# User-Specific Preference Prediction on Generated Images

Anonymous Submission

## Abstract

*User preference prediction requires a comprehensive and accurate understanding of individual tastes. This includes both surface-level attributes, such as color and style, and deeper content-related aspects, such as themes and composition. Existing methods focus on superficial features instead of deep semantic content, or rely on general human preferences while ignoring individual user variations. To address these limitations, we propose an approach built upon Vision-Language Models, introducing contrastive preference loss and preference tokens to learn personalized user preferences from historical interactions. The contrastive preference loss is designed to effectively distinguish between user "likes" and "dislikes", while the learnable preference tokens capture shared interest representations among existing users, enabling the model to generalize and transfer learned knowledge to new users with similar preferences. Extensive experiments demonstrate our model outperforms other methods in preference prediction accuracy, effectively identifying users with similar aesthetic inclinations and providing more precise guidance for generating images that align with individual tastes.*

## 1. Introduction

Recent work in generative models [2, 6, 7, 11, 25, 28, 29, 31, 33, 37] has significantly advanced the field of text-to-image generation. However, these models often produce generic outputs that may not align with the diverse and nuanced preferences of each individual user. A particularly promising direction within this domain is user preference prediction based on generated images, which has garnered increasing attention due to its capability to provide guidance to generative models tailored to individual preferences. By aligning generated content with specific user interests, this direction holds the potential to deliver unique user experiences, thereby enhancing user satisfaction and engagement.

In user preference prediction, the task is to identify preferences, such as color and contents that align with a user's tastes using reference data, typically a set of liked and disliked images. Fig. 1 provides an illustrative example.

Existing preference prediction models such as



Figure 1. Our task aims to predict target images that align with users' preferences based on their history preferences. For example, based on a set of liked and disliked images, user $A$ and user $B$ demonstrate unique preferences but also share common interests.

PickScore [13], ImageReward [41], and HPS [39, 40] are designed to evaluate human preferences at a general level, they lack the capability to effectively capture individual-level preferences. Moreover, recent individual-level personalized preference modeling [32, 36] presents three primary issues: (1) focus on superficial attributes like color and style, which limits their ability to capture the essence of a deep content-level user's preference and (2) overlook the significance of users' disliked images, which provide valuable negative feedback and relative preference signals for refining preference understanding, (3) fail to utilize the fact that users with similar tastes might share preferences for certain types of images.

To address these challenges, based on Vision-Language Models (VLMs) [12, 14, 15, 17, 18, 22, 23, 43, 44], we propose an approach that introduces contrastive preference loss and preference tokens. Our model extracts content-driven patterns from user history preferences, employs contrastive preference loss to differentiate "like" or "dislike" between contents. Additionally, learnable preference tokens that represent shared interests among users, allow us to incorporate similar preferences as reference points, thereby aiding in the

identification and modeling of new users' preferences. Our key contributions are summarized as follows:

- We introduce a VLM-based contrastive learning framework that enables the model to learn discriminative features from users' liked and disliked data, effectively capturing fine-grained user preferences by modeling relative preference relationships among samples.
- We leverage learnable preference tokens to capture shared interests among users, allowing the model to generalize better across users with similar tastes.
- Experimental results demonstrate that our model outperforms existing methods in preference recognition accuracy. It is able to identify users with similar tastes and effectively generalizes to new users with similar preferences. Furthermore, it provides more precise guidance for generating personalized content.

## 2. Related Work

Modeling user preferences in text-to-image generation is essential for improving alignment with human aesthetics and expectations. Existing research in this area can be broadly categorized into two main categories: general preference modeling, which focuses on capturing collective human judgments to enhance overall image quality, and user-specific preference modeling, which personalizes image generation based on individual tastes and behaviors.

**General Preference Modeling for Human-Aligned Image Generation.** Researchers have explored various strategies to improve alignment, categorized into three approaches: (1) Filtering Training Data with Preference Scores. By selecting training data based on human feedback scores or automated metrics, models can benefit from high-quality examples that reflect specific user demands. For instance, Liang *et al.* [20] demonstrates how filtering data based on feedback scores leads to improved model performance, as it ensures that only the most relevant examples are used for fine-tuning. Similarly, HPS [39, 40] builds upon this concept by introducing a scoring mechanism to prioritize image-text pairs closely aligned with user preferences, making the model more responsive to varied user expectations. (2) Reward-Weighted Fine-Tuning for Human-Aligned Models. In this approach, models are fine-tuned using reward signals that weigh heavily on user satisfaction. Lee *et al.* [16] exemplifies this by incorporating feedback-based rewards during training, which generates outputs aligned with user preferences. Furthermore, ImageReward [41] provides a structured method for translating human judgments into reward functions, which guides the model's fine-tuning process. By giving greater importance to rewards that capture user satisfaction, these methods tailor the model's outputs to reflect diverse and nuanced user tastes. (3) Reinforcement Learning for Preference Optimization [4, 10, 21, 24]. Recent work [10, 24] uses rein-

forcement learning to optimize the input prompts for high-quality images. DiffusionDPO [38] leverages user preferences to iteratively fine-tune the model, improving its ability to generate images that reflect user choices. Similarly, D3PO [42] introduces a dynamic update mechanism driven by continuous user feedback, allowing the model to maintain robustness, even as user interests evolve.

**User-Specific Preference Modeling and Personalized Image Generation.** In recent advancements in personalized image generation, several approaches have emerged to better align generative models with individual needs. DreamBooth [30] and Textual Inversion [9] explore personalization by fine-tuning pre-trained models with just a few example images, allowing users to introduce unique characters or styles. This approach, while effective for small datasets, focuses on integrating specific instances rather than broader user behaviors. To improve personalization, Salehi *et al.* [32] proposes a standardized process to collect user preferences using a few query images. User feedback is then systematically incorporated to adjust the preferences extracted from the user during the generation process. Additionally, Shen *et al.* [36] introduces a method to integrate user-specific preferences across different modalities, such as text and images, creating personalized outputs by leveraging historical interactions, such as clicks and conversations. This multimodal approach significantly enhances the models' adaptability to align with user needs. In our work, we introduce a VLM-based contrastive learning framework to model relative preference rankings, ensuring that the model not only distinguishes between liked and disliked samples but also learns their comparative importance. Additionally, we employ preference tokens to capture users' shared preferences, leveraging the preferences of similar users as a reference to enhance the identification of individual user preferences more effectively.

## 3. Method

Our objective is to determine whether a given target item $z = (I, T)$ aligns with a user's preferences based on their historical selections. We define a preference history sequence as $\mathcal{S} = \{s_i | s_i = (I_{\text{pos}}, I_{\text{neg}}, T)_i\}_{i=1}^{N_{\text{ref}}}$, where each entry consists of a text prompt $T$ describing the image content and two generated images, $I_{\text{pos}}$ and $I_{\text{neg}}$, which are produced based on the prompt. The user has labeled $I_{\text{pos}}$ as liked and $I_{\text{neg}}$ as disliked. The total number of reference entries in the user's preference history is denoted by $N_{\text{ref}}$. By leveraging this structured preference data, we aim to model user-specific preferences and improve the alignment of generated images with individual tastes.

### 3.1. Overview

As shown in Fig. 2, we propose a VLM-based contrastive preference learning framework that enables the model to
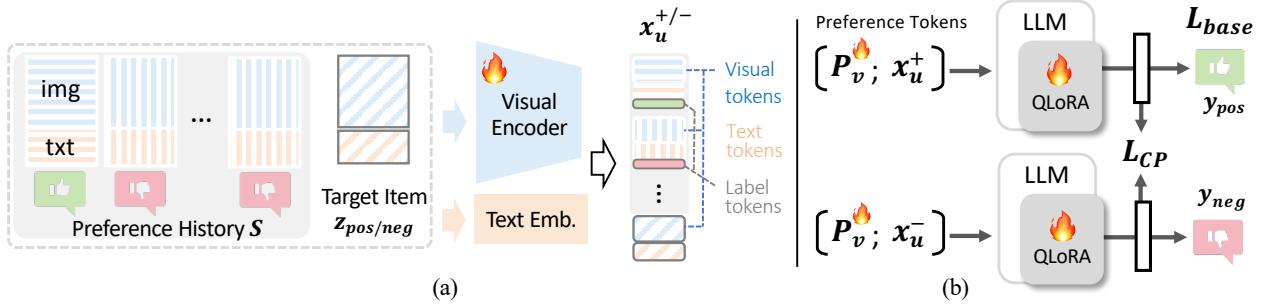
Figure 2. Overview of our VLM-based contrastive learning framework. (a) The model extracts user-specific preference representations by processing a user's preference history $\mathcal{S}$ and a target item $z_{\text{pos/neg}}$. The visual encoder processes image features, while the text embedding module encodes textual descriptions and labels. These processed inputs generate user-specific representations: $x_u^+$ for target images that the user likes and $x_u^-$ for those they dislike. (b) The framework is trained using a base loss $L_{\text{base}}$ that guides the model in predicting user preferences, and a contrastive preference loss $L_{\text{CP}}$ that refines the model by enforcing relative ranking between liked and disliked items. Additionally, learnable preference tokens $P_v$ are introduced to model shared user interests. These tokens are concatenated with user-specific representations, enabling the model to leverage shared user interests as reference points.

distinguish user preferences by learning from liked and disliked images, effectively capturing relative preference relationships and improving prediction accuracy. To further enhance personalization, we introduce preference tokens that capture shared interests among users, allowing the model to better generalize and adapt to diverse user preferences.

### 3.2. Contrastive Preference Learning

We denote our model as $\mathcal{M}$, which conditions on a user's preference history $\mathcal{S}$ to assess the likelihood of a user favoring a particular item $z$. For the target item $z$, we define user preference as $z_{\text{pos}}$ if the user likes the item and $z_{\text{neg}}$ if the user dislikes it. We use $z_{\text{pos/neg}}$ to denote either case generically. We define a comprehensive loss function that combines a base classification loss with a contrastive preference loss, aiming to improve the model's ability to distinguish between "like" and "dislike" predictions.

#### 3.2.1. Base Loss

The base loss, $L_{\text{base}}$, aims to minimize the classification error across both "like" and "dislike" samples. Let $\mathcal{M}^+(\mathcal{S}, z)$ and $\mathcal{M}^-(\mathcal{S}, z)$ represent the logit outputs for predicting "like" and "dislike" outcomes for a sample $z$, respectively. The associated ground-truth labels are represented as $\mathbf{y}_{\text{pos}}$ and $\mathbf{y}_{\text{neg}}$, respectively. The base loss is defined as:

$$L_{\text{base}} = \frac{1}{2}\left(\mathcal{L}(\mathcal{M}^+(\mathcal{S}, z_{\text{pos}}), \mathbf{y}_{\text{pos}}) + \mathcal{L}(\mathcal{M}^-(\mathcal{S}, z_{\text{neg}}), \mathbf{y}_{\text{neg}})\right),$$
(1)

where $\mathcal{L}(\cdot)$ denotes a classification loss function. Additionally, we use $\mathcal{Q}(\mathcal{S}, z)$ to predict how much the user prefers:

$$\mathcal{Q}(\mathcal{S}, z) = \frac{\exp(\mathcal{M}^+(\mathcal{S}, z))}{\exp(\mathcal{M}^+(\mathcal{S}, z)) + \exp(\mathcal{M}^-(\mathcal{S}, z))}.$$
(2)

#### 3.2.2. Contrastive Preference Loss

To complement the base loss, we introduce two contrastive preference loss terms, $L_+$ and $L_-$, which enhance the

model's ability to differentiate between "like" and "dislike" predictions by emphasizing their relative rankings. While the base loss, $L_{\text{base}}$, effectively minimizes classification errors for individual "like" and "dislike" labels, it fails to capture relative preference when only pairwise comparisons are available. Specifically, given a pairwise relationship such as $A \succ B$ (where $\succ$ denotes a preference), $L_{\text{base}}$ struggles to distinguish the relative ranking of $A$ and $B$ as it treats each sample independently without explicitly modeling their comparative importance. Consequently, the model may assign similar preference scores to both $A$ and $B$ in $\mathcal{Q}(\mathcal{S}, z)$, since $L_{\text{base}}$ overlooks the structured ranking constraints that enforce a clear distinction between correctly classified "like" predictions for disliked samples and "dislike" predictions for liked samples. To overcome this limitation, the contrastive preference loss terms, $L_+$ and $L_-$, address this limitation by incorporating pairwise ranking information, thereby refining preference predictions and improving the model's ability to distinguish between closely related "like" and "dislike" cases.

**Positive Preference Loss ($L_+$).** This loss term focuses on ensuring that the model assigns a higher score to positive samples compared to negative ones, encouraging the model to prioritize positive outcomes:

$$L_+ = -\frac{1}{N}\sum_{i=1}^{N}\log\sigma(\mathcal{M}^+(\mathcal{S}, z_{\text{pos}}) - \mathcal{M}^+(\mathcal{S}, z_{\text{neg}})),\quad (3)$$

where $N$ is the number of samples and $\sigma$ is the sigmoid function.

**Negative Preference Loss ($L_-$).** This loss term ensures that the model assigns a higher score to negative samples when predicting a negative outcome, encouraging the model
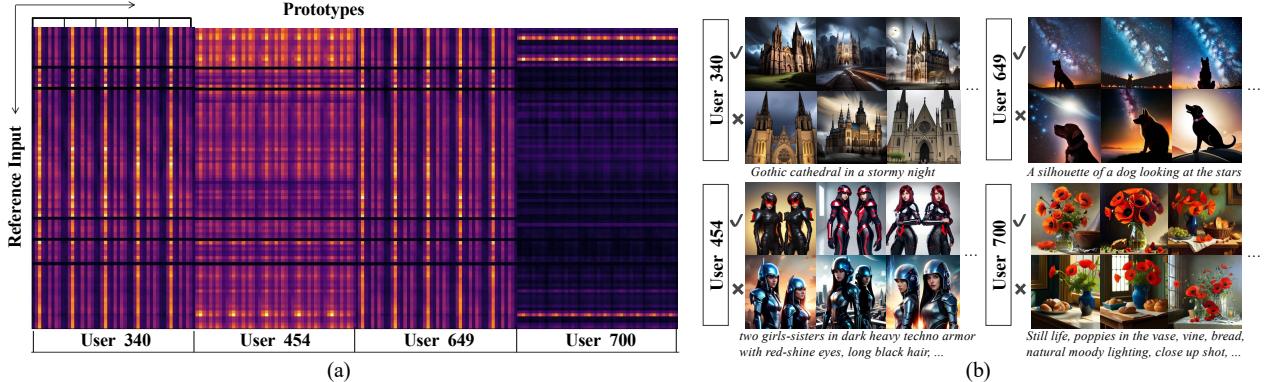
Figure 3. (a) Attention scores $\mathcal{A}$ represent interactions between preference tokens and target image tokens for individual users. Each user has a unique reference history, and we concatenate the same target image to the input sequence across users. For each user, the vertical axis represents tokens from the target image, while the horizontal axis represents the preference tokens. Each user has five different random re-orderings of reference images. (b) Examples of images liked (✓) or disliked (✗) by each user.

to prioritize negative samples appropriately:

$$L_- = -\frac{1}{N} \sum_{i=1}^{N} \log \sigma(\mathcal{M}^-(\mathcal{S}, z_{\text{neg}}) - \mathcal{M}^-(\mathcal{S}, z_{\text{pos}})). \quad (4)$$

Then, we calculate $L_{\text{CP}} = L_+ + L_-$ to obtain the contrastive preference loss.

The combined loss function, which incorporates both the base and contrastive preference loss, enhances the model's ability to distinguish user preferences by refining predictions for both positive and negative outcomes:

$$L_{\text{all}} = L_{\text{base}} + L_{\text{CP}}, \quad (5)$$

facilitating the model to optimize nuanced preference distinctions for more accurate and effective predictions.

### 3.3. Learnable Preference Tokens

To enhance user preference modeling, we introduce the learnable preference tokens $P_v \in \mathbb{R}^{L_p \times D}$ as part of the input sequence to the VLM, where $L_p$ represents the number of preference tokens and $D$ is the embedding dimension. All reference entries in the history sequences $\mathcal{S}$ and target item $z$, except for the target image label token, are encoded and stacked as the user-specific input token sequence, denoted as $x_u$. The preference tokens are then concatenated with this user-specific sequence $x_u$ to form the final input $[P_v; x_u]$, where $x_u \in \mathbb{R}^{L_e \times D}$ represents the embedded input tokens, $L_e$ is the length of input tokens and $[\,;\,]$ denotes the concatenation operation.

**Mining Similar Users' Preferences via the Attention Mechanism.** Our model leverages preference tokens to capture shared interests among users, allowing it to generalize preference prediction beyond individual interaction histories. A key aspect of this mechanism is the attention-based interaction between the user input tokens and the preference tokens within the transformer layers. The attention scores between input tokens and preference tokens are computed as:

$$\mathcal{A} = \text{softmax}\left(\frac{W_q(x_u)W_k(P_v)^T}{\sqrt{D'}}\right), \quad (6)$$

where $W_q$ and $W_k$ are linear projections in the attention mechanism, $D'$ is the embedding dimension, and $\mathcal{A} \in \mathbb{R}^{L_e \times L_p}$ represents the attention scores between each input token and the preference tokens.

To better understand how preference tokens facilitate user similarity modeling and generalization to unseen users, we analyze the learned attention scores $\mathcal{A}$, which capture the interactions between input tokens and preference tokens. Fig. 3 visualizes these interactions, where the same target image is concatenated across users with different reference histories to examine how their preferences are represented. Specifically, Fig. 3 (a) shows Users 340 and 649, who exhibit a highly similar pattern of attention across multiple preference tokens, suggesting that they share a common aesthetic inclination. Notably, User 649 is present in the training set, while User 340 is an unseen user. However, the learned preference tokens effectively bridge this gap by encoding shared thematic patterns, such as an affinity for landscapes with dramatic skies, silhouettes, and nightscapes. This observation supports our claim that preference tokens serve as a structured preference representation that captures common aesthetic traits across users, transfers knowledge to unseen users, ensuring that their preferences are accurately inferred without requiring direct memorization of past interactions. In contrast, Fig. 3 (b) illustrates that Users 454 and 700 exhibit distinct attention patterns, reinforcing that the preference token space does not simply cluster all users together arbitrarily but rather preserves individual differences while leveraging commonalities where applicable. Further details and in-depth anal-

Figure 4. Qualitative comparison of user preference alignment across models. We compare our model to ViPer [32], PickScore [13], ImageReward [41], CLIP [27], and Aesthetic Score [35]. Subfigures (a) and (b) illustrate user-specific preferences for style and content, respectively. Each subfigure contains images categorized as "like" and "Dislike" based on user reference preferences. The green boxes represent the desired outputs that align with the user's preference. Our model consistently demonstrates a higher accuracy in predicting the user's preferences than other models, showcasing its superior alignment with user-specific preferences.

ysis of this experiment can be found in the appendix.

## 4. Experiments

### 4.1. User-Specific Preference Prediction

**Datasets.** We process Pick-a-Pic v2 dataset [13] to obtain user-specific preference datasets based on user IDs. Since Pick-a-Pic is collected through real user interactions, our model learns from actual user choices rather than artificial or automated metrics. This large-scale, diverse dataset captures a broad spectrum of aesthetic preferences, making it a strong benchmark for user-specific preference modeling. The processed dataset includes $224,952$ images and $2,267$ users in the training set, $1,707$ images and $89$ users in the validation set, and $2,234$ images and $70$ users in the test set. To better evaluate our proposed framework, we divide test data into two parts: a 'seen' dataset and an 'unseen' dataset. 'Seen' refers to users who appear in the training set but have different images in the test set, while 'unseen' refers to users who do not appear in the training set at all. The test set includes $459$ images from seen users and $1,775$ images from unseen users.

**Implementation Details.** Following the approach of [32], we use IDEFICS2-8B [15] as our VLM. To conserve memory, each prompt is truncated to a maximum length of $100$ tokens, and input images are resized to $512 \times 512$ pixels. We employ a batch size of $64$, training on 8 A100 (80GB) GPUs with a local batch size of $2$ pairs and gradient ac-

cumulation over $4$ steps. Following the setup of [32], we set the length of each user's preference history sequence, $N_{\text{ref}}$, to $8$. The learning rate is set to $1 \times 10^{-5}$, with a weight decay of $1 \times 10^{-2}$. The language model is fine-tuned using QLoRA [5], while the vision encoder is trained simultaneously. The input tokens template for the VLM is "<image>The prompt is <prompt>. Score for this image?<label>". Initially, the VLM is trained with our custom loss function for 5k steps, after which the model weights are fixed, and only the learnable preference tokens are further tuned for an additional 16k steps. To prevent the model from learning a fixed pattern, we randomly shuffle the order of reference history sequences when training.

**Evaluation Metric.** For evaluating user-specific preference prediction, we assess our method using top-$K$ accuracy, which determines whether the liked image is ranked among the top $K$ candidates. Among all candidates, only one "like" image is provided. When comparing one liked image against one disliked image, we use top-1 accuracy.

**Comparison to Other Methods.** In our study, we compare our method with several existing approaches to better understand its efficacy: (1) ViPer proxy model [32], which predicts user preferences by analyzing reference images without text descriptions, treats each preference independently rather than modeling their relative ranking. (2) PickScore [13], (3) ImageReward [41], (4) CLIP [27], and (5) LAION Aesthetic Score Predictor [35]. PickScore and ImageReward focus on learning general human preferences

| Model | Aes Score | CLIP Score | ImageReward | PickScore* | PickScore | IDEFICS | ViPer | ViPer* | Ours |
|---|---|---|---|---|---|---|---|---|---|
| $N_{\text{ref}}$ | 0 | 0 | 0 | 0 | 0 | 8 | 8 | 8 | 8 |
| Top-1 acc (%) | 49.96 | 53.13 | 55.64 | 57.72 | 61.82 | 50.27 | 55.15 | 57.39 | **61.68** |

\* Trained with the same settings as our model.

Table 1. Quantitative comparison between a liked and a disliked case.

| Model | $N_{\text{ref}}$ | Top-1 Acc | Top-2 Acc | Top-3 Acc |
|---|---|---|---|---|
| Random | 0 | 25.0 | 50.0 | 75.0 |
| Aes Score | 0 | 28.11 | 54.12 | 78.33 |
| CLIP Score | 0 | 30.04 | 55.82 | 76.05 |
| ImageReward | 0 | 31.42 | 58.01 | 78.47 |
| IDEFICS | 8 | 24.40 | 51.88 | 78.33 |
| ViPer | 8 | 31.20 | 56.45 | 78.65 |
| ViPer* | 8 | 33.62 | 59.49 | 80.84 |
| w/o $P_v$ | 8 | 35.72 | 61.64 | 83.44 |
| Ours | 8 | **37.47** | **62.85** | **84.74** |

Table 2. A quantitative comparison between the liked case and three disliked cases. We report the top-1 to top-3 accuracy (%).

| | IDEFICS | ViPer | ViPer* | w/o $P_v$ | Ours |
|---|---|---|---|---|---|
| Seen | 51.41 | 54.03 | 58.17 | 60.35 | **61.44** |
| Unseen | 50.31 | 55.38 | 57.18 | 61.63 | **61.75** |
| $|\triangle|$ (%) | 1.10 | 1.35 | 0.99 | 1.28 | **0.31** |

\* Trained with the same settings as our model.

Table 3. Top-1 accuracy on seen-unseen data with $N_{\text{ref}} = 8$.

and consider relative preferences between images. CLIP and Aesthetic Score are designed to evaluate generic text-image alignment and aesthetic quality respectively. To ensure a fair comparison, we carefully set the hyperparameters for all baseline models to align with their original implementation guidelines and relevant prior work. Specifically, for ViPer and PickScore, we follow their reported training configurations and conduct additional tuning to optimize their performance under our experimental setup, and these extended versions are marked by '*' in the results.

**Qualitative User-Specific Preference Prediction.** In Fig. 4, our model effectively aligns with user-specific preferences by distinguishing styles and content according to user reference data. For instance, in Fig. 4 (a), our method accurately captures the user's preference for anime-style imagery with specific attributes such as color, theme, and character features. In Fig. 4 (b), our method alleviates semantic ambiguity, particularly in cases such as the interpretation of "grey cat," ensuring that the generated images better reflect the user's intended preferences. More qualitative comparisons are in the appendix.

**Quantitative User-Specific Preference Prediction.** Tab. 1 and Tab. 2 compare different models in terms of top-$K$ accuracy. Our model outperforms baselines, such as ViPer, PickScore, CLIP score, and ImageReward. In Tab. 1, it achieves the highest top-1 accuracy in like-dislike pairs. Since PickScore (shown in gray) is trained on the full Pick-a-Pic dataset, it is less comparable to our method. For a fair comparison, we focus on asterisk-marked models, which use the same training settings as ours. Tab. 2 further highlights our model's advantage, particularly in scenarios involving multiple disliked cases. These results indicate that

our approach better aligns with user preferences. Also, generic metrics, including Aesthetic score and CLIP score, report the worst accuracy, indicating that user-specific preferences may differ significantly from general preferences. PickScore, designed for general human preference modeling rather than personalized user-specific preferences, lacks the ability to capture fine-grained individual variations, resulting in lower accuracy in user-specific preference prediction tasks. The ViPer proxy model, which only relies on reference images without text descriptions, treats each preference independently and struggles with comprehending image content and modeling relative rankings. In contrast, our approach leverages contrastive preference learning and preference tokens to effectively capture both individual relative rankings and shared preference characteristics, leading to more accurate and robust preference prediction.

**Preference Tokens Enable Better Generalization to New Users.** To accurately model user preferences, a robust system should maintain consistent accuracy across both seen and unseen users. Our framework leverages learnable preference tokens to effectively transfer learned user preference structures, significantly improving generalization. As shown in Tab. 3, our method reduces the performance gap between seen and unseen users to 0.31%, outperforming ViPer and other baselines. Furthermore, we compare the difference in top-1 accuracy across different $N_{ref}$ settings, where a smaller difference indicates more accurate predictions for unseen users and improved model generalization. According to Fig. 5, our approach maintains stable accuracy even when fewer historical preferences are available. The key factor behind this improved generalization is that preference tokens effectively capture shared user interests and allow the model to leverage information from similar users as Fig. 3 shows. This aligns with findings in context optimization methods [19, 45], where learnable tokens improve generalization in both vision and language tasks by capturing shared patterns across seen and unseen categories.
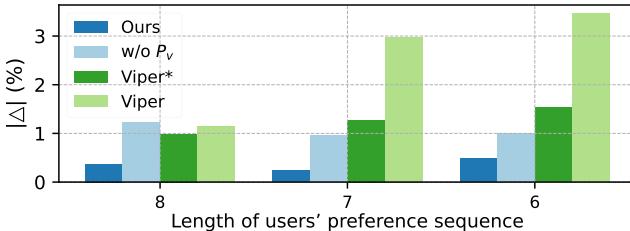
Figure 5. Difference between Top-1 accuracy on seen-unseen data with different $N_{\text{ref}}$.

| Model | Anim. | C-Art | Paint. | Photo | Avg. |
|---|---|---|---|---|---|
| Baseline | 27.75 | 26.86 | 27.06 | 27.36 | 27.26 |
| Openjourney | 27.85 | 27.18 | 27.25 | 27.53 | 27.45 |
| ChilloutMix | 27.92 | 27.29 | 27.32 | 27.61 | 27.54 |
| Ours | **28.42** | **27.47** | **27.79** | **28.14** | **27.95** |

Table 4. Performance on HPSv2 benchmark.

**Number of User Reference Preferences.** As demonstrated in Fig. 6, our method consistently achieves the highest top-1 accuracy even as the length of preference sequences decreases. This indicates that our model effectively preserves accuracy with less reference data, demonstrating its robustness. In contrast, other methods show a noticeable decline in accuracy as the sequence length shortens, highlighting the stability and adaptability of our approach in scenarios with limited user reference information.
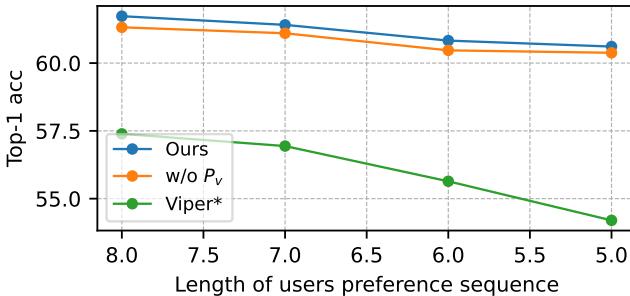


Figure 6. Top-1 accuracy on test dataset with different $N_{\text{ref}}$.

## 4.2. Enhancing Generation with User Preferences

**Datasets.** We evaluate our model's ability to learn more detailed attribute preferences and guide image generation based on user reference data using the HPSv2 benchmark [39]. The HPSv2 benchmark is a large-scale evaluation framework designed to assess human-aligned preferences in text-to-image generation. It includes 3,200 prompts covering four distinct styles: Animation, Concept Art, Painting, and Photo.

**Experimental Setup.** Following the approach of [8], we generate images guided by our model, incorporating both user likes and dislikes as feedback. To optimize GPU memory usage, we limit the number of generated reference images per user to six. Based on the prompts from the HPSv2 benchmark, we generate six images: three samples from FLUX.1-schnell [3] and three from Stable Diffusion 1.5 Turbo [34], using their default hyperparameter settings. We designate the images generated by FLUX.1-schnell as "liked" samples and those from Stable Diffusion 1.5 Turbo as "disliked" samples, forming the user's preference data. Additionally, we generate an image using Stable Diffusion

1.5 Turbo as the target item and apply Eq. (2) to obtain the guidance signal, which is then iteratively refined using this guidance to optimize the target image. Since FLUX has better generative capabilities than SD-Turbo, especially in aesthetic quality, text-image alignment, and attribute consistency, this setup allows us to assess whether our model can accurately capture these nuanced user preferences and use them to improve image generation. More experimental details are provided in the appendix.

**Evaluation Metric.** We utilize the HPSv2 score on the HPSv2 benchmark. This metric is trained on $798,090$ human-annotated comparisons and measures image quality in relation to human preferences.

**Enhancing Image Generation with User Preferences.** As shown in Tab. 4, our method effectively enhances Stable Diffusion 1.5 Turbo (Baseline) by learning richer attribute preferences from user feedback. Compared to the baseline, our model shows a notable increase across all four categories, with the most significant gains in Concept Art ($+0.61$) and Animation ($+0.67$). These improvements highlight the effectiveness of incorporating user feedback in enhancing artistic coherence and stylistic alignment. Furthermore, with only six history images and no additional model training, our approach enhances image generation quality, surpassing ChilloutMix [1] and Openjourney [26], two well-established models fine-tuned for realistic image generation. This suggests that our model can effectively utilize user preferences to refine and guide image generation. It also indicates that user preference signals can be as impactful as model-level enhancements, providing a flexible and efficient way to improve generated outputs. Fig. 7 shows several generated examples from the HPSv2 benchmark. Our method provides guidance for aligning generated images with user preferences without requiring additional model training, making it an efficient and adaptable approach for preference-based generation.

**Comparison of Like and Dislike Predictions as Guidance for Enhancing Generation.** Figure 8 illustrates how different guidance settings influence the generated images. The "w/ Like" setting increases the likelihood of producing images with a blue-themed style that aligns with user preferences, while the "w/ Dislike" setting emphasizes an unwanted pink hue. These results show that our model effectively captures both positive and negative user pref-

*A helmet-wearing monkey skating.*  *An elderly lady pours some cups of tea on a tray.*  *A bathroom has pink tiles and a black toilet.*  *This is a red bike on a dirt path.*  *Portrait of Archduke Franz Ferdinand by Charlotte Grimm, depicting his detailed face.*  *A painting by Greg Manchess depicting an anime woman.*  *A close-up anime portrait of Sailor Moon against a grey background with Russian panel housing in bokeh.*  *A portrait of a beautiful anime girl with pink hair wearing a white t-shirt and looking directly at the viewer.*

Figure 7. Some examples of images generated on HPSv2 benchmark.



*light blue haired anime girl ocean themed anime isopod antenna twintails*

Figure 8. Using the approach from [8], we generate images guided by our model, incorporating both user likes and dislikes as feedback. Adjusting the guidance based on "Like" or "Dislike" produces distinct variations, demonstrating the impact of user feedback on image generation results. The images in each row are generated using the same random seed.

erences, allowing for a more comprehensive and nuanced modeling of user preferences.

### 4.3. Ablation Study

We conduct ablation experiments on Pick-a-Pic v2 dataset to examine the effects of contrastive preference loss, learnable preference tokens and the length of preference tokens.

As shown in Tab. 5, incorporating the contrastive preference loss term, $L_{CP}$, further improves the accuracy by 0.9%, suggesting that the contrastive preference loss term aids in refining the model's alignment with the target outputs. Lastly, the full model, which incorporates learnable preference tokens in addition to the previous components, achieves an accuracy of 61.68%, representing a cumulative performance gain over the baseline model. In Tab. 3, we observe that removing $P_v$ significantly increases the performance gap between seen and unseen users. The difference between seen and unseen user accuracy expands from

|  | Baseline (w/ Pmpt) | w/ $L_{CP}$ | Full (w/ Tokens) |
|---|---|---|---|
| Top-1 Acc (%) | 60.47 | 61.37 (+0.9) | 61.68 (+0.31) |

Table 5. Ablation Study for different Settings.

| Number of Preference Tokens | 5 | 10 | 20 |
|---|---|---|---|
| Top-1 Acc (%) | 61.41 | 61.68 | 61.19 |

Table 6. Ablation study for preference tokens numbers.

0.31% (with $P_v$) to 1.28% (without $P_v$), a 0.97% increase in the generalization gap. This suggests that preference tokens help the model generalize better to new users by capturing shared preference structures.

In Tab. 6, the results indicate that using 10 preference tokens yields the highest top-1 accuracy, slightly outperforming configurations with 5 and 20 preference tokens. This finding suggests that an optimal number of preference tokens is crucial for effectively modeling user preferences without overfitting or under-representing data variation.

### 5. Conclusion

In this paper, we propose a novel approach for user-specific preference prediction in generated images by leveraging Vision-Language Models (VLMs). To address the limitations of existing methods that focus primarily on superficial attributes or general human preferences, we introduce contrastive preference loss and learnable preference tokens. The contrastive preference loss enables the model to distinguish between users' "likes" and "dislikes" more effectively, while the preference tokens capture shared interests across users, facilitating better generalization to new users with similar preferences. Extensive experiments demonstrate that our model outperforms existing methods in preference prediction accuracy, effectively identifying users with similar aesthetic inclinations and providing more precise guidance for personalized content generation.

# References

[1] ChilloutMix: A high-quality anime image generation model, 2023. 7

[2] Fan Bao, Shen Nie, Kaiwen Xue, Yue Cao, Chongxuan Li, Hang Su, and Jun Zhu. All are worth words: A vit backbone for diffusion models. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 22669–22679, 2023. 1

[3] Black Forest Labs. FLUX.1-schnell: High-performance image generation model, 2024. 7

[4] Chaofeng Chen, Annan Wang, Haoning Wu, Liang Liao, Wenxiu Sun, Qiong Yan, and Weisi Lin. Enhancing diffusion models with text-encoder reinforcement learning. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2024. 2

[5] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. 5

[6] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *NeurIPS*, 2021. 1

[7] 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. *arXiv preprint arXiv:2403.03206*, 2024. 1

[8] Luca Eyring, Shyamgopal Karthik, Karsten Roth, Alexey Dosovitskiy, and Zeynep Akata. Reno: Enhancing one-step text-to-image models through reward-based noise optimization. *Neural Information Processing Systems (NeurIPS)*, 2024. 7, 8, 1

[9] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit Haim Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. 2

[10] Yaru Hao, Zewen Chi, Li Dong, and Furu Wei. Optimizing prompts for text-to-image generation. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. 2

[11] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020. 1

[12] Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang, Yuxiang Huang, Weilin Zhao, Xinrong Zhang, Zhen Leng Thai, Kai Zhang, Chongyi Wang, Yuan Yao, Chenyang Zhao, Jie Zhou, Jie Cai, Zhongwu Zhai, Ning Ding, Chao Jia, Guoyang Zeng, Dahai Li, Zhiyuan Liu, and Maosong Sun. Minicpm: Unveiling the potential of small language models with scalable training strategies. *CoRR*, abs/2404.06395, 2024. 1

[13] Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-a-pic: An open dataset of user preferences for text-to-image generation. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. 1, 5

[14] Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander M. Rush, Douwe Kiela, Matthieu Cord, and Victor Sanh. OBELICS: an open web-scale filtered dataset of interleaved image-text documents. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. 1

[15] Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. What matters when building vision-language models? *arXiv preprint arXiv:2405.02246*, 2024. 1, 5

[16] Kimin Lee, Hao Liu, Moonkyung Ryu, Olivia Watkins, Yuqing Du, Craig Boutilier, Pieter Abbeel, Mohammad Ghavamzadeh, and Shixiang Shane Gu. Aligning text-to-image models using human feedback. *CoRR*, abs/2302.12192, 2023. 2

[17] Bo Li, Kaichen Zhang, Hao Zhang, Dong Guo, Renrui Zhang, Feng Li, Yuanhan Zhang, Ziwei Liu, and Chunyuan Li. Llava-next: Stronger llms supercharge multimodal capabilities in the wild (2024). *URL https://llava-vl. github. io/blog/2024-05-10-llava-next-stronger-llms*. 1

[18] Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li. Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. *CoRR*, abs/2407.07895, 2024. 1

[19] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 4582–4597. Association for Computational Linguistics, 2021. 6

[20] Youwei Liang, Junfeng He, Gang Li, Peizhao Li, Arseniy Klimovskiy, Nicholas Carolan, Jiao Sun, Jordi Pont-Tuset, Sarah Young, Feng Yang, Junjie Ke, Krishnamurthy Dj Dvijotham, Katherine M. Collins, Yiwen Luo, Yang Li, Kai J. Kohlhoff, Deepak Ramachandran, and Vidhya Navalpakkam. Rich human feedback for text-to-image generation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*, pages 19401–19411. IEEE, 2024. 2

[21] 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. *CoRR*, abs/2406.04314, 2024. 2

[22] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*, pages 26286–26296. IEEE, 2024. 1

[23] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge. 2024. 1

[24] Wenyi Mo, Tianyu Zhang, Yalong Bai, Bing Su, Ji-Rong Wen, and Qing Yang. Dynamic prompt optimizing for text-to-image generation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*, pages 26617–26626. IEEE, 2024. 2

[25] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. In *ICML*, 2022. 1

[26] PromptHero. Openjourney: An open-source journey to ai art generation, 2023. 7

[27] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, pages 8748–8763. PMLR, 2021. 5

[28] Yuxi Ren, Xin Xia, Yanzuo Lu, Jiacheng Zhang, Jie Wu, Pan Xie, Xing Wang, and Xuefeng Xiao. Hyper-sd: Trajectory segmented consistency model for efficient image synthesis. *arXiv preprint arXiv:2404.13686*, 2024. 1

[29] Robin Rombach, A. Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10674–10685, 2022. 1

[30] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*, pages 22500–22510. IEEE, 2023. 2

[31] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L. Denton, Seyed Kamyar Seyed Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, Jonathan Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. In *NeurIPS*, 2022. 1

[32] Sogand Salehi, Mahdi Shafiei, Teresa Yeo, Roman Bachmann, and Amir Zamir. Viper: Visual personalization of generative models via individual preference learning. *CoRR*, abs/2407.17365, 2024. 1, 2, 5

[33] Axel Sauer, Frederic Boesel, Tim Dockhorn, Andreas Blattmann, Patrick Esser, and Robin Rombach. Fast high-resolution image synthesis with latent adversarial diffusion distillation. *arXiv preprint arXiv:2403.12015*, 2024. 1

[34] Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin Rombach. Adversarial diffusion distillation. In *Computer Vision - ECCV 2024 - 18th European Conference, Milan, Italy, September 29-October 4, 2024, Proceedings, Part LXXXVI*, pages 87–103. Springer, 2024. 7

[35] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021. 5

[36] Xiaoteng Shen, Rui Zhang, Xiaoyan Zhao, Jieming Zhu, and Xi Xiao. PMG : Personalized multimodal generation with large language models. In *Proceedings of the ACM on Web Conference 2024, WWW 2024, Singapore, May 13-17, 2024*, pages 3833–3843. ACM, 2024. 1, 2

[37] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, 2015. 1

[38] 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 *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*, pages 8228–8238. IEEE, 2024. 2

[39] 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. *CoRR*, abs/2306.09341, 2023. 1, 2, 7

[40] Xiaoshi Wu, Keqiang Sun, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score: Better aligning text-to-image models with human preference. In *IEEE/CVF International Conference on Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023*, pages 2096–2105. IEEE, 2023. 1, 2

[41] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. 1, 2, 5

[42] 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 *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*, pages 8941–8951. IEEE, 2024. 2

[43] Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, Qianyu Chen, Huarong Zhou, Zhensheng Zou, Haoye Zhang, Shengding Hu, Zhi Zheng, Jie Zhou, Jie Cai, Xu Han, Guoyang Zeng, Dahai Li, Zhiyuan Liu, and Maosong Sun. Minicpm-v: A GPT-4V level MLLM on your phone. *CoRR*, abs/2408.01800, 2024. 1

[44] Yuanhan Zhang, Bo Li, haotian Liu, Yong jae Lee, Liangke Gui, Di Fu, Jiashi Feng, Ziwei Liu, and Chunyuan Li. Llava-next: A strong zero-shot video understanding model. 2024. 1

[45] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *Int. J. Comput. Vis.*, 130(9):2337–2348, 2022. 6

# User-Specific Preference Prediction on Generated Images

## Supplementary Material

In this supplementary material, we provide comprehensive additional resources to further support our research. These include representative training samples, additional qualitative results to illustrate the model's behavior. Furthermore, we provide an in-depth description of experimental setups for reproducibility, as well as extended discussions to offer deeper insights into the implications and potential improvements of our approach.

## A. Examples of Training Data

Our dataset, based on Pick-a-Pic v2 dataset [13], focuses on image pairs annotated with user preferences. To ensure reliability, we filtered entries to include only users with at least 11 unique liked images. Fig. 11 and Fig. 12 present a selection of the training set from the dataset, providing valuable insights into how user-specific preferences. Patterns distinguishing a user's likes and dislikes are evident.

## B. More Qualitative Analysis Results

We present a focused comparison between our model and ViPer [32], supported by qualitative results in Fig. 9, where target images with green borders indicate preferences aligned with the user. Unlike ViPer, which primarily relies on explicit features from reference images, our method leverages Vision-Language Models (VLMs) to capture deeper semantic relationships in user preferences. By leveraging learnable preference tokens, our approach captures both shared and individual preferences, enhancing prediction accuracy and robustness across seen and unseen users. Unlike ViPer, it integrates attention-based interactions and a tailored loss design, improving alignment with nuanced user preferences and boosting generalization.

To further demonstrate the effectiveness of our approach, we present additional qualitative results on the HPSv2 benchmark, as shown in Fig. 10. These results validate its ability to generate high-quality, user-aligned images, surpassing existing methods in both semantic relevance and visual realism.

**Preference Tokens Interpretation.** As shown in Tab. 7, the learned preference tokens are not inherently tied to natural language semantics. Instead, they exist as multimodal representations optimized for capturing user preferences in a vision-language space. The relatively uniform distances suggest that the tokens are distributed in a structured but non-textual manner, emphasizing their role as functional embedding constructs rather than interpretable linguistic elements.

| Index | Top 1 | Top 2 | Top 3 |
|---|---|---|---|
| 1 | N/A (36.664) | N/A (36.664) | badly (36.664) |
| 2 | Lewis (36.598) | sock (36.598) | differently (36.599) |
| 3 | N/A (36.934) | consider (36.935) | N/A (36.935) |
| 4 | N/A (37.094) | exposed (37.095) | N/A (37.095) |
| 5 | N/A (36.996) | SUPPORT (36.997) | therapy (36.997) |
| 6 | N/A (37.221) | judge (37.222) | N/A (37.222) |
| 7 | N/A (36.649) | EqualTo (36.650) | N/A (36.651) |
| 8 | N/A (37.181) | nullptr (37.182) | dull (37.182) |
| 9 | eaten (36.466) | N/A (36.467) | N/A (36.468) |
| 10 | N/A (37.089) | SUPPORT (37.089) | N/A (37.089) |

Table 7. Visualization of learnable preference tokens with the length of 10. We derive the words by measuring the Euclidean distances between word embeddings and preference token embeddings, and the quantified distances are shown in parentheses. N/A represents non-Latin characters.

## D. More Experimental Details

**Image Generation Guided by Our Models.** Following the method outlined in [8], we assign the weight 0.75 to our model. The initial image is optimized over 30 steps. For our model, we replace non-differentiable components of the vision preprocessor such as numpy-based resizing and similar operations with PyTorch operations. The preprocessed image is then integrated into the model's input for optimization, ensuring that gradients flow seamlessly from the output score back to the initial image.

**Visualization of Attention Scores.** After applying the softmax operation in the self-attention mechanism, we extract attention weights, which are used to compute the weighted average within the self-attention heads. For visualization, we use the attention scores from head No. 28.

## E. Discussion

As shown in Fig. 8, the model effectively leverages the "like" or "dislike" signal to refine outputs that capture user preferences. These qualitative improvements highlight the potential of integrating advanced reward structures driven by user-specific feedback. Future work could focus on expanding the framework to incorporate dynamic preference modeling, enabling it to adapt to evolving user tastes over time. Furthermore, enhancing the multimodal capabilities of VLMs to include temporal data could improve the system's ability for more context-aware personalization. By leveraging sequential user interactions and historical behavior, the system could provide a deeper understanding of nuanced preferences, paving the way for even greater alignment between generated content and user expectations.
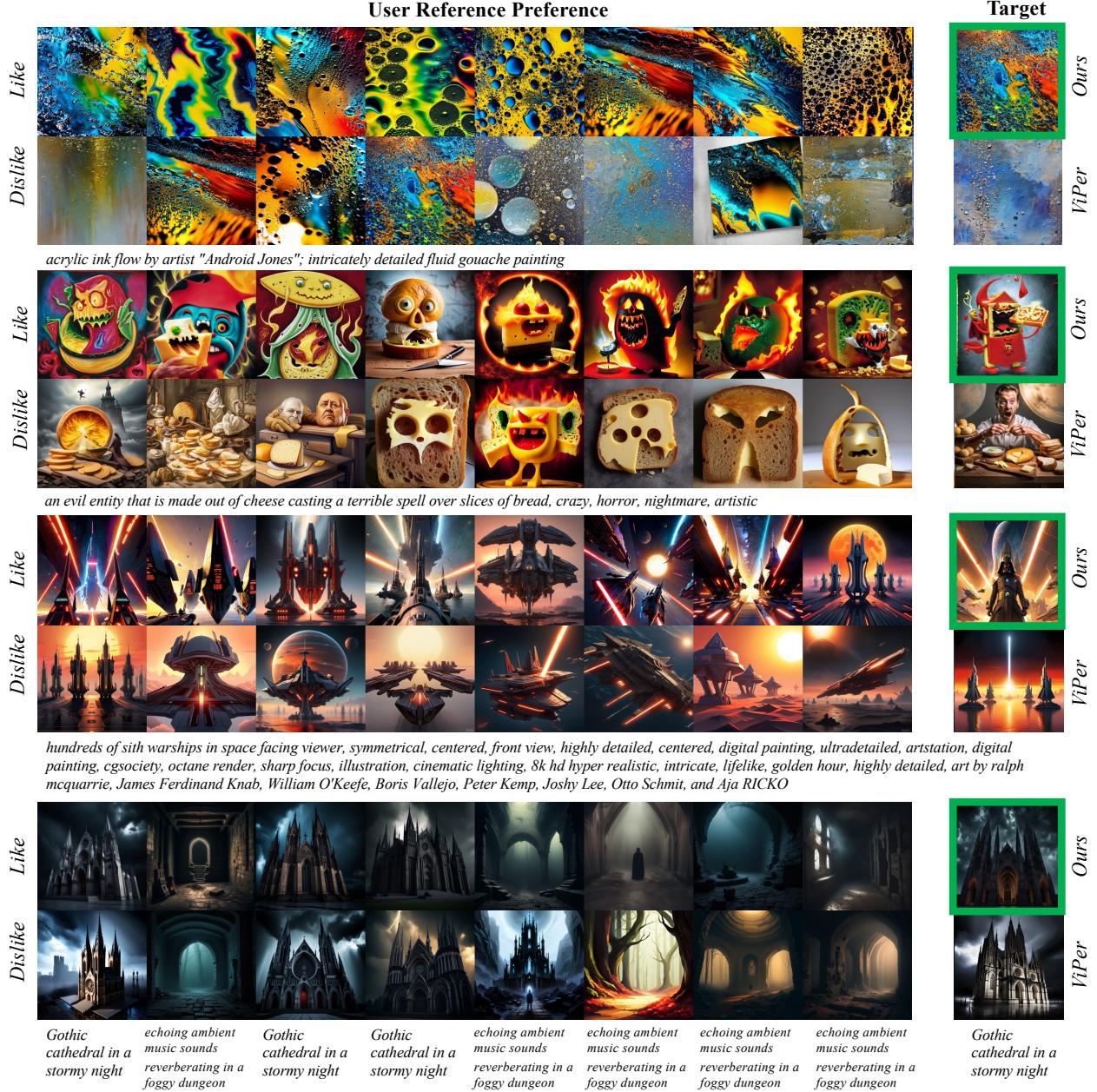
Figure 9. Visual comparison of user-specific preference alignment between our model and ViPer [32] across varying preferences. Target images with green borders indicate preferences aligned with the user. Our method demonstrates effective capture of user-specific personalized results.
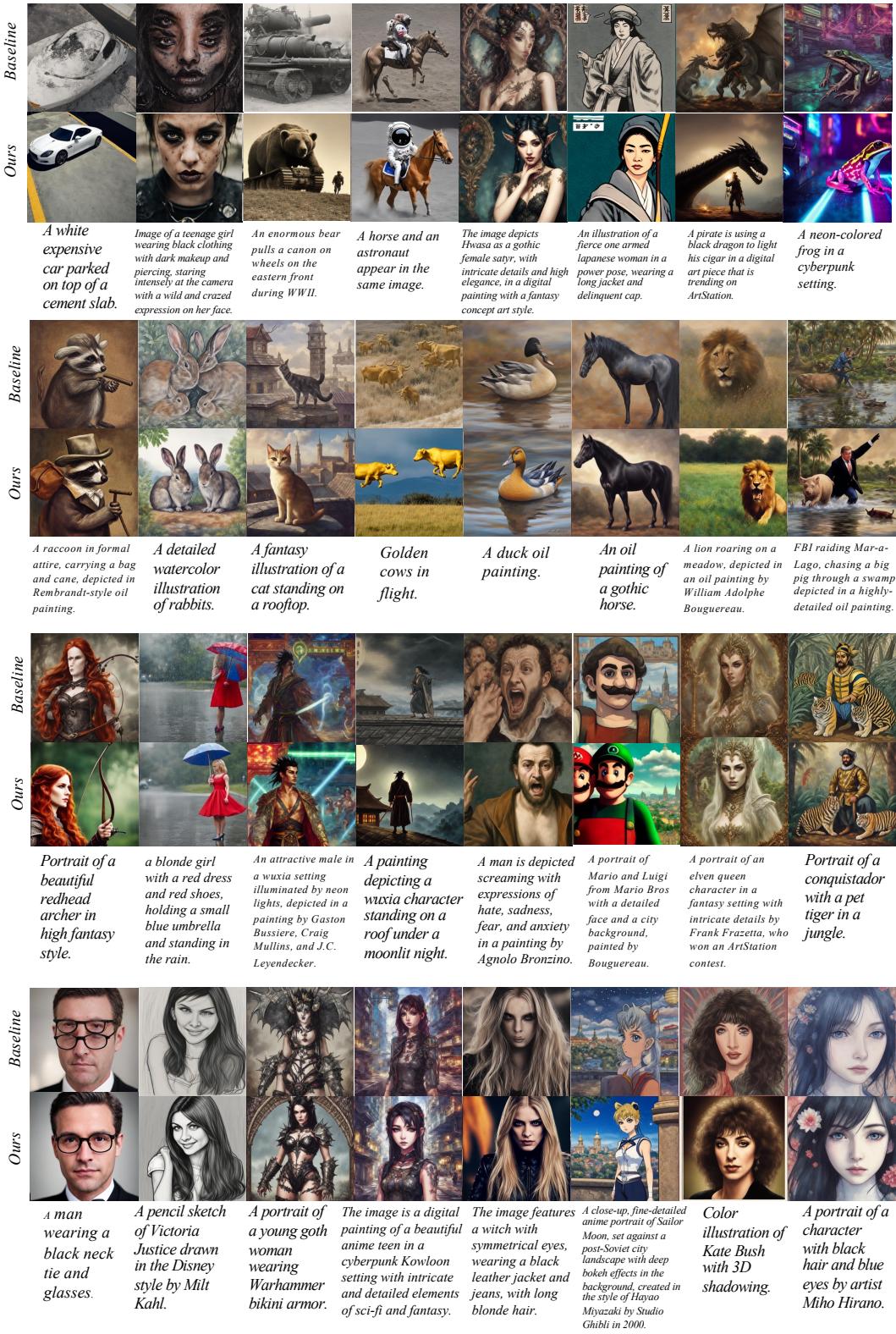
Figure 10. Some examples of images generated on HPSv2 benchmark.

Figure 11. Some examples of the training data.

**User Reference Preference**      **Target**



*She wears lilacs in her hair, and picks roses and picks daisys by artist Ralph Horsley*

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| *motion blur, heavy rain, street photo of an 30 yo Asian woman with short hair, she is laughing* | *a candid shot of Ian McKellen as Gandalf eating soft icecream cone* | *Photo of a blonde girl, intricate cyberpunk respirator and armor* | *a tall woman with purple hair in leather, alcohol, bar, tatooed, neon light* | *Beautiful woman standing in armour, Futuristic Cyberpunk city* | *still shot from a cyberpunk western, girl fedora firing a handgun* | *Movie still of starwars princess leah cworking as a waitress in a dinner, extremely detailed, intricate, high resolution, hdr, trending on artstation* | *SF movie, movie still of a young astronaut fighter pilot, round helmet, life support system, surrounded by instruments, inside a spaceship cockpit cinematic, epic, volumetric light, award winning photography, intricate details* | *Movie still of starwars princess leah cworking as a waitress in a dinner, extremely detailed, intricate, high resolution, hdr, trending on artstation* |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| *A logo of ai laptop and surveillance camera* | *A logo for a computer vision lady developer* | *A logo for a computer vision lady developer* | *A logo of a laptop and cameras* | *A logo for a computer vision lady developer* | *A logo for a computer vision lady developer* | *A logo of laptop with woman* | *A logo of a woman vollyball playe* | *A logo of laptop camera and woman* |

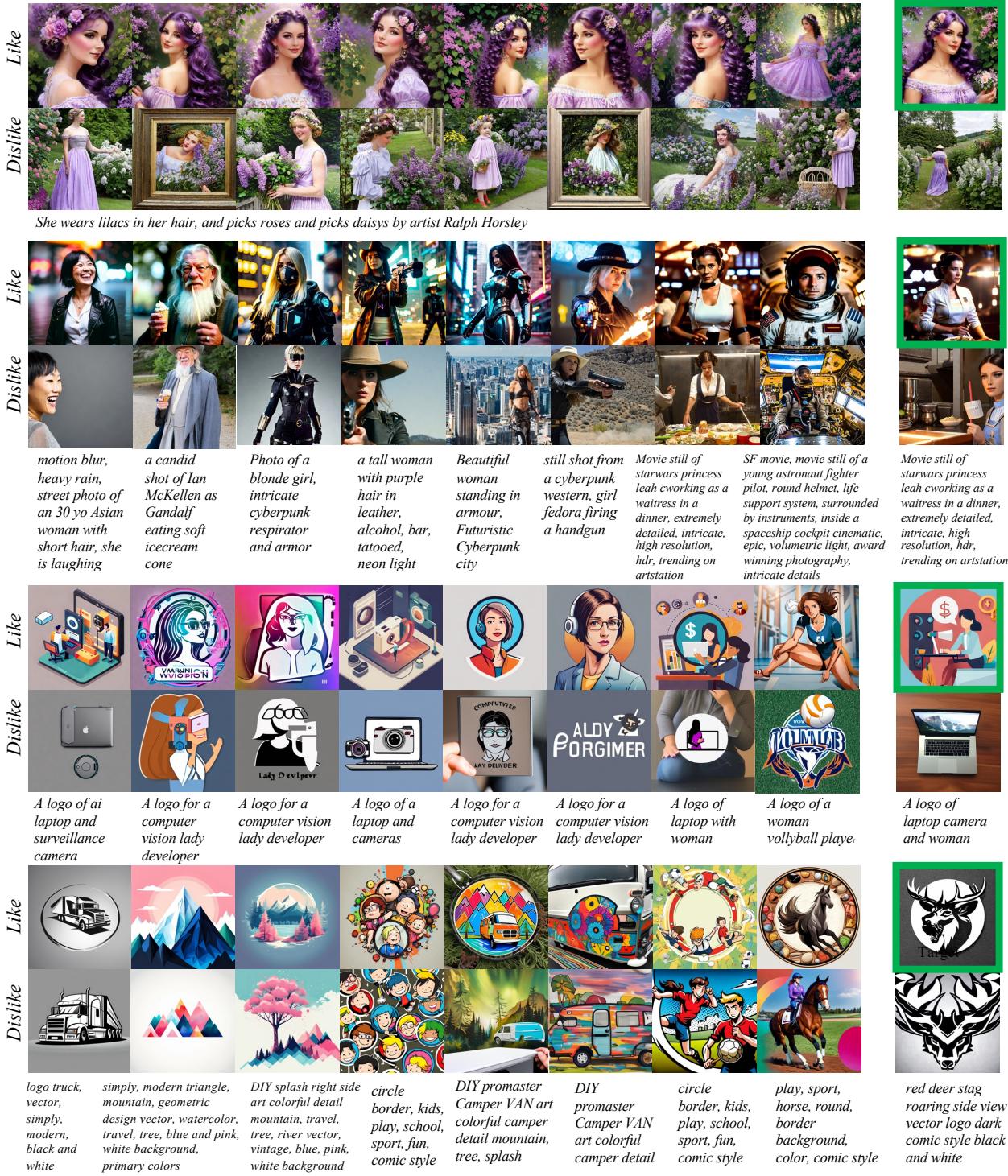| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| *logo truck, vector, simply, modern, black and white* | *simply, modern triangle, mountain, geometric design vector, watercolor, travel, tree, blue and pink, white background, primary colors* | *DIY splash right side art colorful detail mountain, travel, tree, river vector, vintage, blue, pink, white background* | *circle border, kids, play, school, sport, fun, comic style* | *DIY promaster Camper VAN art colorful camper detail mountain, tree, splash* | *DIY promaster Camper VAN art colorful camper detail* | *circle border, kids, play, school, sport, fun, comic style* | *play, sport, horse, round, border background, color, comic style* | *red deer stag roaring side view vector logo dark comic style black and white* |

Figure 12. Some examples of the training data.