# Adaptive Preference Learning for Personalized Image Generation with Vision-Language Understanding

Anonymous Submission

## Abstract

*Predicting individual preferences for personalizing image generation presents a captivating challenge. Current methods often focus on superficial features such as color and style, failing to capture the deeper content-driven connections that make personalization meaningful. Our approach reimagines this process by utilizing the semantic understanding of Vision-Language Models (VLMs) to move beyond surface-level attributes, capturing the essence of what users genuinely value in visual content. To develop a richer and more nuanced representation of individual preferences, we identify shared interests among users by leveraging latent preference prototypes. This helps distinguish each user's unique tastes and draws insights from users with similar preferences, enhancing the personalization experience. We construct a personalized preference dataset according to user IDs from the Pick-a-Pic dataset, encapsulating detailed user preferences. Experimental results show that our model outperforms PickScore by 3.96% in preference recognition accuracy and effectively identifies users with similar tastes, providing more accurate guidance for generative models to produce images aligned with specific user preferences.*
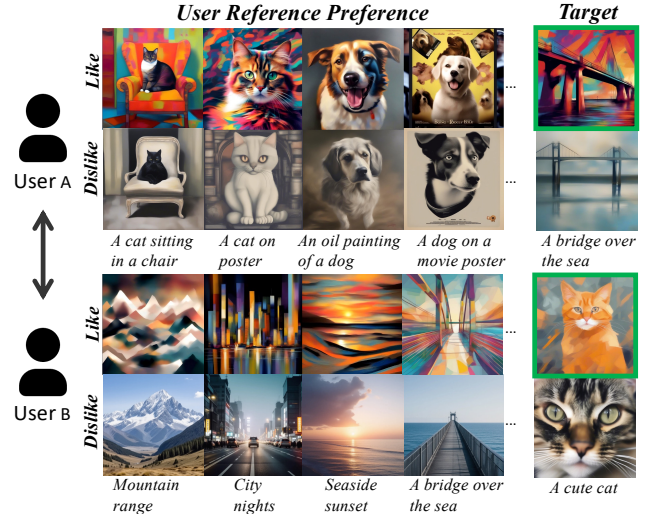


Figure 1. Illustrating the task of personalized image generation by leveraging individual and shared preferences. For each user, a set of liked and disliked images serves as reference data to model distinct tastes. Users $A$ and $B$ demonstrate unique preferences but also share some common interests, allowing the model to capture both personalized and overlapping preferences. Our task aims to predict target images that align with each user's tastes while considering these shared interests.

## 1. Introduction

Recent advancements in generative models [1, 4, 5, 9, 22, 24, 26, 27, 29, 31, 35] have brought significant progress to content creation across various domains, ranging from text to images. Among these developments, personalized content generation is increasingly attracting attention due to its ability to produce outputs tailored to individual preferences. Personalized generation has the potential to offer unique user experiences, as it enables generative models to align their outputs with specific user interests, thereby improving user satisfaction and engagement.

In personalized image generation, the objective is to generate content that is consistent with the user's tastes by using user-provided reference data, typically a set of liked and disliked images. Our task is to predict individual preferences with user reference preference and shared relations among users who have similar preferences. As shown in Fig. 1, each user possesses distinct tastes, yet a shared preference emerges among the group. While existing models such as PickScore [11], ImageReward [39], and HPS [37, 38] are designed to evaluate human preferences at a general level, they fall short when it comes to understanding individual-level preferences. Moreover, recent individual-level personalized preference modeling [30, 34] presents two 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) fail to recognize the connections between users, despite the likelihood that users with similar tastes might share preferences for certain types of images.

To address these challenges, we propose an approach that leverages Vision-Language Models (VLMs) [10, 12,

1

13, 15, 16, 19–21, 41, 42] to understand and extract deeper content-driven commonalities from user reference images. By utilizing the contextual multimodal understanding capabilities of VLMs, we move beyond superficial attributes and identify meaningful connections within the visual content. Additionally, we introduce latent preference prototypes, which models user-user relationships by utilizing prototypes to represent shared interests among users with similar preferences. This module allows us to incorporate the preferences of similar users as reference points, thereby aiding in the identification and modeling of new users' preferences. Our main contributions are as follows:

- We design a VLM-based latent preference prototype learning framework to utilize the multimodal contextual abilities of VLM in user preference modeling. Moreover, this is the first work to capture shared interests among users to enhance user preference prediction by introducing latent preference prototypes.
- We construct a personalized preference dataset based on user IDs from the Pick-a-Pic dataset [11], facilitating more accurate preference recognition and effectively guiding generative models to produce images aligned with user-specific tastes.
- Experimental results demonstrate that our model outperforms existing methods, such as PickScore, in preference recognition accuracy and successfully identifies users with similar tastes, offering a robust solution for personalized image generation.

## 2. Related Work

Our work focuses on user preference learning and personalized image generation. Previous work can be divided into two categories, one focusing on general human preferences, and the other on user-specific preferences.

**Personalizing Image Generation Based on General Preferences.** (1) Fine-tuning generative models with examples filtered by scores is an effective way to enhance their alignment with user preferences. By selecting training data based on human feedback scores or automated metrics, models are exposed to high-quality examples that reflect specific user demands. For instance, Liang *et al.* [17] 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 [37, 38] refine this idea by introducing a scoring mechanism to prioritize image-text pairs that closely align with user preferences, making the model more responsive to varied user expectations. (2) Reward-weighted fine-tuning is another promising strategy for aligning text-to-image models with human preferences. In this approach, models are fine-tuned using reward signals that weigh more heavily on user satisfaction. Lee *et al.* [14] exemplifies this by incorporating feedback-based rewards

during training, which enables the model to generate outputs that align with user preferences. Furthermore, ImageReward [39] 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 help tailor the model's outputs to reflect diverse and nuanced user tastes. (3) Reinforcement learning has been widely applied in generative models to better align outputs with human preferences [2, 8, 18, 23]. Recent work [8, 23] uses reinforcement learning to optimize the input prompts to get high-quality images. DiffusionDPO [36] leverages user preferences to iteratively fine-tune the model, improving its ability to generate images that reflect user choices. Similarly, D3PO [40] utilizes dynamic updates based on evolving user feedback, ensuring the model remains adaptable to changing preferences.

**Personalizing Image Generation Based on User-Specific Preferences.** In recent advancements in personalizing image generation based on user-specific preferences, several approaches have emerged to refine how generative models align with individual needs. DreamBooth [28] and Textual Inversion [7] 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 address personalization, Salehi *et al.* [30] proposes a standardized process for collecting 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.* [34] introduces a method for integrating user-specific preferences across different modalities, such as text and images, creating personalized outputs by leveraging historical interactions like clicks and conversations. This multimodal approach significantly enhances the adaptability of models, enabling them to better align with user needs. In our work, we extend this by utilizing VLMs to understand both textual and visual content deeply. By mining the deeper commonalities in image content through VLM capabilities, we aim to uncover key patterns that reflect user preferences. Additionally, we employ prototypes to capture shared preferences across users, using the preferences of similar users as a reference to help identify individual user preferences more effectively. This method personalizes content by leveraging collective data to improve its ability to generalize across users with similar preferences.

## 3. Method

Given a user's preference history sequences $\mathcal{S} = \{s_i | s_i = (I_{\text{pos}}^{(i)}, I_{\text{neg}}^{(i)}, T^{(i)})\}_{i=1}^{N_{\text{ref}}}$, we aim to evaluate whether a new target item $z = (I, T)$ aligns with a user's preferences. Each
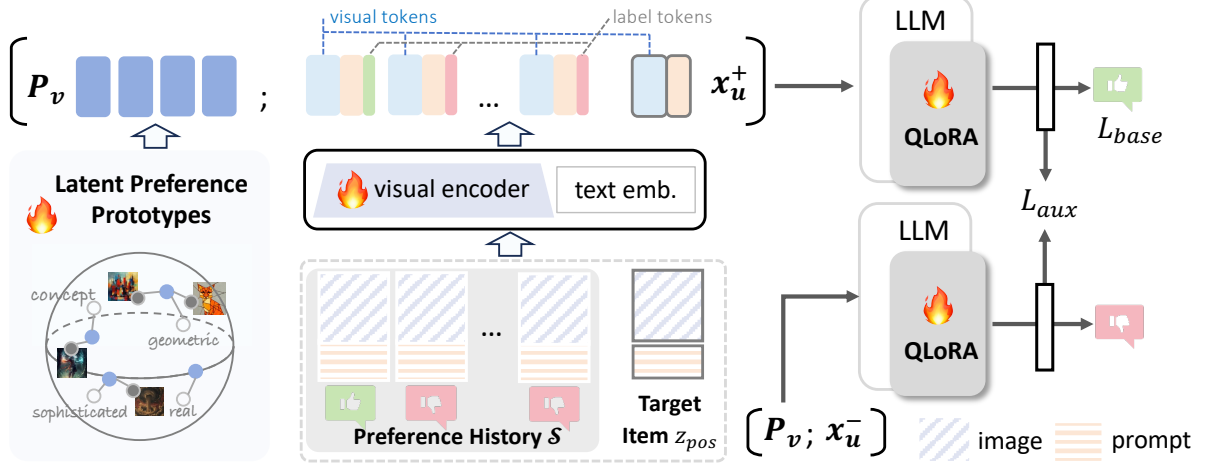
Figure 2. Our proposed VLM-based latent preference prototype learning framework leverages user preference history $\mathcal{S}$ to capture nuanced preferences through visual and textual encodings. User preference history samples, labeled as "like" (+) or "dislike" (−), along with target items $z$, are processed to generate user-specific representations $x_u$. The framework is trained with a base loss $L_{base}$ that guides the model in predicting user preferences, and an auxiliary loss $L_{aux}$ that focuses on relative ranking among preferences. Additionally, latent preference prototypes $P_v$ capture and model shared interests across users, enriching personalization and enabling more tailored preference prediction.

reference entry in $\mathcal{S}$ consists of an image that the user likes, denoted as $I_{\text{pos}}^{(i)}$, and dislikes, denoted as $I_{\text{neg}}^{(i)}$, with an optional text prompt $T^{(i)}$ describing the image content. For the target item $z$, user preference is represented by $z_{\text{pos}}$ (if the user likes it) and $z_{\text{neg}}$ (if the user dislikes it). The length of the user's preference history sequence is denoted by $N_{\text{ref}}$.

### 3.1. Overview

As shown in Fig. 2, to predict individual visual preferences, we propose a VLM-based latent preference prototypes learning framework. The transformer-based VLM [13] is utilized to capture deeper, meaningful representation that reflects what users genuinely value in visual content. Furthermore, latent preference prototypes are introduced to model user-user relationships and extract shared interests, thereby enhancing the personalization process by drawing insights from other users with similar preferences.

The overall architecture consists of three main components: a visual encoder, a connector, and a language model. Given an input image $I \in \mathbb{R}^{H \times W \times C}$, the visual encoder and connector extract visual tokens $\mathbf{x}_v$. We use two labels "+" and "−" to represent "like" and "dislike", respectively. Visual tokens $x_v$, text tokens $x_t$ and a label token $x_{\text{label}}$ are concatenated to form the input sequence for one reference entry. All reference entries in the history sequences $\mathcal{S}$, except the last label token, are stacked as the user-specific input token sequence, denoted by $x_u$. The last label token will be the objective of next-token prediction.

### 3.2. Latent Preference Prototypes

**Attention-based Interaction with Prototypes.** To integrate user preferences into the generated content, we use the

learnable latent preference prototypes $P_v \in \mathbb{R}^{L_p \times D}$ as part of the input sequence to the VLM, where $L_p$ represents the number of prototypes and $D$ is the embedding dimension. These prototypes and the user-specific input token sequence are combined to form the final input sequence:

$$x'_p = [P_v; x_u] \tag{1}$$

where $x_u \in \mathbb{R}^{L_e \times D}$ represents the embedded input tokens, $L_e$ is the length of input tokens and $[P_v; x_u]$ denotes the concatenation of the preference prototypes with the input embeddings.

**Analysis of Attention Mechanism.** The core of the interaction lies in the attention mechanism of transformer layers, where the interaction between the input tokens and the preference prototypes is carried out. Let the attention weights related to prototypes be represented by:

$$\mathcal{A} = \text{softmax}\left(\frac{W_q(x_u)W_k(P_v)^T}{\sqrt{D'}}\right) \tag{2}$$

where $W_q$, $W_k$ and $W_v$ are linear projections in the attention mechanism, $\mathcal{A} \in \mathbb{R}^{L_e \times L_p}$ represents the attention scores between each input token and the prototypes. These attention scores are used to compute a weighted representation of the preference prototypes:

$$\tilde{x}_p = \mathcal{A} \times W_v(P_v) \tag{3}$$

The attended features $\tilde{x}_p \in \mathbb{R}^{L_e \times D'}$ are then used to adjust the input representation.

This feature-to-prototype similarity mapping, which is implemented implicitly in the basic attention mechanism in
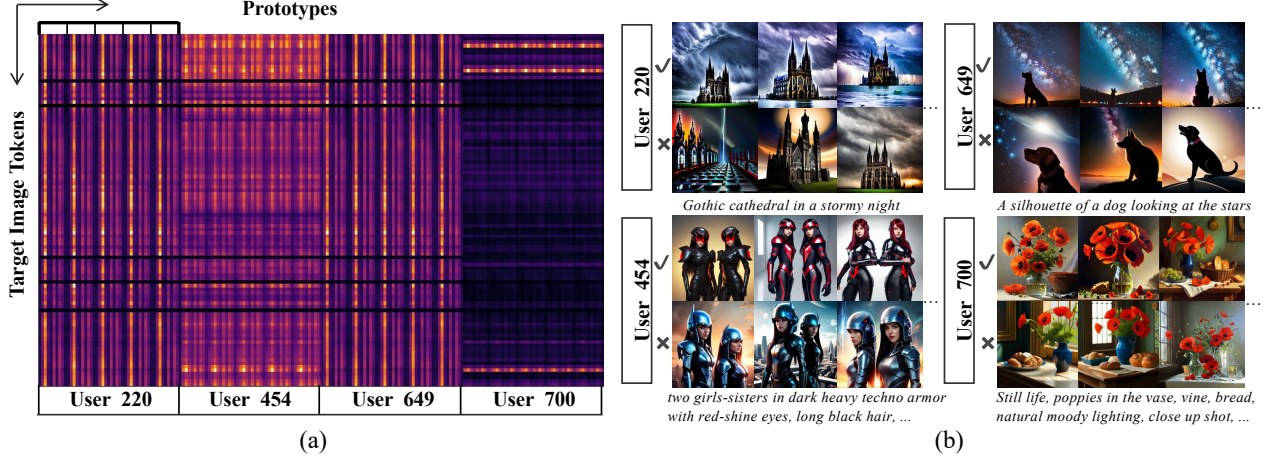
Figure 3. (a) Attention scores $\mathcal{A}$ display the interactions between prototypes and target image tokens for individual users, showing how our prototypes capture user-specific features within a shared semantic space. Each user has a unique reference history, and we concatenate the same target image to the input sequence across users to highlight these individualized interactions. For each user, the vertical axis represents tokens from target image, while the horizontal axis represents the prototypes. Each user has five different random reorderings of reference images. User 220 and User 649 share certain features, which distinguish them from other users. (b) Examples of images liked (✓) or disliked (✗) by each user.

VLMs, projects user-specific features onto a shared space, effectively modeling user preferences as a combination of multiple prototypes. This shared representation captures both common traits and individual distinctions among users, enhancing our understanding of user preferences. To validate our assumption, we provide a visual presentation of the attention scores $\mathcal{A}$ in Fig. 3, which reflects the interactions between target image and the prototypes. The reference histories of these users are different, and the same target image is concatenated to the input sequence. Specifically, Fig. 3 (a) shows similarities between Users 220 and 649, who share a distinct pattern of attention across several prototypes. Fig. 3 (b) further supports this, showing that both Users 220 and 649 prefer similar themes, including landscapes with dramatic skies, silhouettes, and nightscapes, pointing to their shared visual aesthetic. In contrast, Users 454 and 700 exhibit distinct differences in their alignment with prototypes, highlighting the diversity of preferences captured by the shared space. These differences illustrate how the model can differentiate users with varying tastes while retaining unique characteristics. Further details of the experiment can be found in the appendix.

### 3.3. Adaptive Preference Prediction with Vision-Language Understanding

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 $z$. We define a loss function that combines a base classification loss with auxiliary losses to improve the model's ability to distinguish between "like" and "dislike" predictions.

**Base Loss.** The base loss, $L_{\text{base}}$, aims to minimize the classification error across both "like" and "dislike" samples. Let $\mathcal{M}^+(\mathcal{S}, z_i)$ and $\mathcal{M}^-(\mathcal{S}, z_i)$ represent the logit outputs for predicting "like" and "dislike" outcomes for a sample $z_i$, respectively, and let $\mathbf{y}_{\text{pos}}$ and $\mathbf{y}_{\text{neg}}$ be their corresponding labels. The base loss is defined as:

$$L_{\text{base}} = \frac{1}{2}\left(\mathcal{L}(\mathcal{M}^+(\mathcal{S}, z_{\text{pos}}), \mathbf{y}_{\text{pos}}) + \mathcal{L}(\mathcal{M}^-(\mathcal{S}, z_{\text{neg}}), \mathbf{y}_{\text{neg}})\right) \tag{4}$$

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

$$\mathcal{Q}(\mathcal{S}, z_i) = \frac{\exp(\mathcal{M}^+(\mathcal{S}, z_i))}{\exp(\mathcal{M}^+(\mathcal{S}, z_i)) + \exp(\mathcal{M}^-(\mathcal{S}, z_i))} \tag{5}$$

**Auxiliary Losses for Preference Refinement.** To complement the base loss, we introduce two auxiliary 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 falls short when only pairwise comparisons are available. Specifically, given pairwise relationships 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 relationship. Consequently, the model may predict similar outcomes for both $A$ and $B$ in terms of $\mathcal{Q}(\mathcal{S}, z_i)$, since $L_{\text{base}}$ overlooks the interplay between correct "like" predictions for disliked samples and "dislike" predictions for liked samples. To overcome this limitation, the auxiliary loss terms, $L_+$ and $L_-$, address this limitation by incorporating pairwise ranking information, thereby refining preference predictions and improving
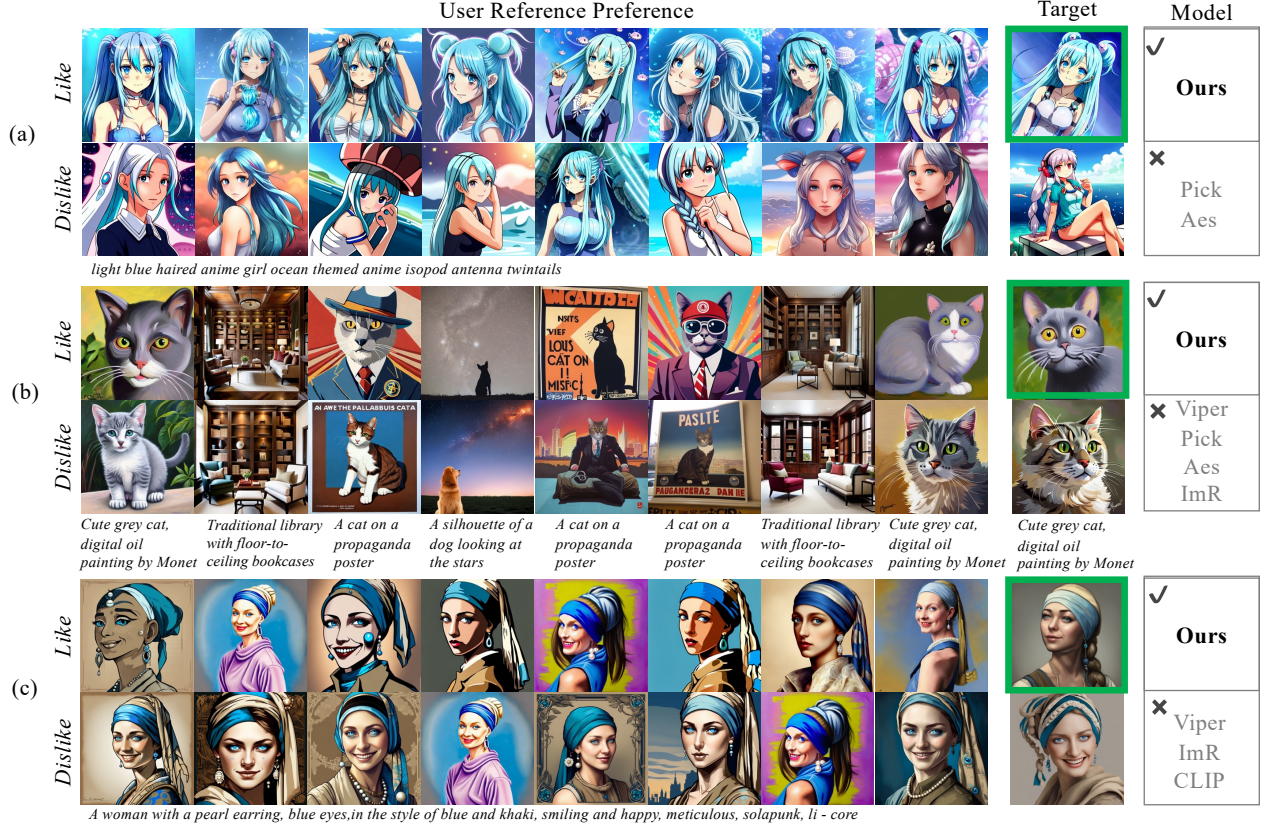
Figure 4. Qualitative comparison of user preference alignment across models. We compare our model to ViPer [30], PickScore [11], ImageReward [39], CLIP [25], and Aesthetic Score [33]. Subfigures (a), (b), and (c) illustrate user-specific preferences for style, content, and pose, respectively. Each subfigure 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.

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:

$$L_+ = -\frac{1}{N} \sum_{i=1}^{N} \log \sigma(\mathcal{M}^+(\mathcal{S}, z_{\text{pos}}) - \mathcal{M}^+(\mathcal{S}, z_{\text{neg}})) \quad (6)$$

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 to prioritize negative samples appropriately:

$$L_- = -\frac{1}{N} \sum_{i=1}^{N} \log \sigma(\mathcal{M}^-(\mathcal{S}, z_{\text{neg}}) - \mathcal{M}^-(\mathcal{S}, z_{\text{pos}})). \quad (7)$$

Then, we calculate $L_{\text{aux}} = L_+ + L_-$ to obtain the auxiliary losses. The combined loss function, which incorporates

both the base and auxiliary losses, enhances the model's ability to distinguish user preferences by refining predictions for both positive and negative outcomes:

$$L_{\text{all}} = L_{\text{base}} + L_{\text{aux}}, \quad (8)$$

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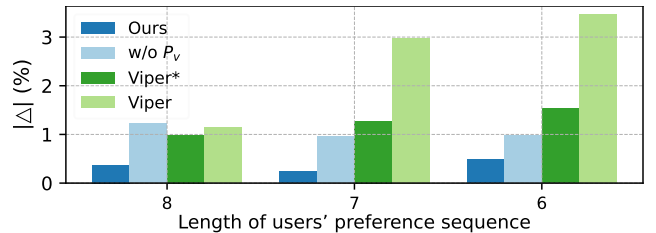
## 4. Experiments



Figure 5. Top-1 accuracy on seen-unseen data with different $N_{\text{ref}}$.

| 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.

|  | 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 2. Top-1 accuracy on seen-unseen data with $N_{\text{ref}} = 8$.

## 4.1. Experimental Setup

**Datasets.** We construct a user-specific dataset from the Pick-a-Pic v2 dataset [11], in which each user's image pairs and corresponding preference annotations are available. For each user, we filter the data to include only entries with at least 11 unique liked images, ensuring that our dataset has sufficient distinct preferences to be considered reliable. Our 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 the 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.** We use IDEFICS2-8B [13] 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 GPUs with a local batch size of 2 pairs and accumulation of gradients over 4 steps. Following the setup of [30], 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 [3], while the vision encoder is trained simultaneously. The input tokens template for the VLM is "<image>The prompt is <prompt>. Score for this image?". Initially, the VLM is trained with our custom loss function for $5,000$ steps, after which the model weights are fixed, and only the learnable prototypes are further tuned for an additional $16,000$ steps. To prevent the model from learning a fixed pattern, we randomly shuffle the order of reference history sequences during training.

**Evaluation Metric.** We evaluate our method using top-$K$ accuracy, which assesses whether the liked image is ranked among the top $K$ candidates. Among all candidates, only one "like" image is provided. In cases where one liked image is compared against one disliked image, we only employ 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 [30], (2) PickScore [11], (3) ImageReward [39], (4) CLIP [25], and (5) LAION Aesthetic Score Predictor [33]. ViPer proxy model predicts user preferences by analyzing reference images. PickScore and ImageReward focus on learning general human preferences 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 train ViPer and PickScore with the same settings as our model, and these extended versions are marked by '*' in the results.

## 4.2. Evaluation and Analysis

### 4.2.1 Quantitative Analysis.

**Quantitative User-Specific Preference Prediction.** Tab. 1 and Tab. 3 present a quantitative comparison of different models in terms of top-$K$ accuracy on our dataset. Our model consistently outperforms all other approaches, including baselines such as ViPer, PickScore, CLIP score, and ImageReward. Specifically, in Tab. 1, our model achieves the highest top-1 accuracy in like-dislike pairs, surpassing other methods. Note that PickScore (shown in gray text) is trained on the full Pick-a-Pic dataset, making it less comparable to our method. For a fairer comparison, we focus on models with an asterisk next to their names, which are trained under the same settings as our model. In Tab. 3, our model demonstrates superior performance compared to other methods, particularly in scenarios involving multiple disliked cases. These results indicate that our approach better aligns with user preferences and achieves higher overall prediction accuracy. Also, generic metrics, including Aesthetic score and CLIP score, report the worst accuracy, indicating that user-specific preferences may differ significantly from general preferences.

**Seen vs. Unseen.** For a model that can fully represent the user-specific characteristics, there should be a basically consistent accuracy on seen and unseen users. Our proposed framework utilizes learnable prototypes to build

6

Figure 6. We generate images using SD-Turbo [32], incorporating different reward models based on the methodology described in [6]. The "Like" column represents user-preferred characteristics, while the "Dislike" column indicates features that users do not prefer. The green check marks highlight the results that best align with user preferences. Our method demonstrates a higher probability of generating spacecraft flight directions that match user preferences effectively.

the basis of user preferences from the relationship between users, which has stronger robustness. As shown in Tab. 2, our latent preference prototypes learning framework effectively shrinks the performance gap between seen and unseen datasets. Moreover, from Fig. 5, our method shows a slower increase in the difference between unseen and seen accuracies with the reduction of preference sequence length. This implies that our prototypes effectively generalize across users, maintaining stable performance even with fewer history preference records. In contrast, other methods, such as ViPer, exhibit a more significant gap, which highlights their limitations in preserving user relationship modeling when user data is limited.
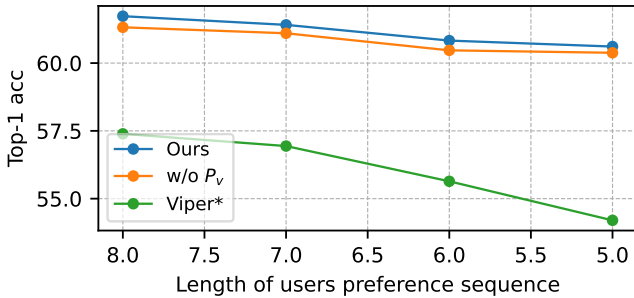
| Model | $N_{\text{ref}}$ | Top-1 Acc | Top-2 Acc | Top-3 Acc |
|---|---|---|---|---|
| Random | 0 | 25.0 | 50.0 | 75.0 |
| Aes Score | 0 | 28.11 | 54.12 | 78.33 |
| CLIP Score | 0 | 30.04 | 55.82 | 76.05 |
| ImageReward | 0 | 31.42 | 58.01 | 78.47 |
| IDEFICS | 8 | 24.40 | 51.88 | 78.33 |
| ViPer | 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 3. A quantitative comparison between the liked case and three disliked cases. We report the top-1 to top-3 accuracy (%).

| | Baseline | w/ Pmpt | w/ Pmpt & $L_{\text{aux}}$ | Full |
|---|---|---|---|---|
| Top-1 Acc (%) | 57.39 | 60.47 | 61.37 | 61.68 |

Table 4. Ablation Study for different Settings.

### 4.2.2 Qualitative Analysis.

**Qualitative User-Specific Preference Prediction.** In Fig. 4, our model effectively aligns with user-specific preferences by distinguishing styles, content, and poses 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.

**Comparative Analysis of Reward Models in Image Generation.** We use SD-Turbo [32] incorporating different reward models based on the method described in [6] to generate images combined with our model. For more information regarding the image generation settings, please refer to the supplementary materials. As shown in Fig. 6, incorporating user preference feedback significantly improves image gen-



Figure 7. Top-1 accuracy on test dataset with different $N_{\text{ref}}$.

**Number of User Reference Preferences.** As demonstrated in Fig. 7, 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.

*Like    Dislike    w/o Guidance    w/ Like    w/ Dislike*

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

(a)

*A woman with a pearl earring, blue eyes,in the style of blue and khaki, smiling and happy, meticulous, solapunk, li - core*
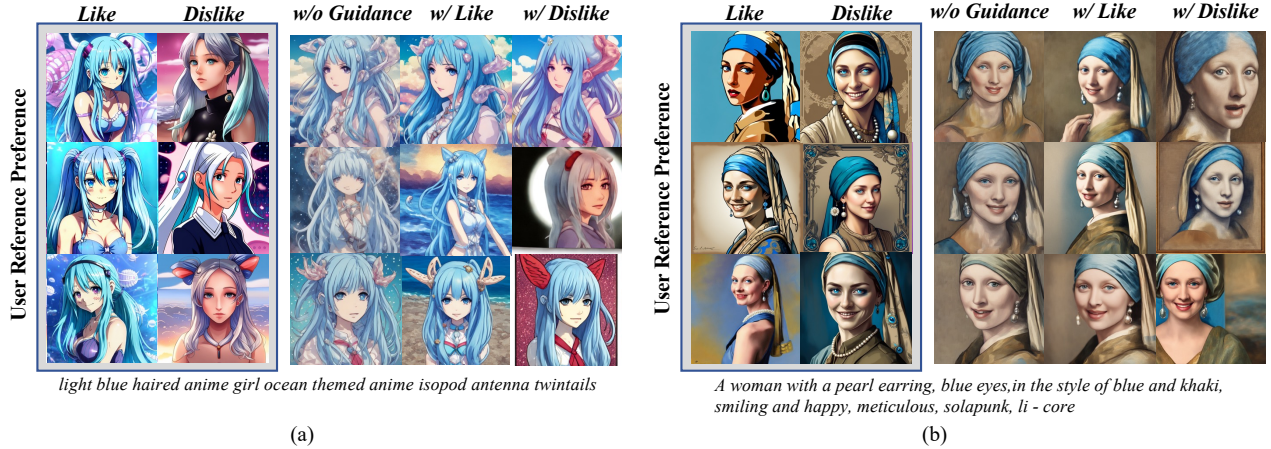
(b)

Figure 8. Using the approach from [6], we generate images guided by our model, incorporating both user likes and dislikes as feedback. By adjusting the guidance based on "Like" or "Dislike" feedback, we achieve distinct variations in the generated outputs, demonstrating the impact of user feedback on image generation results.

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

Table 5. Ablation Study for Prototype Numbers.

eration quality across different reward models. Specifically, our proposed method effectively captures user preferences by modeling both the like and dislike references, leading to noticeable improvements in the generated images. Compared to other methods, our approach demonstrates a better alignment with user preferences, particularly in maintaining consistent features, as indicated by the green check marks for our results. This experiment highlights the importance of leveraging user feedback in refining the reward structure to generate more appealing and preference-aligned images.

**Comparison of Like vs. Dislike Predictions as Reward Feedback.** In Fig. 8, we adopt the approach from [6] to generate images combined with our model's guidance. In Fig. 8 (a), the "w/ Like" setting enhances the probability of producing images with a blue-themed style that aligns with user preferences, whereas the "w/ Dislike" setting tends to emphasize an unwanted pink hue. Similarly, in Fig. 8 (b), incorporating "Like" feedback increases the chances of generating poses that match user preferences. These results demonstrate that our model is capable of effectively capturing both user preferences and dislikes, allowing for a more comprehensive modeling of user preferences.

### 4.3. Ablation Study

**Effect of Textual Description, Auxiliary Loss and Learnable Prototypes.** As shown in Tab. 4, incorporating prompts results in a 3.08% accuracy improvement, highlighting the beneficial effect of providing textual context to enhance the model's understanding and predictions. The in-

clusion of the auxiliary loss term, $L_{aux}$, further improves the accuracy by 0.9%, suggesting that the auxiliary loss term aids in refining the model's alignment with the target outputs. Lastly, the full model, which incorporates learnable prototypes in addition to the previous components, achieves an accuracy of 61.68%, representing a cumulative performance gain over the baseline model.

**Ablation Analysis of Prototype Hyper-Parameters.** As shown in Tab. 5, the results indicate that using 10 prototypes yields the highest top-1 accuracy, slightly outperforming configurations with 5 and 20 prototypes. These findings suggest that an optimal number of prototypes is crucial for effectively modeling user preferences without overfitting or underrepresenting the variation in data.

## 5. Conclusion

In this paper, we advance the prediction of user-specific preferences by harnessing the contextual understanding capabilities of Vision-Language Models to capture more nuanced, content-driven user preferences. By introducing latent preference prototypes, our approach enhances the modeling of user-to-user preference connections, allowing for more accurate predictions of individual tastes and significantly improving the quality of content personalization. Comprehensive experiment results validate the effectiveness of our method, demonstrating that it surpasses existing models in both preference recognition accuracy and the depth of content personalization it provides.

## References

[1] 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

[2] 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

[3] 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. 6

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

[5] 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

[6] 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, 2

[7] 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

[8] 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

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

[10] 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

[11] 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, 2, 5, 6

[12] 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

[13] Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. What matters when building vision-language models? *arXiv preprint arXiv:2405.02246*, 2024. 2, 3, 6

[14] 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

[15] 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*. 2

[16] 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. 2

[17] 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

[18] 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

[19] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *ArXiv*, abs/2304.08485, 2023. 2

[20] 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.

[21] 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. 2

[22] Xiaoxiao Ma, Mohan Zhou, Tao Liang, Yalong Bai, Tiejun Zhao, Huaian Chen, and Yi Jin. STAR: scale-wise text-to-image generation via auto-regressive representations. *CoRR*, abs/2406.10797, 2024. 1

[23] 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

[24] 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

[25] 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, 6

[26] 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

[27] 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

[28] 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

[29] 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

[30] 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, 6, 3

[31] 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

[32] 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

[33] 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, 6

[34] 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

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

[36] 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

[37] 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

[38] 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

[39] 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, 6

[40] 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

[41] 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. 2

[42] 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. 2

# Adaptive Preference Learning for Personalized Image Generation with Vision-Language Understanding

## 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, and a detailed analysis of failure cases to highlight challenges and limitations. 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 [11], 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. 12 and Fig. 13 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

**User-Specific Preference Prediction Comparison with ViPer.** In this section, we present a focused comparison between our model and ViPer [30], supported by qualitative results in Fig. 11, 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 introducing latent preference prototypes, our approach effectively models shared and individual preferences, achieving notable improvements in both prediction accuracy and robustness across seen and unseen users. As a concurrent work with ViPer, our method takes a distinct approach by incorporating attention-based interactions with learnable prototypes and a tailored loss design, which enhance alignment with nuanced user preferences and improve generalization.

**t-SNE Visualization of Feature Representations.** In Fig. 9, the t-SNE visualization showcases the feature representations of all target images liked by different users in the test set, as processed by our model. These embeddings are extracted from the transformer outputs before the final linear layer. The visualization reveals that images liked by the same user tend to form tighter clusters, demonstrating the model's effectiveness in capturing user-specific preferences. Among them, there are associations between differ-
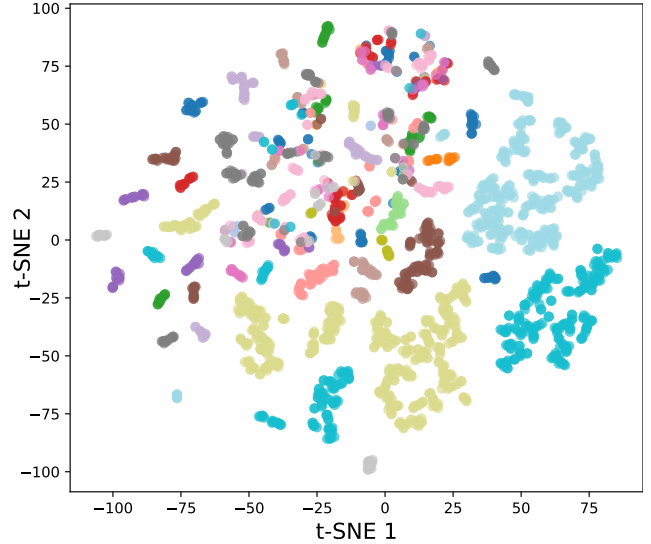


Figure 9. t-SNE visualization of feature representations for all target images liked by different users in the test set. Each point represents the feature representation of a target image, and different colors indicate different users. The clusters reflect how the feature representations for target images liked by individual users are grouped based on similarity in user preferences.

ent users, reflecting possible common preference patterns. This distribution emphasizes the ability of the model to represent and distinguish between user-specific characteristics and shared in the embedded space of learning.

## C. Analysis of Bad Cases

We conduct an analysis of the model's failure cases to gain deeper insights into potential areas for improvement and to inform future development directions.

**Inconsistencies between Prompts and Images.** When users evaluate the images, they may prefer images that are unrelated to the prompts, which undermines the model's ability to effectively learn the correlation between text and images. In Fig. 10 (a), the "liked" images fail to accurately reflect the content described in the prompts. These discrepancies introduce noise into the dataset. This misalignment ultimately hampers the model's performance in tasks requiring precise text-image associations.

**Malformed Low-quality Input Data.** As the dataset is synthetically generated, some images may be perceived as invalid from a human perspective, as illustrated in Fig. 10 (b). For instance, structural inconsistencies in hu-
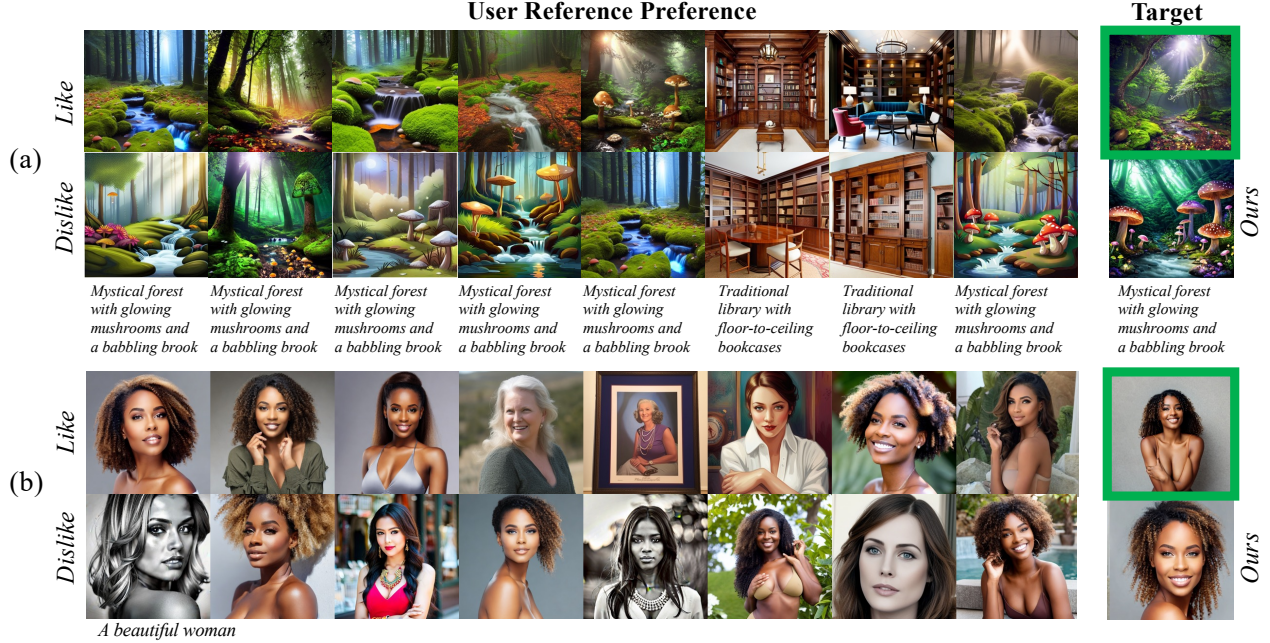
Figure 10. Examples of failure cases from our model preference prediction. Subfigure (a) demonstrates conflicts between user preferences and prompt alignments, where visually appealing but textually unrelated images are preferred. Such cases introduce noise that hinders the model's ability to learn meaningful text-image correlations. Subfigure (b) showcases malformed human faces and body structures in synthetically generated images, which negatively impact user satisfaction and lead to model confusion. These cases illustrate the difficulties in aligning generated content with both user preferences and prompt consistency, which causes bad cases.

man hands or faces occasionally occur, leading users to dislike such images. Importantly, this aversion is unrelated to the user's other preferences, such as stylistic choices. These inconsistencies can confuse the model, resulting in suboptimal performance in certain cases.

Our analysis of failure cases reveals limitations stemming from both model design and dataset quality. To address these issues, we plan to incorporate additional references based on historical user preference data to mitigate inconsistencies between prompts and images. Additionally, from the dataset perspective, introducing a data evaluation module during the collection process can effectively minimize the negative impact of low-quality data.

## D. More Experimental Details

**Image Generation Guided by Different Reward Models.**
Following the method outlined in [6], we assign the following weights to the reward models respectively: 1.0 for ImageReward, 1.0 for Aesthetic Score, 0.05 for PickScore, 1.0 for CLIP Score, and 0.75 for ViPer consistent with our approach. The initial image is optimized over 30 steps. For our method and ViPer, 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. To address GPU memory constraints, we use 3 like-dislike image pairs for both our method and ViPer.

**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.

   

*acrylic ink flow by artist "Android Jones"; intricately detailed fluid gouache painting*

*an evil entity that is made out of cheese casting a terrible spell over slices of bread, crazy, horror, nightmare, artistic*

*hundreds of sith warships in space facing viewer, symmetrical, centered, front view, highly detailed, centered, digital painting, ultradetailed, artstation, digital painting, cgsociety, octane render, sharp focus, illustration, cinematic lighting, 8k hd hyper realistic, intricate, lifelike, golden hour, highly detailed, art by ralph mcquarrie, James Ferdinand Knab, William O'Keefe, Boris Vallejo, Peter Kemp, Joshy Lee, Otto Schmit, and Aja RICKO*

*Gothic cathedral in a stormy night* — *echoing ambient music sounds reverberating in a foggy dungeon* — *Gothic cathedral in a stormy night* — *Gothic cathedral in a stormy night* — *echoing ambient music sounds reverberating in a foggy dungeon* — *echoing ambient music sounds reverberating in a foggy dungeon* — *echoing ambient music sounds reverberating in a foggy dungeon* — *echoing ambient music sounds reverberating in a foggy dungeon*     *Gothic cathedral in a stormy night*
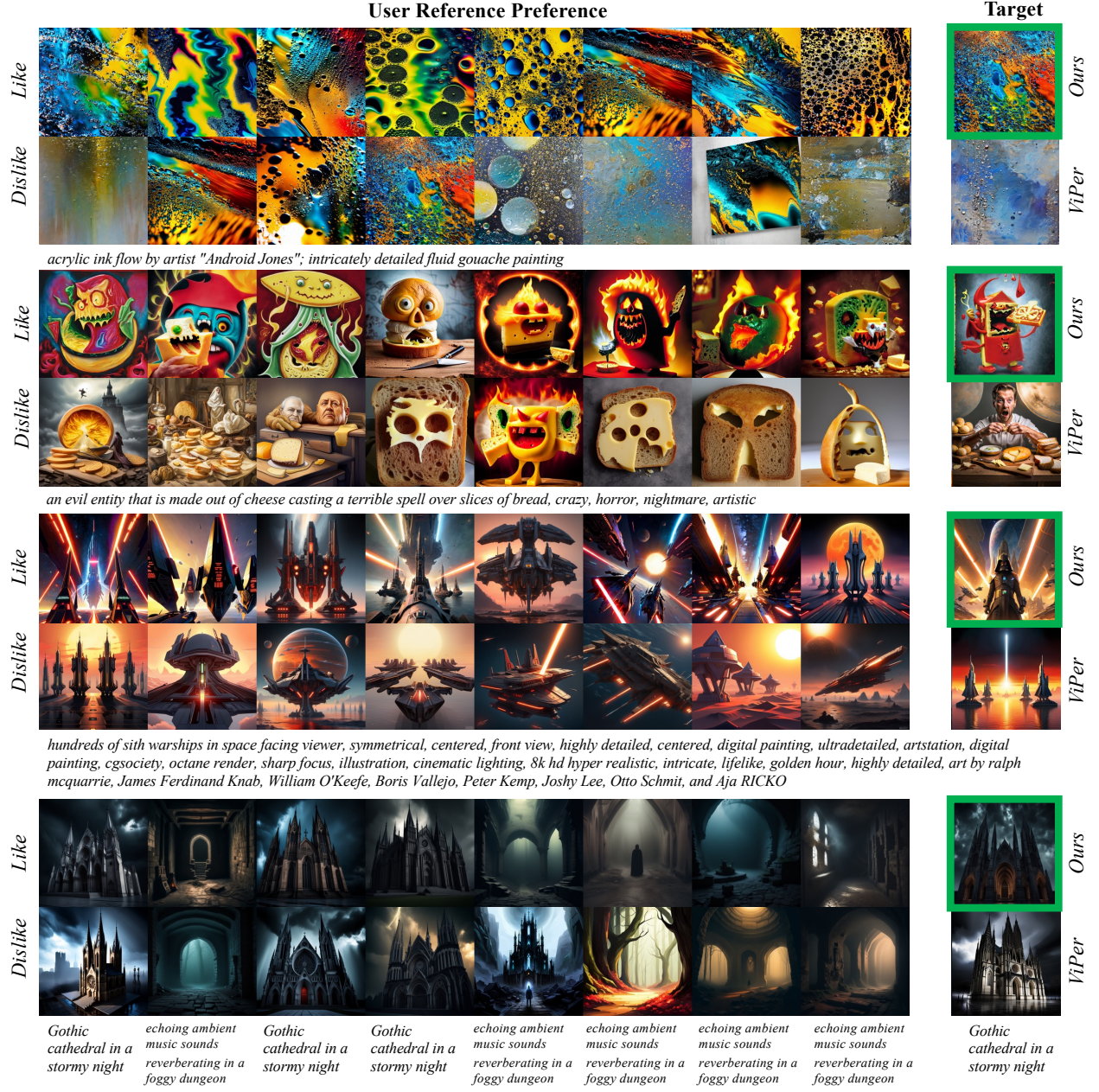
Figure 11. Visual comparison of user-specific preference alignment between our model and ViPer [30] across varying preferences. Target images with green borders indicate preferences aligned with the user. Our method demonstrates effective capture of user-specific personalized results.
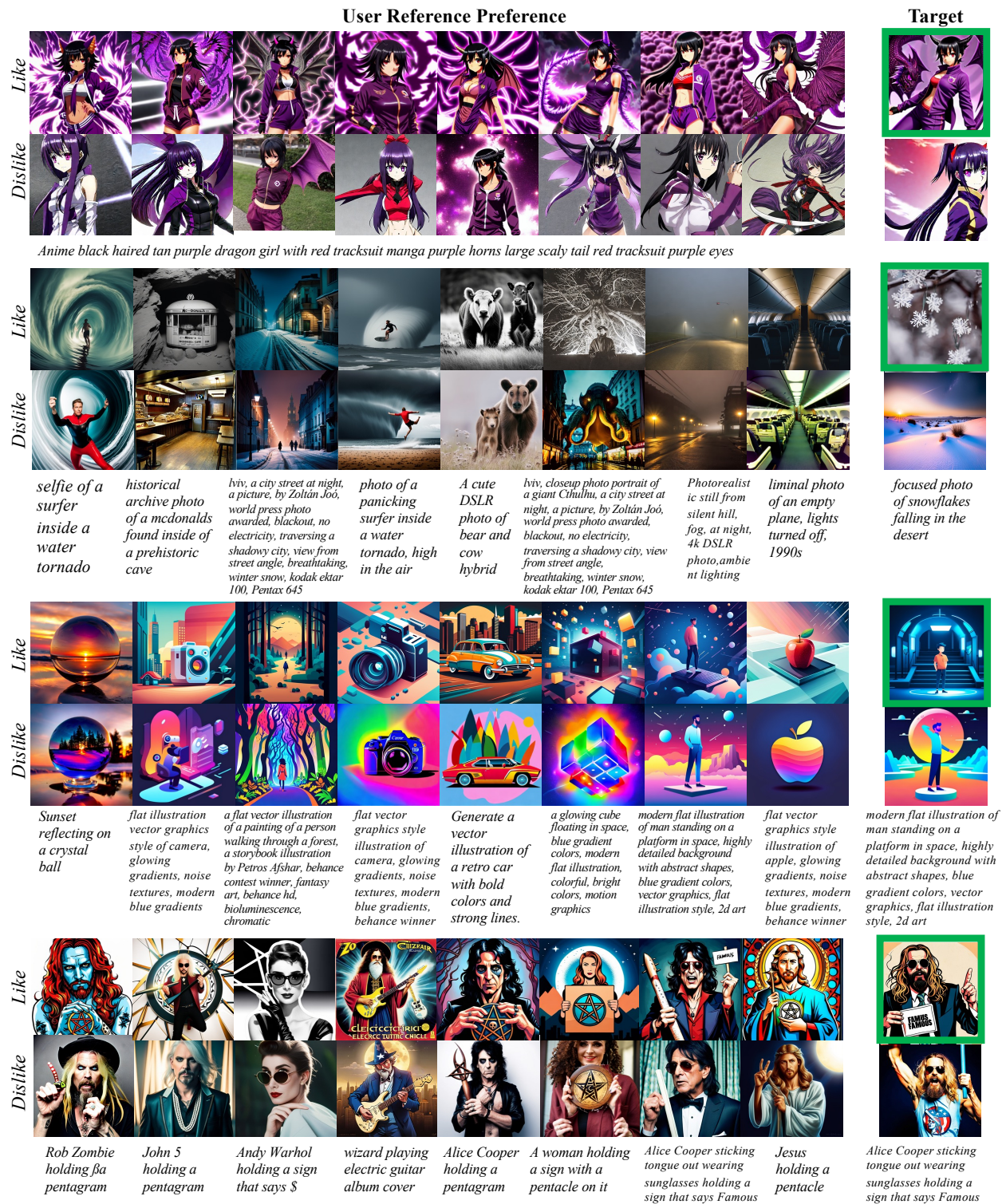
Figure 12. Some examples of the training data.

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

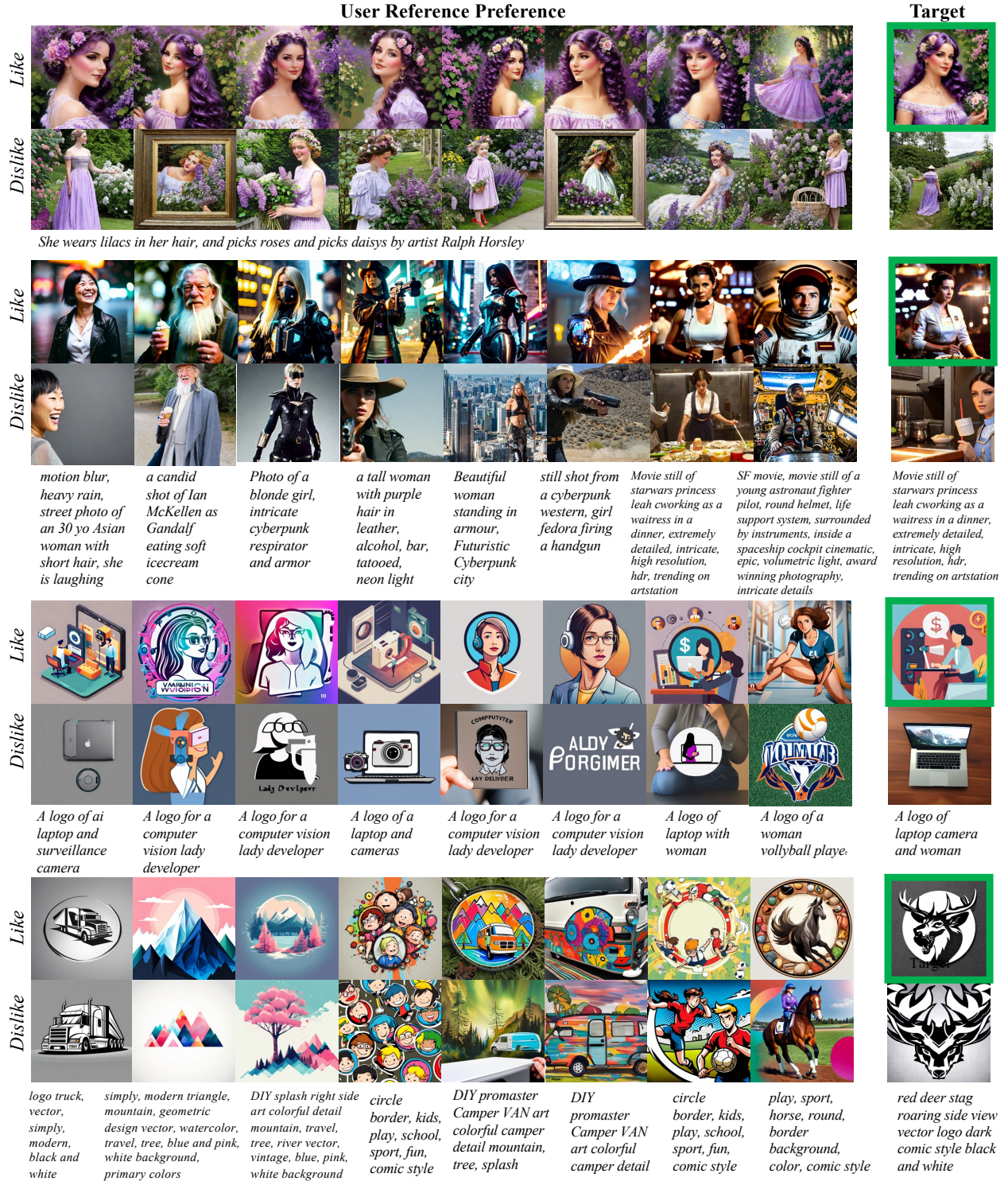*She wears lilacs in her hair, and picks roses and picks daisys by artist Ralph Horsley*

Figure 13. Some examples of the training data.