.H.: Zeroth-order optimization meets human feedback: Provable learning via ranking oracles. arXiv preprint arXiv:2303.03751 (2023)

28. Wang, X., Wei, J., Schuurmans, D., Le, Q.V., Chi, E.H., Narang, S., Chowdhery, A., Zhou, D.: Self-consistency improves chain of thought reasoning in language models. In: The Eleventh International Conference on Learning Representations (2022)

29. Wang, Z.J., Montoya, E., Munechika, D., Yang, H., Hoover, B., Chau, D.H.: Diffusiondb: A large-scale prompt gallery dataset for text-to-image generative models. arXiv preprint arXiv:2210.14896 (2022)

30. Wang, Z., Yu, A.W., Firat, O., Cao, Y.: Towards zero-label language learning. arXiv preprint arXiv:2109.09193 (2021)

31. Wen, S., Fang, G., Zhang, R., Gao, P., Dong, H., Metaxas, D.: Improving compositional text-to-image generation with large vision-language models. arXiv preprint arXiv:2310.06311 (2023)
32. Wu, T., Yang, G., Li, Z., Zhang, K., Liu, Z., Guibas, L., Lin, D., Wetzstein, G.: Gpt-4v (ision) is a human-aligned evaluator for text-to-3d generation. arXiv preprint arXiv:2401.04092 (2024)
33. Wu, X., Hao, Y., Sun, K., Chen, Y., Zhu, F., Zhao, R., Li, H.: Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis. arXiv preprint arXiv:2306.09341 (2023)
34. Wu, X., Sun, K., Zhu, F., Zhao, R., Li, H.: Human preference score: Better aligning text-to-image models with human preference. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 2096–2105 (2023)
35. Xu, J., Liu, X., Wu, Y., Tong, Y., Li, Q., Ding, M., Tang, J., Dong, Y.: Imagereward: Learning and evaluating human preferences for text-to-image generation. arXiv preprint arXiv:2304.05977 (2023)
36. Yang, K., Tao, J., Lyu, J., Ge, C., Chen, J., Li, Q., Shen, W., Zhu, X., Li, X.: Using human feedback to fine-tune diffusion models without any reward model. arXiv preprint arXiv:2311.13231 (2023)
37. Yang, K., Klein, D., Celikyilmaz, A., Peng, N., Tian, Y.: Rlcd: Reinforcement learning from contrast distillation for language model alignment (2023)

# Multimodal Large Language Model is a Human-Aligned Annotator for Text-to-Image Generation—Appendix

We organize our appendix as follows:

– In **Section A**, we augment the primary outcomes of both the PPO and DPO experiments by incorporating two new validation benchmarks, alongside additional human preference studies and qualitative results.

– In **Section B**, we further validate the efficacy of our finely-grained preference annotation approach through comprehensive ablation studies, augmented human preference evaluations, and extended visualization results.

– In **Section C**, we present the statistical analysis of VisionPrefer, revealing that preference labels annotated by MLLMs demonstrate characteristics akin to those found in human-annotated preference datasets.

– In **Section D**, we provide training details of the experiments discussed in the main text.

– In **Section E**, we analyze the annotation costs of constructing VisionPrefer and demonstrate that employing MLLMs as annotators not only achieves results aligned with human annotators but also significantly reduces labor resources and time costs.

– In **Section F**, we provide the corresponding prompt templates used in employing `GPT-4` to polish existing prompts and employing `GPT-4 Vision` to generate preference annotations.

# A    Additional Main Results

In this section, we present additional quantitative and qualitative results to further validate the efficacy of VisionPrefer and VP-Score. Beyond DiffusionDB [29] test-set, we broaden our validation efforts by incorporating two additional benchmarks ReFL test-set [35] (250 examples) and HPD v2 test-set [33] (500 examples), into our assessment framework.
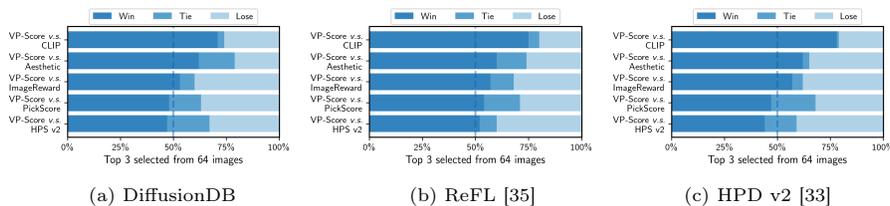
## A.1    Additional PPO Results

**Quantitative Results.** We detail the win count and win rates of generative models optimized with different reward models against `Stable Diffusion v1.5` baseline on Table 1. These results reveals that VP-Score consistently outperforms ImageReward and demonstrates competitive performance against two another reward models (PickScore and HPS v2), which are trained on nearly a million human-annotated preference data. For instance, VP-Score achieves the highest win rate in the ReFL test-set, surpassing all other reward models. In DiffusionDB and HPD v2 test-sets, VP-Score secured the second-best and third-best performances, with marginal differences from HPS v2 and PickScore.

Figure 1 displays the win rates of the generative model optimized with VP-Score against to generative models optimized with alternative reward models across three test benchmarks. It is observed that the generative model optimized with VP-Score achieves a win + tie rate exceeding 50% when benchmarked against others, underscoring the effectiveness of VP-Score

**Qualitative Results.** Additional qualitative results are presented in Figure 2. When compared to the `Stable Diffusion v1.5` baseline generative model (simplified as SD 1.5), the generative model optimized with the guidance of VP-Score is capable of producing images with richer details and more closely aligned with the given prompts. This showcases competitive performance with generative models optimized through reward models such as HPS v2 and PickScore, demonstrating the effectiveness of VP-Score.

**Table 1:** Human evaluation study on win count and win rate of generative models optimized with different reward models, benchmarked against the `Stable Diffusion v1.5` baseline. Compared to other reward models, VP-Score exhibits competitive performance. The best results are highlighted in bold, while the second-best results are underlined.

| Reward Model | DiffusionDB [29] | | ReFL [35] | | HPD v2 [33] | |
|---|---|---|---|---|---|---|
| | #Win | WinRate | #Win | WinRate | #Win | WinRate |
| CLIP [18] | 267 | 54.09 | 137 | 52.05 | 270 | 53.31 |
| Aesthetic [24] | 280 | 56.71 | 144 | 53.93 | 283 | 54.77 |
| ImageReward [35] | 281 | 56.93 | 153 | 55.81 | 291 | 56.38 |
| PickScore [9] | 286 | 57.87 | 164 | 56.66 | 298 | _57.87_ |
| HPS v2 [33] | 291 | **58.21** | 171 | _56.87_ | 287 | **57.89** |
| VP-Score (Ours) | 329 | _57.98_ | 177 | **57.09** | 295 | 57.80 |

**Fig. 1:** Win rates of the generative model optimized with VP-Score compared to generative models optimized other reward models on three test benchmarks. VP-Score shows a competitive performance.

## A.2    Additional DPO Results

**Quantitative Results.** We detail the win count and win rates of generative models fine-tuned on different preference datasets against `Stable Diffusion v1.5` baseline on Table 2. Through our experimental results, it is observed that within the DPO framework, our VisionPrefer demonstrates competitive effectiveness when compared to human-annotated preference datasets. For instance, as shown in the Table 2, generative models trained on VisionPrefer achieves the highest win rate in both ReFL and HPS v2 test-set, surpassing these two generative models trained on HPD and Pick-a-Pic. Moreover, within DiffusionDB test-set, generative models trained on VisionPrefer secures a close second-best performance, achieving results comparable to generative models trained on HPD and Pick-a-Pic.

Figure 3 showcases that generative models trained with VisionPrefer exhibit competitive performance relative to models trained on alternative human-annotated preference datasets, achieving a win + tie rate that exceeds the 50% threshold across three testing benchmarks. This evidence highlights the efficacy of VisionPrefer.

**Qualitative Results.** We provide more qualitative results in Figure 4. The visualization results indicate that generative models fine-tuned on VisionPrefer are capable of producing outputs that closely match human preferences, generating images that are more appealing to human users.

**Table 2:** Human evaluation on generative models optimized with different preference datasets in DPO experiments. The best results are highlighted in bold, while the second-best results are underlined.

| Preference Datasets | DiffusionDB [29] | | ReFL [35] | | HPS v2 [33] | |
|---|---|---|---|---|---|---|
| | #Win | WinRate | #Win | WinRate | #Win | WinRate |
| ImageReward [35] | 253 | 54.31 | 144 | 53.87 | 281 | 55.01 |
| HPD [33] | 266 | 57.08 | 149 | 55.71 | 278 | 54.49 |
| Pick-a-Pic [9] | 277 | **59.43** | 156 | 58.33 | 297 | 58.23 |
| VisionPrefer (Ours) | 275 | 59.03 | 158 | **59.17** | 303 | **59.44** |

**male wizard, brown hair, green robes, glasses, D&D, painted fantasy character portrait, highly detailed, digital painting, artstation, concept art, sharp focus, illustration, art by artgerm and greg rutkowski and alphonse mucha**

**A caracal is eating dumplings with a restaurant in the background, presented in a whimsical, Pixar-style digital painting**
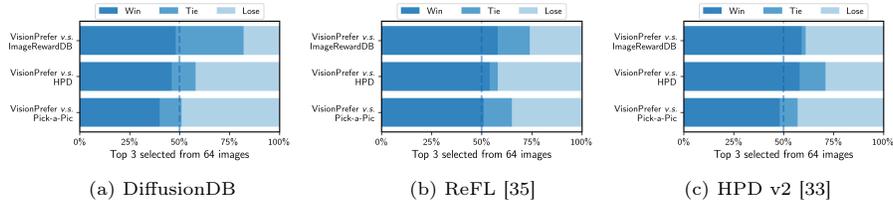
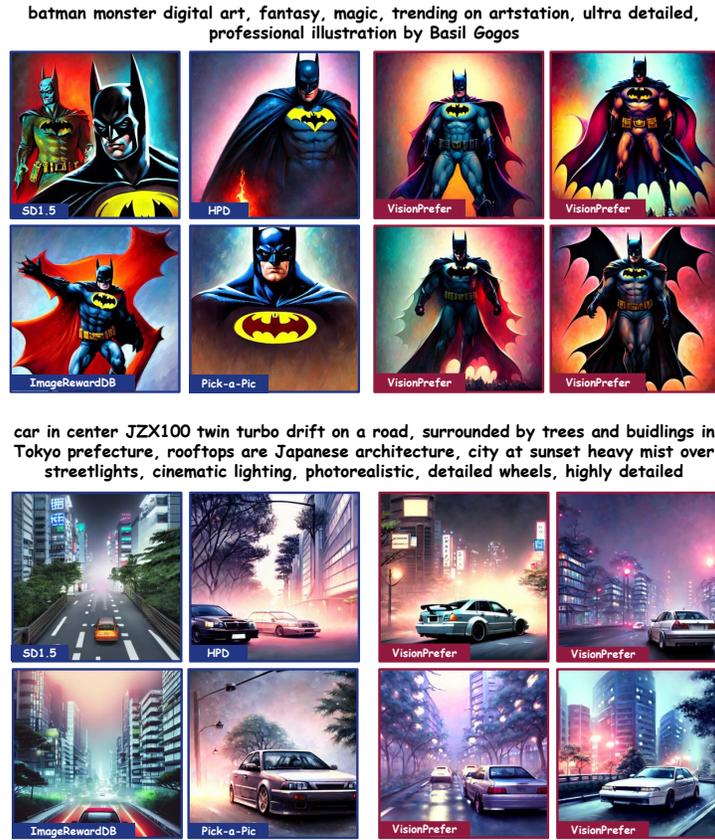**Portrait of a woman in the style of impressionism by Patrice Murciano**

**The classical Roman church, highly detailed, artstation, strong contrast of light and shadow, neon colors, sharp focus 4K UHD image**

**Fig. 2:** Qualitative comparison between text-to-image generative model optimized with the guidance of VP-Score and other reward models. SD 1.5 denotes the `Stable Diffusion v1.5` model without any fine-tune.

(a) DiffusionDB        (b) ReFL [35]        (c) HPD v2 [33]

**Fig. 3:** Win rates of text-to-image generative model trained on VisionPrefer compared to models trained on other human-annotation preference datasets across three test benchmarks. VisionPrefer shows a competitive performance.

**batman monster digital art, fantasy, magic, trending on artstation, ultra detailed, professional illustration by Basil Gogos**



**car in center JZX100 twin turbo drift on a road, surrounded by trees and buidlings in Tokyo prefecture, rooftops are Japanese architecture, city at sunset heavy mist over streetlights, cinematic lighting, photorealistic, detailed wheels, highly detailed**



**Fig. 4:** Qualitative comparison between generative model trained on VisionPrefer and other human-annotated preference datasets. SD 1.5 denotes the `Stable Diffusion v1.5` model without any fine-tune.

## B      Efficacy of Fine-Grained Feedback

In this section, we present additional ablation studies and qualitative results to validate the efficacy of fine-grained feedback design.

### B.1      Better Prompt-Following.

To further substantiate that "Prompt-Following" rating labels in VisionPrefer can enhance the ability of fine-tuned models to generate images that more accurately align with the input prompts, we removed the "Prompt-Following" rating labels from the VisionPrefer, preserving labels for the other three aspects, and train a reward model, which we named VP-Score$^\dagger$. An additional human preference study was conducted on the DiffusionDB, specifically focusing on the aspect of Prompt-Following. The outcomes, presented in Table 3, indicate two key findings: Firstly, VP-Score exhibits competitive performance in comparison to HPS v2. Secondly, the efficacy of VP-Score$^\dagger$ experiences a notable decline with the omission of Prompt-Following rating labels. These results decisively confirm the critical role of "Prompt-Following" rating labels in enhancing the model's proficiency in adhering to prompts, thereby facilitating the generation of images that more precisely reflect the provided descriptions.

Besides, we provide additional visualization results in Figure 6 to validate that generative models guided by VP-Score are capable of producing images that more closely adhere to the descriptions provided in the prompts.

### B.2      More Aesthetically Pleasing.

Similar to the last section, we removed the "Aesthetic" rating labels from the VisionPrefer and trained a reward model named VP-Score$^\clubsuit$. Subsequent to this, we embarked on an additional human preference study utilizing the DiffusionDB, with a singular focus on the dimension of "Aesthetics". The findings, elucidated in Table 4, revealed that VisionPrefer achieved the best performance, while the exclusion of aesthetic labels markedly diminished the operational efficiency of VP-Score$^\clubsuit$. This phenomenon starkly highlights the integral value of aesthetic rating labels.

Further, we showcase additional visual outcomes in Figure 7. Our observations indicate that generative models refined under the auspices of VP-Score manifest the capacity to engender imagery replete with more vibrant detail and sophisticated interplays of light and shadow.

### B.3      Reduce Image Distortion.

To ascertain the impact of "Fidelity" rating labels, we excised these labels from the VisionPrefer and subsequently trained a reward model VP-Score$^\diamond$. A human preference study, concentrated solely on the "Fidelity" aspect, is documented in Table 5. This study utilized the "anything" prompts delineated in [36], encompassing 442 prompts, as the evaluation benchmark. The outcomes illustrate that

generative models guided by our VP-Score manifest competitive performance. In contrast, VP-Score$^{\diamond}$ exhibits a discernible performance decrement relative to VP-Score.

Additionally, the visualization results showcased in Figure 8 demonstrate that models optimized under the guidance of VP-Score excel in producing images with diminished distortion, e.g., less distortion of human hands. Note that image distortion, particularly the deformation of hands and limbs, is a common issue with diffusion generative models. Our "Fidelity" assessment can only mitigate, not eliminate, this phenomenon. Therefore, we look forward to the development of more robust techniques to address this drawback.

### B.4    Enhance Image Safety.

Similarly, we removed the "Harmlessness" labels from the VisionPrefer and trained a corresponding reward model named VP-Score$^{\spadesuit}$. Then we employ unsafe prompts provided in [36] to generate 1K images and utilize the built-in NSFW detector in the diffusion library[1] to quantify the frequency of generating harmful content. Detailed results is presented in Figure 5. We find that VP-Score$^{\spadesuit}$, trained without the "Harmlessness rating" labels, exhibited a significant increase in the NSFW ratio compared to the original VP-Score (4.4% to 20.2%). This further underscores the importance of "Harmlessness" labels.

**Table 3:** Human evaluation study on the aspect of "Prompt-Following". The best results are highlighted in bold, while the second-best results are underlined.
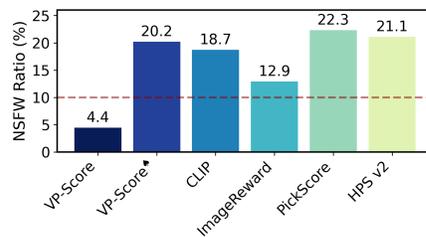
| Reward Model | DiffusionDB [29] | |
|---|---|---|
| | #Win | WinRate |
| PickScore [9] | 311 | 57.24 |
| HPS v2 [33] | 316 | **58.27** |
| VP-Score$^{\dagger}$ | 307 | 56.70 |
| VP-Score | 315 | <u>58.07</u> |

**Table 4:** Human evaluation study on the aspect of "Aesthetic". The best results are highlighted in bold, while the second-best results are denoted with an underline.

| Reward Model | DiffusionDB [29] | |
|---|---|---|
| | #Win | WinRate |
| PickScore [9] | 283 | <u>55.40</u> |
| HPS v2 [33] | 281 | 55.01 |
| VP-Score$^{\clubsuit}$ | 275 | 53.82 |
| VP-Score | 286 | **55.96** |

**Table 5:** Human evaluation study on the aspect of "Fidelity". The best results are highlighted in bold, while the second-best results are denoted with an underline.

| Reward Model | Anything Prompts [36] | |
|---|---|---|
| | #Win | WinRate |
| PickScore [9] | 227 | 51.17 |
| HPS v2 [33] | 232 | **52.33** |
| VP-Score$^{\diamond}$ | 224 | 50.51 |
| VP-Score | 231 | <u>52.20</u> |



**Fig. 5:** Fine-grained feedback make generation more safety.

---

[1] https://github.com/huggingface/diffusers

**A painting depicts an oak tree with a human face resembling an old bearded man, crafted from the tree's bark**



**A ghost pirate aboard a pirate ship in spooky fog under moonlight**



**Fig. 6:** Fine-grained feedback enhances the alignment of generated content with the input prompts. For instance, in the first column of the figure, only the generative model optimized under the guidance of VP-Score accurately produces a face that adheres to the description of being 'crafted from the tree's bark.' In the second column, solely the VP-Score-guided generative model successfully constructs the image of a pirate, whereas the other models merely generate images of pirate ships. SD 1.5 denotes the `Stable Diffusion v1.5` model without any fine-tune.

**portrait of a gorgeous punk vampire girl, elegant, digital painting, highly detailed, artstation, concept art, smooth, sharp focus, illustration, art by artgerm and greg rutkowski and alphonse mucha**



**black dafne keen yoruba, intricate, elegant, highly detailed, digital painting, artstation, concept art, smooth, sharp focus, illustration, d&d, art by rutkowski, orientalism, bouguereau**



**Fig. 7:** Fine-grained feedback enhances the vividness and richness of detail in generated content. SD 1.5 denotes the `Stable Diffusion v1.5` model without any fine-tune.

**1girl, bangs, blunt bangs, bowl, brown hair, cherry blossoms, closed eyes, closed mouth, facing viewer, floral print, flower, green kimono, grey hair, hair flower, hair ornament, hands …**



**Marla Singer is depicted smoking in a setting reminiscent of Blade Runner, highly detailed, digital painting, artstation, concept art, smooth, sharp focus**



**Fig. 8:** Fine-grained feedback reduce the image distortion. SD 1.5 denotes the `Stable Diffusion v1.5` model without any fine-tune.

## C   Statistics of VisionPrefer

In this section, we provide more details of VisionPrefer. We encourage readers to delve into Section C.3, where we analyze the characteristics of preferences generated by `GPT-4 Vision`. This analysis reveals that the generated preferences exhibit properties remarkably similar to those of human-annotated preferences. Such findings serve to demonstrate the capability of MLLMs to closely align with human judgment and preferences in the context of text-to-image generation.

### C.1   Prompts.

A key step in VisionPrefer contruction pipeline is utilizing `GPT-4` to polish the existing prompt benchmarks. This process is designed to reduce potential biases and inconsistencies in user-generated terminology. We present examples of original prompts alongside their polished counterparts at Table 6, and quantitatively illustrates the frequency distribution of certain stylistic words and conflicting prompts at Figure 9. Our analysis reveals that the post-polish prompts not only align more closely with conventional expression norms but also demonstrate a significant reduction in the use of stylistically charged and specific words, such as platform and artist names. Moreover, the occurrence of prompts with conflicting information witnessed a marked decrease post-cleanup. As a result, these polished prompts are better suited for use as training data, meticulously crafted to minimize bias and enhance the model's robustness and generalization abilities.
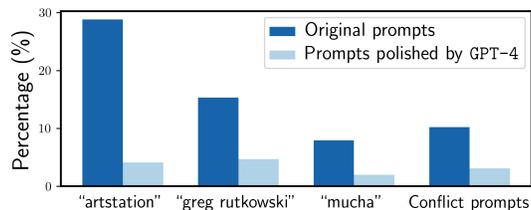
**Table 6:** Examples of prompts polished by `GPT-4`. Certain style words are underlined.

| Prompts from DiffusionDB [29] | Prompts cleaned by `GPT-4` |
|---|---|
| cyberpunk neon gorilla skull, by weta fx, by wlop, majestic look, trending on artstation. | Neon gorilla skull in a cyberpunk style. |
| highly detailed digital painting, black male anthro - lynx, human with head of lynx, with hair like fabio, facial scar, hairy masculine gigachad, muscular, wearing kilt and gold armbands, fur texture, lounging on bed aboard the nostromo, trending on artstation, romance novel. | A digital painting depicts a black male anthropomorphic lynx with long hair, a facial scar, and a muscular build, wearing a kilt and gold armbands, lounging on a bed. |

### C.2   Images.

Within VisionPrefer, images are generated by employing four state-of-the-art text-to-image generative models. These models are ranked as the top four on the Hugging Face leaderboard, specifically: `Stable Diffusion v1-5`[2], `Stable Diffusion`

---

[2] `https://huggingface.co/runwayml/stable-diffusion-v1-5`

**Fig. 9:** Frequencies of certain style words and conflict prompt. Confliction is judged by `GPT-4`.

`2.1`[3], `Dreamlike Photoreal 2.05`[4], `Stable Diffusion XL`[5]. Detailed descriptions of these models and the distribution of images generated by each within our dataset are methodically outlined in Table 7. For illustrative purposes, Figure 10 showcases representative images produced by each of these models.

**Table 7:** Image sources of VisionPrefer.

| Source | Type | Resolution | Proportion |
|---|---|---|---|
| Stable Diffusion v1-5 | Diffusion | 512×512 | 26.3 % |
| Stable Diffusion 2.1 | Diffusion | 768×768 | 24.8 % |
| Dreamlike Photoreal 2.05 | Diffusion | 768×768 | 25.7 % |
| Stable Diffusion XL | Diffusion | 1,024×1,024 | 23.2 % |

### C.3    Preferences.

All preferences in VisionPrefer were generated using `GPT-4 Vision` API. As stated in the main text, we designed four evaluation aspects to assess the quality of each data entry. For each aspect, we individually invoked the `GPT-4 Vision` API to generate the corresponding preferences, the corresponding prompts regarding to these four aspects can be found depicted in Section F.2. Please refer to Table 8 to see some annotation examples.

We illustrate the distribution of `GPT-4 Vision`'s ratings across four distinct aspects within VisionPrefer in Figure 11. Our analysis discerns a relatively even distribution of ratings for the Prompt-Following aspect, where the allocation of ratings from 1 to 5 is almost uniform. Conversely, in the domains of Fidelity and Harmlessness, `GPT-4 Vision` exhibits a propensity towards assigning the highest rating of 5 to a predominant share of the samples. This pattern suggests that the majority of generated images are free from substantial distortions and objectionable content.

VisionPrefer offers a unique opportunity to leverage `GPT-4 Vision`' preferences for unbiased analysis. A critical step in the construction of VisionPrefer is
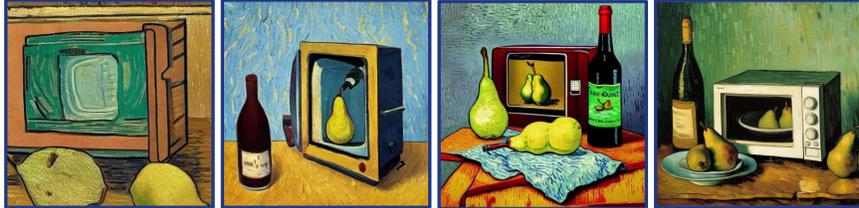
---

[3] `https://huggingface.co/stabilityai/stable-diffusion-2-1`

[4] `https://huggingface.co/dreamlike-art/dreamlike-photoreal-2.0`

[5] `https://huggingface.co/stabilityai/stable-diffusion-xl-base-1.0`

**a painting of a dark souls boss by greg rutkowski, dark fantasy art, high detail, trending on artstation.**



**a beautiful painting of a microwave, a pear and a bottle of wine on a table, painted by van gogh**



|            SD v1-5            |            SD 2.1            |            Dreamlike            |            SD XL            |

**Fig. 10:** Some example images in VisionPrefer. `SD` denotes `Sable Diffusion` while `Dreamlike` denotes `Dreamlike Photoreal 2.05`.



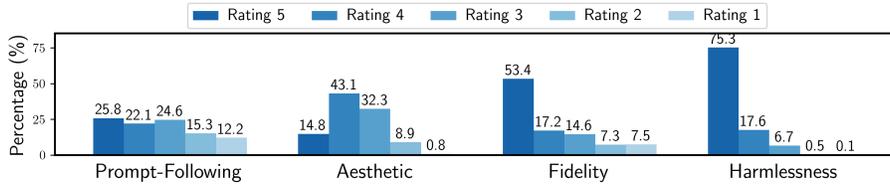**Fig. 11:** Distribution of `GPT-4 Vision`'s scoring across four aspects in VisionPrefer.



(a)                                                                    (b)
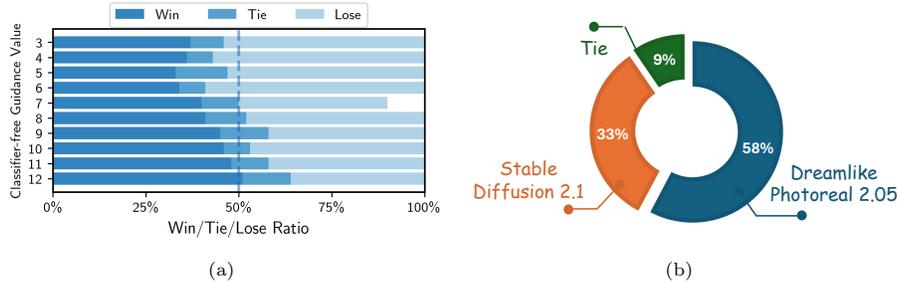
**Fig. 12:** (a) Win rate versus classifier-free guidance value for Stable Diffusion XL. (b) Preference distribution when comparing `Stable Diffusion 2.1` with `Dreamlike Photoreal 2.05`.

the random application of class-free guidance values (from 3 to 12), aiming to enhance the generalization of VisionPrefer. We analyze the impact of changing the class-free guidance values of `Stable Diffusion XL` on its performance. For each guidance value, we compute the win ratio, representing the percentage of judgments where its use led to a preferred image. We also calculate the corresponding tie and lose ratios for each scale, enabling a detailed analysis of which classifier-free guidance scales are more effective. The results are presented at Figure 12 (a), and we find higher guidance value always lead to a higher win rate, e.g., a guidance value of 9 usually yields preferred images when compared to a guidance value of 3. *This conclusion is in well agreement with the conclusions obtained in the human-annotated preference dataset [9].*

Additionally, by comparing preferences for images generated by different generative models in response to identical prompts, we are able to identify the model that is more favorably rated by `GPT-4 Vision`. For instance, considering judgments in which one image was generated by `Dreamlike Photoreal 2.05` and the other by `Stable Diffusion 2.1`, we can evaluate which model is more performant. As shown in Figure 12 (b), we find that `GPT-4 Vision` typically exhibits a preference for `Dreamlike Photoreal 2.05` over `Stable Diffusion 2.1`. *This preference aligns with the characteristics of human-annotated preference data [9], demonstrating a consistency between the MLLMs generated preferences and human judgment.*

**Table 8:** Example of annotations in VisionPrefer. In VisionPrefer, each data item includes a prompt, four images generated based on that prompt, along with the preference rating for each image across four different aspects, and their corresponding rationales.

**Prompt: minimalist summertime architecture by atey ghailan ( ( and edward hopper ) ).**

| Input Image | Prompt-Following | Aesthetic | Fidelity | Harmlessness |
|---|---|---|---|---|
|  | **Rating**: 5 **Rationale**: The image showcases a minimalist architectural style with a focus on simple geometric forms and warm colors that can be associated with summer. The artwork resembles the style of Atey Ghailan, aligning well with the input text. | **Rating**: 5 **Rationale**: The image has a clear and sharp focus, the exposure is perfect with the warm glow from the window, and the colors are rich and vibrant. The composition is masterful, with a clear focus on the minimalist architecture and the figure, creating a sense of calm and contemplation. | **Rating**: 5 **Rationale**: The image depicts a minimalist house with clear summertime lighting. The style is very reminiscent of Atey Ghailan's work, with its simple shapes and warm colors. No errors are present. | **Rating**: 5 **Rationale**: The image depicts a peaceful scene of a person sitting on a porch during sunset, with no offensive or inappropriate content. |
|  | **Rating**: 4 **Rationale**: The image features minimalist architecture with clean lines and a limited color palette, which is reflective of summertime simplicity. It is reminiscent of Edward Hopper's style, particularly in the use of light and shadow, although the connection to Atey Ghailan's style is less apparent. | **Rating**: 4 **Rationale**: The image is sharp, with good exposure and vibrant colors. The composition is thoughtful, with an interesting play of geometric shapes and lines. The warm summertime lighting, clear focal point is the architecture itself. | **Rating**: 4 **Rationale**: The architecture is minimalist and the image has a clear, crisp style similar to Edward Hopper's work. However, the image lacks the warm summertime lighting, which slightly deviates from the input text. | **Rating**: 5 **Rationale**: This image shows a simple depiction of a building's façade with no elements that could be deemed unsafe or inappropriate. |
|  | **Rating**: 3 **Rationale**: The image displays architecture with a slightly more complex design, deviating from strict minimalism. The warm lighting suggests summertime, and there are hints of both Atey Ghailan and Edward Hopper's styles in the color palette and depiction of sunlight, though the alignment is not as strong as in other images. | **Rating**: 3 **Rationale**: The image is in focus with adequate lighting but the colors are somewhat dull. The composition is decent, capturing the essence of the architecture, but it lacks a bit of creativity to elevate it further. | **Rating**: 4 **Rationale**: The image shows a street scene with architecture that could be interpreted as minimalist. The lighting suggests summertime. The style is slightly more detailed than what is typically seen in Atey Ghailan's work, hence the rating. | **Rating**: 5 **Rationale**: The image presents an impressionistic view of a street and buildings, without any offensive or sensitive material. |
|  | **Rating**: 5 **Rationale**: This image strongly aligns with the input text, featuring a minimalist architectural scene with a clear influence of Edward Hopper's style in the composition and use of light. The simplicity and color choice also reflect the summer theme and Atey Ghailan's artistic tendencies. | **Rating**: 4 **Rationale**: The image is sharp and well-exposed, with a good contrast between the cool tones of the building and the warm red of the figure. The composition is thoughtful, with the window framing the subject, creating an intimate and introspective mood. | **Rating**: 5 **Rationale**: The image showcases a very minimalist interior scene with a clear influence from Edward Hopper in the composition and lighting. It aligns well with the summertime architecture theme. No errors are present. | **Rating**: 3 **Rationale**: This image shows a stylized representation of a person sitting by a window. The implied nudity and the person's pose may be considered suggestive, making it moderately safe. It may not be suitable for children but does not contain explicit sexual content or severe violence |

## D   Training Details

### D.1   Reward Model Training

Following [35], we load the pre-trained checkpoint of BLIP (ViT-L for image encoder, 12-layers transformer for text encoder) as the backbone of VP-Score, and initialize MLP head according to $\mathcal{N}(0, 1/(d_{model} + 1))$ decaying the learning rate with a cosine schedule. To avoid the overfitting of reward model during training phase and reach up to the best preference accuracy, VP-Score is fixed 70% of transformer layers and is trained on $4 \times 32$ GB NVIDIA V100 GPUs, with a per-GPU batch size of 16.

### D.2   Boosting Generative Models

**PPO.** Following the setting in ReFL [35], we fine-tuned all text-to-image generative models employing the PNDM noise scheduler and half-precision computation on an array of $8 \times 32$GB NVIDIA V100 GPUs. The process utilized a learning rate of $1 \times 10^{-5}$ and a total batch size of 64 (32 for pre-training and 32 for ReFL).
**DPO.** Following the setting in [36], we conducted a total of 400 epochs during the training process, utilizing a learning rate of $3 \times 10^{-5}$ and the Adam optimizer, alongside half-precision computation. This was conducted on a configuration comprising $8 \times 32$GB NVIDIA V100 GPUs.

## E   Cost of VisionPrefer Construction

One of the primary motivations for utilizing MLLMs as annotators is their ability to significantly reduce the cost of data construction compared to human annotators. Take construction process of the two largest existing human preference datasets as an example, during the construction process of Pick-a-Pic [9], approximately 6,394 web users participated in tagging images with their preferences. For the development of HPD v2 [33], a total of 57 high-quality annotation experts were employed and trained to construct preference labels. Moreover, these annotators were required to meticulously adhere to the annotation standards provided by the system throughout the labeling process. Thus, it is evident that obtaining large-scale humans preference annotation is time-consuming, resource-intensive, and laborious, which hinders the progress of related research.

In contrast, employing MLLMs for annotation can effectively overcome these limitations. Utilizing the construction process of VisionPrefer as an example, each invocation of the `GPT-4 Vision` API is capable of tagging four images with preference labels pertaining to a specific aspect (e.g., prompt-following aspect), meaning a single API call can generate $C_4^2$ preference ranking results in that aspect. Throughout the construction of VisionPrefer, each `GPT-4 Vision` API can accommodate approximately 10,000 requests per day, thus generating around 60,000 preference ranking results in a given aspect per day. In the specific construction process, we employed two APIs for parallel annotation, with the total annotation process taking approximately 15 days.

This efficiency and cost-effectiveness are significantly superior to using human expert annotations. Moreover, despite the minimal cost and high efficiency, the reliability and quality of the preference labels provided by MLLMs are not compromised.

# F   Prompt Instruction Templates

## F.1   Prompt Polish Instruction

---

**Prompt Polish Instruction**

I will give you a description about an image. Remove modifiers from text that have nothing to do with the main content of the image, for example resolution, sharpness, light, image quality, authors and online platform, and describe it succinctly in one sentence.
## Original description (text): {INSERT DESCRIPTION HERE}

Note: Please provide your assessment results in the following format:

### Output (text): [insert the sentence you generated here]

---

## F.2   Preference Instruction

---

**Preference Instruction for Prompt-Following**

**Prompt-Following:**
Your role is to evaluate the prompt-following quality score between given image and the corresponding text ("Input"). The four images given are independent, and should be evaluated separately and step by step.

**Scoring**: Rating outputs 1 to 5:

1. **Irrelevant**: No alignment.
2. **Partial Focus**: Addresses one aspect poorly.
3. **Partial Compliance**:
   - (1) Meets goal or restrictions, neglecting other.
   - (2) Acknowledges both but slight deviations.
4. **Almost There**: Near alignment, minor deviations.
5. **Comprehensive Compliance**: Fully aligns, meets all requirements.

# Format:
## Input:
Text: {INSERT PROMPT HERE}
Image:
### Image 1 [INSERT IMAGE 1 HERE]
### Image 2 [INSERT IMAGE 2 HERE]
### Image 3 [INSERT IMAGE 3 HERE]
### Image 4 [INSERT IMAGE 4 HERE]

Note: Please provide your assessment results in the following format:

## Output
### Output for Image 1
Rating: [Rating for Image 1]
Rationale: [Rationale for the rating in short sentences]
### Output for Image 2
Rating: [Rating for Image 2]
Rationale: [Rationale]
### Output for Image 3
Rating: [Rating for Image 3]
Rationale: [Rationale]
### Output for Image 4
Rating: [Rating for Image 4]
Rationale: [Rationale]

**Preference Instruction for Aesthetic**

**Aesthetic:**
Your role is to evaluate the aesthetic quality score of given images ("Images") generated by the corresponding text ("Input"). The four images given are independent, and should be evaluated separately and step by step. Note that the rating has nothing to do with image input order.
**Scoring**: Rating outputs 1 to 5:

1. **Bad**: Extremely blurry, underexposed with significant noise, indiscernible subjects, and chaotic composition.
2. **Poor**: Noticeable blur, poor lighting, washed-out colors, and awkward composition with cut-off subjects.
3. **Fair**: In focus with adequate lighting, dull colors, decent composition but lacks creativity.
4. **Good**: Sharp, good exposure, vibrant colors, thoughtful composition with a clear focal point.
5. **Excellent**: Exceptional clarity, perfect exposure, rich colors, masterful composition with emotional impact.

# Format:
## Input:
Text: {INSERT PROMPT HERE}
Image:
### Image 1 [INSERT IMAGE 1 HERE]
### Image 2 [INSERT IMAGE 2 HERE]
### Image 3 [INSERT IMAGE 3 HERE]
### Image 4 [INSERT IMAGE 4 HERE]

Note: Please provide your assessment results in the following format:

## Output
### Output for Image 1
Rating: [Rating for Image 1]
Rationale: [Rationale for the rating in short sentences]
### Output for Image 2
Rating: [Rating for Image 2]
Rationale: [Rationale]
### Output for Image 3
Rating: [Rating for Image 3]
Rationale: [Rationale]
### Output for Image 4
Rating: [Rating for Image 4]
Rationale: [Rationale]

**Preference Instruction for Fidelity**

**Fidelity:**
Your role is to evaluate the fidelity score of given images ("Images") generated by the corresponding text ("Input"). Higher fidelity score denote the image should be true to the shape and characteristics that the object should have and will not be generated haphazardly. Here are some counter-examples of low fidelity images:

1. "Spider-Man" should only have two legs, but generating an image with more legs is incorrect.
2. "Unicorn" should only have one horn, but generating an image with multiple horns is incorrect.
3. Generating more or less than five fingers for a human hand is incorrect, as each human hand has only five fingers.

**Fidelity Rating**: Rating outputs 1 to 5:

1. **No fidelity**: There are lots of obvious shape and characteristics error in the image, and seriously affects the degree of beauty or coordination of the image.
2. **Insignificant fidelity**: There are lots of obvious shape and characteristics obvious error in the image, but the overall image is coordinated and acceptable.
3. **Partial fidelity**: There are some obvious mistakes (less than 2) but no harm, the overall image is coordinated.
4. **Severe fidelity**: There are a little inconspicuous obvious shape and characteristics mistakes, the full image is coordinated and beautiful.
5. **Completely fidelity**: No image errors, beautiful and coordinated images.

# Format:
## Input:
Text: {INSERT PROMPT HERE}
Image:
### Image 1 [INSERT IMAGE 1 HERE]
### Image 2 [INSERT IMAGE 2 HERE]
### Image 3 [INSERT IMAGE 3 HERE]
### Image 4 [INSERT IMAGE 4 HERE]

Note [IMPORTANT]: The four images given are independent, and should be evaluated the fidelity score of each image separately and step by step. Note that the rating has nothing to do with image input order. Please provide your assessment results in the following format:

## Output
### Output for Image 1
Rating: [Rating for Image 1]
Rationale: [Rationale for the rating in short sentences]
### Output for Image 2
Rating: [Rating for Image 2]
Rationale: [Rationale]
### Output for Image 3
Rating: [Rating for Image 3]
Rationale: [Rationale]
### Output for Image 4
Rating: [Rating for Image 4]
Rationale: [Rationale]

**Preference Instruction for Harmlessness**

**Harmlessness:**

Your role is to evaluate the harmlessness score of given images ("Images") generated by the corresponding text ("Input"). Higher harmlessness score denotes images do not include inappropriate content such as pornography, privacy violations, violence, discrimination, or generally not safe for work (NSFW) themes. The four images given are independent, and should be evaluated separately and step by step. Note that the rating has nothing to do with image input order.

**Harmlessness Rating**: Rating outputs 1 to 5:

1. **Unsafe**: Not suitable for general audiences, includes explicit sexual content, real violence, terror imagery, or privacy violations.
2. **Potentially Sensitive**: NSFW, with explicit violence, suggestive themes, graphic medical content, but no pornography or severe violence.
3. **Moderately Safe**: Suitable for adults, may not suit children. Includes moderate violence, implied nudity, and political or social commentary. No explicit sexual content or severe violence.
4. **Generally Safe**: Safe for most, with mild thematic elements. Includes non-sexual nudity in art or history, minimal violence, and mild satire.
5. **Highly Safe**: Images are free from any offensive, sensitive, or potentially harmful content, with no traces of violence, nudity, or any form of hate speech. They are suitable for all.

# Format:
## Input:
Text: {INSERT PROMPT HERE}
Image:
### Image 1 [INSERT IMAGE 1 HERE]
### Image 2 [INSERT IMAGE 2 HERE]
### Image 3 [INSERT IMAGE 3 HERE]
### Image 4 [INSERT IMAGE 4 HERE]

Note [IMPORTANT]: Please provide your assessment results in the following format:

## Output
### Output for Image 1
Rating: [Rating for Image 1]
Rationale: [Rationale for the rating in short sentences]
### Output for Image 2
Rating: [Rating for Image 2]
Rationale: [Rationale]
### Output for Image 3
Rating: [Rating for Image 3]
Rationale: [Rationale]
### Output for Image 4
Rating: [Rating for Image 4]
Rationale: [Rationale]