

DISCRIMINATIVE MULTIMODAL PREFERENCE MODELS AS GUIDANCE FOR PERSONALIZED IMAGE GENERATION

Anonymous authors

Paper under double-blind review

ABSTRACT

User preference prediction requires a deep 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. However, existing methods typically rely on general human preferences or assume static user profiles, often neglecting individual variability and the dynamic, multifaceted nature of personal taste. To address these limitations, we propose an approach built upon Multimodal Large 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 activate group-specific preferences and enhance consistency across similar users. 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 (Ho et al., 2020; Dhariwal & Nichol, 2021; Sohl-Dickstein et al., 2015; Nichol et al., 2022; Saharia et al., 2022; Rombach et al., 2022; Ren et al., 2024; Esser et al., 2024; Sauer et al., 2024a; Mo et al., 2025; Zhou et al., 2025; Zhang et al., 2025; Ba et al., 2025) has significantly advanced the field of 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 guide 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.

The feasibility of such personalized approaches is supported by psychological research, which suggests that aesthetic preference is not arbitrary but often reflects a mixture of low-level visual features (e.g., color, contrast) and high-level semantic content (e.g., subject matter, composition) (Iigaya et al., 2021). Such findings support the assumption that individual taste can be inferred from observable image properties, laying the foundation for data-driven modeling of personalized visual preference.

Building on this foundation, the task of user preference prediction becomes well-defined: given reference data, typically a set of liked and disliked images, the task of user preference prediction is to identify preferences, such as color and content, that align with a user’s tastes. Fig. 1 provides an illustrative example. Existing preference prediction models such as PickScore (Kirstain et al., 2023), ImageReward (Xu et al., 2023), and HPS (Wu et al., 2023b;a) evaluate human preferences at a general level, without granular individual-specific adaptation. Moreover, recent individual-level personalized preference modeling (Salehi et al., 2024; Shen et al., 2024) 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 preference and (2) overlook the significance of users’ disliked images, which provide

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

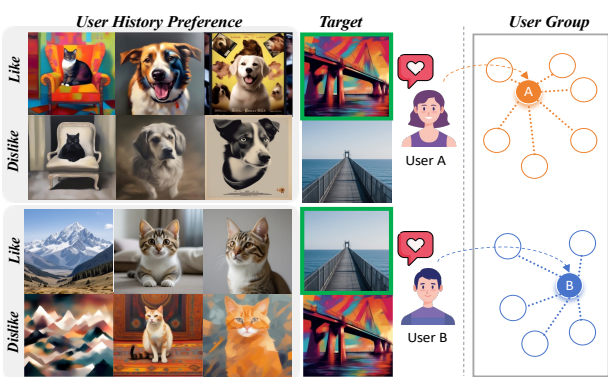


Figure 1: Our task aims to predict target images that align with users’ tastes based on their history data. Users within each group exhibit similar preference distributions and behavioral patterns, while users across different groups may display conflicting or complementary preferences.

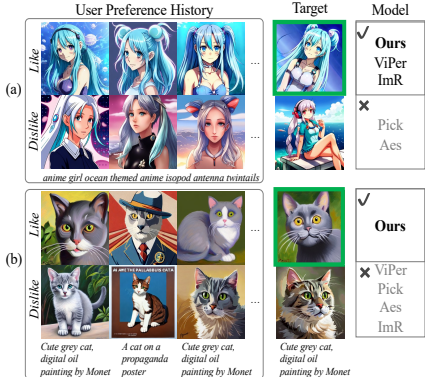


Figure 2: Qualitative comparison of user preference prediction. (a) and (b) illustrate user-specific preferences for style and content, respectively. The green boxes represent the desired outputs that match the user’s preferences.

valuable defeatist 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.

Learning user preferences, while ostensibly analyzing individual historical reference data, fundamentally requires global modeling—capturing both inter-user divergence, which defines individualized content needs, and cross-user commonality, which enables structured learning across similar users. While prior works in text-to-image generation often treat users as independent units, we posit that users frequently exhibit shared preference patterns. This motivates us to formulate a user group structure, where intra-group consistency and inter-group divergence guide the preference modeling process. Our approach is also inspired by recent developments in recommendation systems, where contrastive clustering has been successfully employed to learn group-level behavior representations (Lan et al., 2024). Analogously, we design a multimodal preference learning framework built upon Multimodal Large Language Models (MLLMs) (Laurençon et al., 2023; 2024; Hu et al., 2024; Yao et al., 2024; Li et al., 2024; Liu et al., 2024b;a; An et al., 2025), in which we introduce contrastive preference loss terms to sharpen decision boundaries, and incorporate learnable preference tokens to dynamically encode cluster-specific attributes. This design not only enhances individual preference discrimination but also promotes alignment within user groups, leading to more consistent and structured preference modeling. Our contributions are summarized as follows:

- We introduce a MLLM-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 METHOD

Our approach develops a discriminative preference model that aligns with user-specific tastes. We leverage each user’s preference history $\mathcal{S} = \{(I_{\text{pos}}, I_{\text{neg}}, T)_i\}_{i=1}^{N_{\text{ref}}}$ containing N_{ref} liked/disliked images for prompt T . For any target image pair (z_1, z_2) , we define $D_u(z_1, z_2) = \mathbf{1}[Q(\mathcal{S}_u, z_1) > Q(\mathcal{S}_u, z_2)]$, which equals 1 if user u prefers z_1 over z_2 and 0 otherwise, where $Q(\mathcal{S}, z)$ is a preference scoring function.

Our approach aims to achieve global modeling of user preferences. Therefore, we formalize the user preference structure through the following assumption:

Assumption 1. (User Preference Group Structure). We assume users partition into K groups $\{\mathcal{U}_k\}_{k=1}^K$ satisfying:

Intra-group homogeneity: For users $i, j \in \mathcal{U}_k$:

$$d(\mathcal{S}_i, \mathcal{S}_j) \leq \rho_k, \mathbb{E}[|Q(\mathcal{S}_i, z) - Q(\mathcal{S}_j, z)|] \leq \epsilon_k \quad (1)$$

$$\mathbb{P}[D_i(z_1, z_2) = D_j(z_1, z_2)] \geq 1 - \alpha_k \quad (2)$$

Inter-group heterogeneity: For users $i \in \mathcal{U}_k, j \in \mathcal{U}_{l \neq k}$:

$$d(\mathcal{S}_i, \mathcal{S}_j) \geq \delta_{kl}, \mathbb{E}[|Q(\mathcal{S}_i, z) - Q(\mathcal{S}_j, z)|] \geq \max(\epsilon_k, \epsilon_l) \quad (3)$$

$$\mathbb{P}[D_i(z_1, z_2) \neq D_j(z_1, z_2)] \geq 1 - \beta_{kl} \quad (4)$$

Here, $d(\cdot, \cdot)$ is a distance metric on user preference histories, ϵ_k controls the similarity of preference scores within group; $1 - \alpha_k$ guarantees the consistency of intra-group decisions; $1 - \beta_{kl}$ ensures the divergence of inter-group decisions. This assumption illustrates that: intra-group homogeneity ensures users within the same group have similar preferences and consistent decisions, while inter-group heterogeneity guarantees significant preference differences and decision discrepancies between different groups.

Motivated by Assumption 1, we propose a multimodal large language model-based contrastive preference learning framework. As shown in Fig. 7, our method learns user preferences through contrastive learning on image pairs and employs learnable preference tokens to encode individual aesthetic patterns, enabling personalized preference modeling.

2.1 PREFERENCE LEARNING OBJECTIVE

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_{pos} if the user likes the item and z_{neg} if the user dislikes it. 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.

2.1.1 BASE LOSS.

The base loss, $\mathcal{L}_{\text{base}}$, 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}_{pos} and \mathbf{y}_{neg} , respectively. The base loss is defined as:

$$\mathcal{L}_{\text{base}} = \frac{1}{2} (\mathcal{L}(\mathcal{M}^+(\mathcal{S}, z_{\text{pos}}), \mathbf{y}_{\text{pos}}) + \mathcal{L}(\mathcal{M}^-(\mathcal{S}, z_{\text{neg}}), \mathbf{y}_{\text{neg}})), \quad (5)$$

where $\mathcal{L}(\cdot)$ denotes a classification loss function.

Further Enhancing Preference Extraction. While $\mathcal{L}_{\text{base}}$ capture basic preference patterns, they fail to enforce discriminative separability between liked and disliked items, thus violating the group structure constraints in Assumption 1. Consider the model’s preference prediction function:

$$\mathcal{Q}(\mathcal{S}, z) = \frac{\exp(\mathcal{M}^+(\mathcal{S}, z))}{\exp(\mathcal{M}^+(\mathcal{S}, z)) + \exp(\mathcal{M}^-(\mathcal{S}, z))}. \quad (6)$$

Crucially, $\mathcal{L}_{\text{base}}$ lacks mechanisms to explicitly maximize the logit margin of “like” and “dislike” predictions. When $\mathcal{M}^+(\mathcal{S}, z) \approx \mathcal{M}^-(\mathcal{S}, z)$, we have $\mathcal{Q}(\mathcal{S}, z) \approx 0.5$, leading to ambiguous decision boundaries. Such predictions fundamentally violate Assumption 1 in two ways: (1) users within the same group \mathcal{U}_k have similar preference scores $Q(\mathcal{S}_i, z_{\text{pos}}) \approx Q(\mathcal{S}_j, z_{\text{pos}})$, $Q(\mathcal{S}_i, z_{\text{neg}}) \approx Q(\mathcal{S}_j, z_{\text{neg}})$ may make inconsistent pairwise decisions, violating the intra-group decision consistency constraint; and (2) when users from different groups \mathcal{U}_k and \mathcal{U}_l both exhibit ambiguous predictions near 0.5, violating the inter-group score divergence constraint. Simultaneously, their decisions may become unexpectedly similar, thereby violating the inter-group decision divergence requirement.

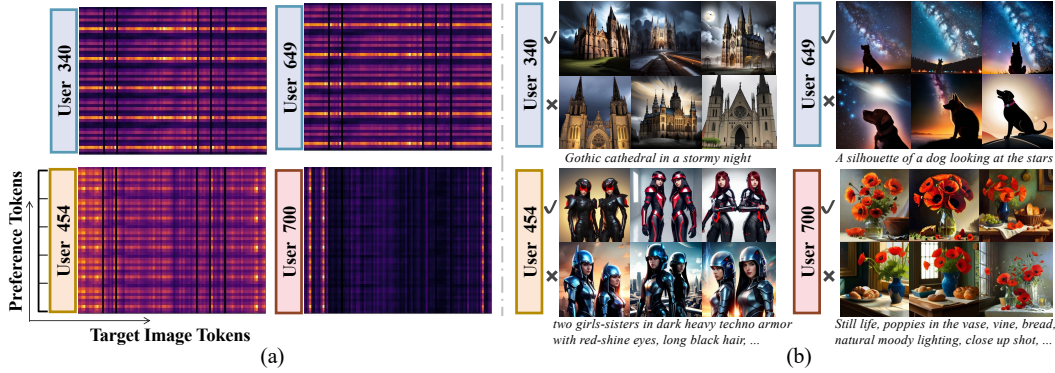


Figure 3: (a) Attention scores \mathcal{A} represent interactions between preference tokens and target image tokens for individual users. Each user has a unique reference history, and we concatenate the same target image to the input sequence across users. For each user, the horizontal axis represents tokens from the target image, while the vertical axis represents the preference tokens. Each user has five random re-orderings of reference images. (b) Examples of images liked (✓) or disliked (✗).

To further enhance preference extraction, we propose a contrastive preference learning approach that enforces $\mathcal{M}^+(\mathcal{S}, z_{\text{pos}}) \gg \mathcal{M}^+(\mathcal{S}, z_{\text{neg}})$ and $\mathcal{M}^-(\mathcal{S}, z_{\text{neg}}) \gg \mathcal{M}^-(\mathcal{S}, z_{\text{pos}})$. This contrastive mechanism pushes preference scores away from the ambiguous boundary, creating decisive preference predictions with higher confidence. A detailed mathematical analysis is provided in Appendix C.

2.1.2 CONTRASTIVE PREFERENCE LOSS

We introduce two contrastive preference loss terms, \mathcal{L}_+ and \mathcal{L}_- , which enhance the model’s ability to differentiate between “like” and “dislike” predictions by emphasizing their relative rankings.

The positive preference loss \mathcal{L}_+ ensures the model’s positive logits favor positive samples over negative samples. Conversely, the negative preference loss \mathcal{L}_- ensures the model’s negative logits favor negative samples over positive samples. Together, these losses push predictions away from ambiguous boundaries:

$$\begin{aligned}\mathcal{L}_+ &= -\frac{1}{N} \sum_{i=1}^N \log \sigma(\mathcal{M}^+(\mathcal{S}, z_{\text{pos}}) - \mathcal{M}^+(\mathcal{S}, z_{\text{neg}})) \\ \mathcal{L}_- &= -\frac{1}{N} \sum_{i=1}^N \log \sigma(\mathcal{M}^-(\mathcal{S}, z_{\text{neg}}) - \mathcal{M}^-(\mathcal{S}, z_{\text{pos}}))\end{aligned}\quad (7)$$

where N is the number of samples and σ is the sigmoid function. The total contrastive preference loss is the sum of these components, $\mathcal{L}_{\text{CP}} = \mathcal{L}_+ + \mathcal{L}_-$.

The final loss function combines the base loss with the contrastive preference loss to enhance the model’s ability to distinguish user preferences: $\mathcal{L}_{\text{all}} = \mathcal{L}_{\text{base}} + \mathcal{L}_{\text{CP}}$. This helps the model optimize for nuanced preference distinctions, leading to more accurate and effective predictions.

2.2 LEARNABLE PREFERENCE TOKENS

Based on Assumption 1, we need to adaptively identify user groups and activate corresponding preference patterns without group labels, ensuring intra-group homogeneity and inter-group heterogeneity. Our key insight is leveraging the inherent soft-clustering properties of attention mechanism for personalized group discovery and preference activation.

Consider the core computation of the attention mechanism: $A(\mathcal{Q}, \mathcal{K}) = \text{softmax}\left(\frac{\mathcal{Q}\mathcal{K}^T}{\sqrt{d_k}}\right)$. For any two users $i, j \in \mathcal{U}_k$ within the same group, satisfying the preference consistency constraint $d(\mathcal{S}_i, \mathcal{S}_j) \leq \rho_k$, there exists a Lipschitz continuous mapping $\phi: \mathcal{S} \rightarrow \mathcal{Q}$ such that: $|\phi(\mathcal{S}_i) - \phi(\mathcal{S}_j)| \leq L_\phi \cdot \rho_k$, where L_ϕ is the Lipschitz constant of the mapping ϕ . Leveraging the Softmax KL Divergence Bound lemma, the similarity of attention responses for users within the same group is constrained

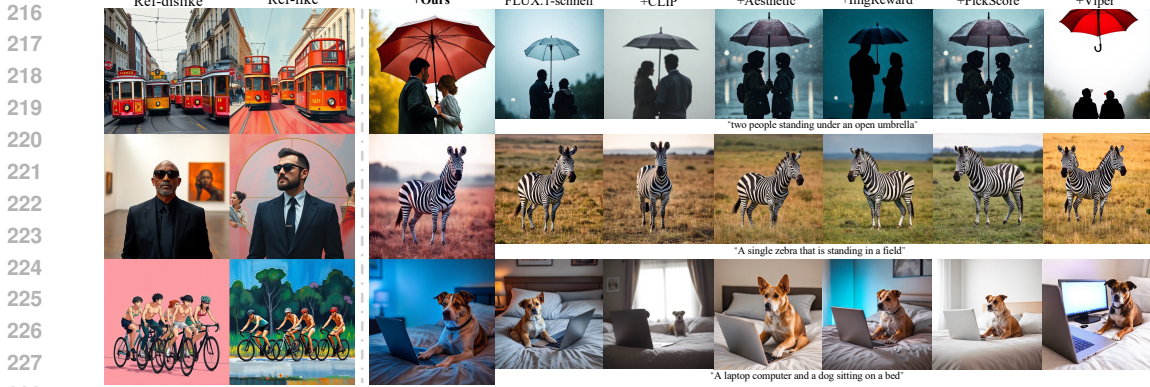


Figure 4: Qualitative comparison of text-to-image generation for three users. Each row shows user preferences (Ref-dislike/like) and generation results from our personalized preference model vs. image-text alignment (CLIP Score), aesthetic quality (Aesthetic Score), general human preference (ImageReward, PickScore), and personalized preference (ViPer) models.

below:

$$\mathbb{E}_{\mathcal{K}} [\text{KL}(A(Q_i, \mathcal{K}) \| A(Q_j, \mathcal{K}))] \leq f_{\text{intra}}^{(k)}(\rho_k) \quad (8)$$

where $f_{\text{intra}}^{(k)}$ is a group-specific continuous increasing function. For users from different groups, $i \in \mathcal{U}_k$ and $j \in \mathcal{U}_{l \neq k}$, satisfying the preference separability constraint $d(\mathcal{S}_i, \mathcal{S}_j) \geq \delta_{kl}$, their attention responses are constrained by:

$$\mathbb{E}_{\mathcal{K}} [\text{KL}(A(Q_i, \mathcal{K}) \| A(Q_j, \mathcal{K}))] \geq g_{\text{inter}}^{(kl)}(\delta_{kl}) \quad (9)$$

where $g_{\text{inter}}^{(kl)}$ is an inter-group continuous increasing function.

Following this theoretical intuition, it can be concluded that attention mechanism can adaptively cluster users with different preferences. However, in preference prediction tasks, the context varies for different users, making it difficult for MLLM to formulate groups by adjusting the similarity of attention. Therefore, we introduce additional, shared learnable preference tokens $P_v \in \mathbb{R}^{L_p \times D}$ to provide an extra attention term, where L_p is the number of preference tokens and D is the embedding dimension. This allows group discovery to be achieved by adjusting the attention towards these preference tokens. Given a user’s historical sequence \mathcal{S} and a target item z , we encode all input content (excluding the target image label token) into a user-specific token sequence $x_u \in \mathbb{R}^{L_e \times D}$. These preference tokens are then concatenated with the user sequence to form the complete input = $[P_v; x_u]$. The Transformer uses attention where the user-specific sequence x_u serves as the Query, and the preference tokens P_v serve as both the Key and the Value, enabling selective activation of relevant preference patterns.

Mining Similar Users via Attention Mechanism. 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 User 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, revealing that the preference token space does not simply cluster all users together but rather preserves individual differences while leveraging commonalities where applicable. Further details and analysis can be found in Appendix D.

| Model | Aes Score | CLIP Score | ImageReward | HPS Score | PickScore* | IDEFICS | ViPer | Ours |
|--------------|-----------|------------|-------------|-----------|------------|---------|-------|--------------|
| N_{ref} | 0 | 0 | 0 | 0 | 0 | 8 | 8 | 8 |
| accuracy (%) | 49.96 | 53.13 | 55.64 | 56.85 | 57.72 | 50.27 | 55.15 | 61.68 |

* Trained with the same dataset as our model.

Table 1: Preference classification accuracy on pairwise comparisons.

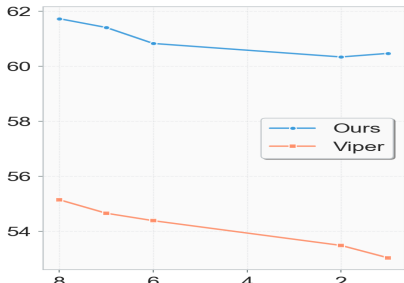


Figure 5: Prediction accuracy with different numbers of images (N_{ref}).

Table 2: Accuracy in one-positive-three-negative evaluation setting. We report the top-1 to top-3 accuracy (%).

| Model | N_{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 |
| Ours (w/o P_v) | 8 | 35.72 | 61.64 | 83.44 |
| Ours | 8 | 37.47 | 62.85 | 84.74 |

3 EXPERIMENTS

3.1 USER-SPECIFIC PREFERENCE PREDICTION

Datasets. We process Pick-a-Pic v2 dataset (Kirstain et al., 2023), which is collected through real user interactions, to obtain user-specific preference datasets based on user IDs. 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.

Implementation Details. Following the approach of (Salehi et al., 2024), we use IDEFICS2-8B (Laurençon et al., 2024) as our MLLM. We employ a batch size of 64, training on 8 A100 (80GB) GPUs with a local batch size of 2 pairs and gradient accumulation over 4 steps.

Evaluation Metric. We evaluate user-specific preference prediction using top- K accuracy, which measures whether the liked image ranks among the top K candidates out of multiple disliked ones.

Comparison to Other Methods. In our study, we compare our method with several existing approaches: (1) CLIP (Radford et al., 2021), designed to evaluate generic text-image alignment, (2) LAION Aesthetic Score Predictor (Schuhmann et al., 2021), which evaluates aesthetic quality, (3) ImageReward (Xu et al., 2023), (4) HPS (Wu et al., 2023a) and (5) PickScore (Kirstain et al., 2023), which focus on learning general human preferences and consider relative preferences between images, and (6) ViPer proxy model (Salehi et al., 2024), which is trained on a preference dataset constructed from 5,000 simulated agents representing diverse individual preferences.

Qualitative User-Specific Preference Prediction. In Fig. 2, we compare our model to ViPer, PickScore, ImageReward, and Aesthetic Score. Our model effectively aligns with user-specific preferences by distinguishing styles and content according to user reference data. For instance, in Fig. 2 (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. 2 (b), Our method alleviates semantic ambiguity, particularly when handling terms like “grey cat” that encompass multiple visual appearances under a single designation, ensuring that the generated images better reflect the user’s intended preferences. More results are in Appendix B.

Quantitative User-Specific Preference Prediction. Tab. 2 and Tab. 1 show that our model achieves the highest accuracy, outperforming baselines like ViPer, CLIP, and ImageReward. It performs especially well in settings with multiple disliked images. Generic metrics perform poorly, reflecting the gap between general and personalized preferences. Unlike ViPer, our model captures fine-grained user-specific patterns through contrastive learning and preference tokens, resulting in more accurate and robust predictions.

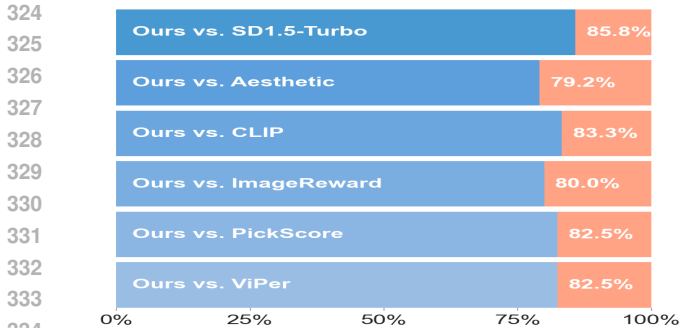


Figure 6: Human expert evaluation of generated images from different methods on SD1.5-Turbo.

Table 3: Evaluation comparison on preference prediction.

| Method | Top-1 Acc (%) |
|-------------------|---------------|
| Claude-3.5-Sonnet | 47.96 |
| Human Expert | 57.60 |
| Ours | 60.45 |

Table 4: Ablation study for preference tokens numbers.

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

Table 5: Quantitative Results. **Bold** and underlined values represent optimal and second-best performance respectively.

| Model | Aesthetic(↑) | CLIP Score(↑) | ImageReward(↑) | PickScore(↑) | HPS Score(↑) | CSD Score(↑) | ViPer (↑) |
|-------------|--------------|---------------|----------------|--------------|--------------|--------------|-------------|
| SD1.5-Turbo | <u>5.81</u> | 35.44 | 0.69 | <u>21.92</u> | <u>28.04</u> | 0.32 | 0.21 |
| + ViPer | 5.37 | 34.08 | 0.31 | <u>21.55</u> | <u>27.91</u> | 0.38 | <u>0.62</u> |
| + Ours | 5.99 | <u>34.45</u> | <u>0.60</u> | 22.28 | 28.49 | 0.42 | 0.84 |

Number of User Reference Preferences. As demonstrated in Fig. 5, our method consistently maintains the highest pairwise classification accuracy as the preference sequence length varies. This indicates that our model can effectively preserve accuracy with limited reference data. In contrast, ViPer shows a decline in accuracy as sequence length shortens, highlighting the stability and adaptability of our approach in scenarios with limited user reference.

Evaluation with Multimodal LLM and Human Experts. To evaluate real-world performance, we conducted a user study on 200 randomly sampled test cases. Claude-3.5-Sonnet (a state-of-the-art multimodal language model) and ten human experts were asked to infer preferences from reference images and select the preferred image from each pair. As shown in Tab. 3, our model outperforms both Claude-3.5-Sonnet (47.96%) and human experts (57.60%) with a top-1 accuracy of 60.45%. This indicates that our method not only captures clear aesthetic signals but also models subtle user preferences more effectively than humans.

3.2 PERSONALIZING GENERATION WITH USER PREFERENCES

Datasets. To evaluate our model’s ability to learn detailed attribute preferences and guide image generation accordingly, we constructed a diverse dataset. To simulate realistic user preferences, we configured 30 Claude agents (Anthropic, 2024) to represent diverse human preferences across 7 key dimensions: Art styles, Color palette, Composition, Skill level, Detail level, and other aesthetic attributes. To ensure diverse preferences, each agent was configured with unique sets of preferred and dispreferred attributes, maintaining at least 80% Jaccard distance between any pair of agents’ preference profiles. For image generation, each agent utilized FLUX.1-schnell (Black Forest Labs, 2024) to generate 10 images aligned with their preferences and 10 images representing their dislikes, resulting in a dataset of 600 images total. Examples of generated images and dataset details are provided in Appendix D.

Experimental Setup. Following the ReNO approach (Eyring et al., 2024), which enhances image generation quality by optimizing initial noise during inference using reward model guidance, we generate images guided by our preference model while incorporating both positive and negative user feedback. We first generate initial target images, then utilize the agent preference data as input to our model, which provides reward signals for iteratively optimizing these target images. Specifically, we apply Eq. (6) to obtain the guidance signal, which is subsequently used in an iterative refinement process to enhance target image adherence based on the learned user preferences. We conduct experiments using two generative models: FLUX.1-schnell (Black Forest Labs, 2024) and Stable Diffusion 1.5 Turbo (Sauer et al., 2024b). The images are generated using the same random seeds. Additional experimental details are provided in Appendix D.

Table 6: Clustering evaluation metrics for the ablation study.

| Method | Silhouette Score (\uparrow) | Davies-Bouldin Score (\downarrow) | Top-1 Acc (%) |
|-------------------------|---------------------------------|---------------------------------------|---------------|
| Base Loss | 0.596 | 0.812 | 60.47 |
| + L_{cp} | 0.635 | 0.812 | 61.37 |
| + $L_{cp} + P_v$ (Ours) | 0.646 | 0.806 | 61.68 |

Evaluation Metric. We employ 7 comprehensive evaluation metrics to assess our model’s performance. In addition to 5 standard general-purpose (Aesthetic, CLIP) and human preference metrics (PickScore, ImageReward, HPS), we introduce two specialized metrics: (1) We assess style-following performance using CSD metric (Somepalli et al., 2024) to evaluate how well the generated images maintain consistency with the specified artistic styles and attributes. (2) We utilize the ViPer proxy metric, which predicts user preferences by analyzing reference images.

Enhancing Image Generation with User Preferences. As shown in Tab. 5, our method effectively enhances Stable Diffusion 1.5 Turbo (base model) across most metrics by learning richer attribute preferences from user feedback, particularly excelling in personalized preference metrics, and significantly outperforms ViPer across all metrics. This demonstrates that our model can effectively utilize user preferences to refine and guide image generation. Fig. 4 shows generation examples for three different preferences on FLUX.1-schnell with different reward models, where our model successfully extracts personalized preferences and ensures superior image generation results.

User Study. To assess the effectiveness of different reward models in guiding personalized image generation, we conducted a user study involving ten human experts. Experts were asked to compare pairs of images generated using different reward models and select the one that better aligns with the reference preferences. As shown in Fig. 6, our method consistently outperforms all baselines, achieving win rates above 79% across all comparisons. This highlights the superiority of our approach in capturing fine-grained user preferences for generation tasks.

3.3 ANALYSIS AND ABLATION STUDY

We perform ablation studies and conduct thorough analysis on the processed Pick-a-Pic v2 dataset to investigate the impact of the contrastive preference loss, learnable preference tokens, and the number of preference tokens. To assess whether our model captures structured user preference patterns, we sample 70 users from the test set, each with 16 historical preference images. All users are paired with the same target image, and we extract the final token embedding from the last layer of the MLLM. These embeddings are then clustered using K-means (MacQueen, 1967), and clustering quality is evaluated using Silhouette Score (Rousseeuw, 1987) and Davies-Bouldin Score (Davies & Bouldin, 1979). As shown in Tab. 6, incorporating L_{cp} and P_v significantly improves cluster compactness and separation. This indicates that our method successfully discovers user groupings, aligning with the group structure assumption in our framework. Besides, the full model achieves a top-1 accuracy of 61.68%, representing a cumulative performance gain over the baseline model. Additionally, we perform ablation experiments to analyze the impact of preference token length. As shown in Tab. 4, the results demonstrate that using 10 preference tokens achieves the highest Top-1 accuracy, slightly outperforming configurations with 5 and 20 preference tokens, while the overall differences remain small. This indicates that the selection of preference token quantity exhibits good robustness, effectively enhancing personalized modeling of user preferences within a reasonable range.

4 CONCLUSION

In this paper, we propose a novel approach for user-specific preference prediction in generated images by leveraging Multimodal Large Language Models (MLLMs). To address the limitations of existing methods that focus primarily on general human preferences or superficial attributes, 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, enabling both personalization and generalization. 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

- 432
433
434 Ruichuan An, Sihan Yang, Renrui Zhang, Zijun Shen, Ming Lu, Gaole Dai, Hao Liang, Ziyu Guo,
435 Shilin Yan, Yulin Luo, Bocheng Zou, Chaoqun Yang, and Wentao Zhang. Unictokens: Boosting
436 personalized understanding and generation via unified concept tokens. *CoRR*, abs/2505.14671,
437 2025. doi: 10.48550/ARXIV.2505.14671. URL [https://doi.org/10.48550/arXiv.
438 2505.14671](https://doi.org/10.48550/arXiv.2505.14671).
- 439 Anthropic. Claude, 2024. URL <https://claude.ai>. Large language model.
- 440 Ying Ba, Tianyu Zhang, Yalong Bai, Wenyi Mo, Tao Liang, Bing Su, and Ji-Rong Wen. Enhancing
441 reward models for high-quality image generation: Beyond text-image alignment, 2025. URL
442 <https://arxiv.org/abs/2507.19002>.
- 443 Black Forest Labs. FLUX.1-schnell: High-performance image generation model, 2024. URL
444 <https://huggingface.co/black-forest-labs/FLUX.1-schnell>.
- 445
446 Chaofeng Chen, Annan Wang, Haoning Wu, Liang Liao, Wenxiu Sun, Qiong Yan, and Weisi Lin.
447 Enhancing diffusion models with text-encoder reinforcement learning. In *Proceedings of the
448 European Conference on Computer Vision (ECCV)*, 2024.
- 449 David L. Davies and Donald W. Bouldin. A cluster separation measure. *IEEE Trans. Pattern
450 Anal. Mach. Intell.*, 1(2):224–227, 1979. doi: 10.1109/TPAMI.1979.4766909. URL [https:
451 //doi.org/10.1109/TPAMI.1979.4766909](https://doi.org/10.1109/TPAMI.1979.4766909).
- 452 Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of
453 quantized llms. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and
454 Sergey Levine (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference
455 on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December
456 10 - 16, 2023*, 2023. URL [http://papers.nips.cc/paper_files/paper/2023/
457 hash/1feb87871436031bdc0f2beaa62a049b-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/1feb87871436031bdc0f2beaa62a049b-Abstract-Conference.html).
- 458
459 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *NeurIPS*,
460 2021.
- 461 Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam
462 Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for
463 high-resolution image synthesis. *arXiv preprint arXiv:2403.03206*, 2024.
- 464 Luca Eyring, Shyamgopal Karthik, Karsten Roth, Alexey Dosovitskiy, and Zeynep Akata. Reno:
465 Enhancing one-step text-to-image models through reward-based noise optimization. *Neural
466 Information Processing Systems (NeurIPS)*, 2024.
- 467
468 Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit Haim Bermano, Gal Chechik, and
469 Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using
470 textual inversion. In *The Eleventh International Conference on Learning Representations, ICLR
471 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL [https://openreview.
472 net/forum?id=NAQvF08TcyG](https://openreview.net/forum?id=NAQvF08TcyG).
- 473 Yaru Hao, Zewen Chi, Li Dong, and Furu Wei. Optimizing prompts for text-to-image generation.
474 In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine
475 (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural
476 Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 -
477 16, 2023*, 2023. URL [http://papers.nips.cc/paper_files/paper/2023/hash/
478 d346d91999074dd8d6073d4c3b13733b-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/d346d91999074dd8d6073d4c3b13733b-Abstract-Conference.html).
- 479 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*,
480 2020.
- 481 Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang,
482 Yuxiang Huang, Weilin Zhao, Xinrong Zhang, Zhen Leng Thai, Kai Zhang, Chongyi Wang, Yuan
483 Yao, Chenyang Zhao, Jie Zhou, Jie Cai, Zhongwu Zhai, Ning Ding, Chao Jia, Guoyang Zeng, Dahai
484 Li, Zhiyuan Liu, and Maosong Sun. Minicpm: Unveiling the potential of small language models
485 with scalable training strategies. *CoRR*, abs/2404.06395, 2024. doi: 10.48550/ARXIV.2404.06395.
URL <https://doi.org/10.48550/arXiv.2404.06395>.

- 486 Kiyohito Iigaya, Sanghyun Yi, Iman A. Wahle, Koranis Tanwisuth, and John P. O’Doherty. Aesthetic
487 preference for art can be predicted from a mixture of low- and high-level visual features. *Nature*
488 *Human Behaviour*, 5:743 – 755, 2021. URL [https://api.semanticscholar.org/
489 CorpusID:235072663](https://api.semanticscholar.org/CorpusID:235072663).
- 490 Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy.
491 Pick-a-pic: An open dataset of user preferences for text-to-image generation. In Alice Oh,
492 Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.),
493 *Advances in Neural Information Processing Systems 36: Annual Conference on Neural In-*
494 *formation Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 -*
495 *16, 2023*, 2023. URL [http://papers.nips.cc/paper_files/paper/2023/hash/
496 73aacd8b3b05b4b503d58310b523553c-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/73aacd8b3b05b4b503d58310b523553c-Abstract-Conference.html).
- 497
498 Wei Lan, Guoxian Zhou, Qingfeng Chen, Wenguang Wang, Shirui Pan, Yi Pan, and Shichao
499 Zhang. Contrastive clustering learning for multi-behavior recommendation. *ACM Transactions*
500 *on Information Systems*, 43:1 – 23, 2024. URL [https://api.semanticscholar.org/
501 CorpusID:273028917](https://api.semanticscholar.org/CorpusID:273028917).
- 502 Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov,
503 Thomas Wang, Siddharth Karamcheti, Alexander M. Rush, Douwe Kiela, Matthieu Cord, and
504 Victor Sanh. OBELICS: an open web-scale filtered dataset of interleaved image-text documents.
505 In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey
506 Levine (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference*
507 *on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA,*
508 *December 10 - 16, 2023*, 2023. URL [http://papers.nips.cc/paper_files/paper/
509 2023/hash/e2cfb719f58585f779d0a4f9f07bd618-Abstract-Datasets_
510 and_Benchmarks.html](http://papers.nips.cc/paper_files/paper/2023/hash/e2cfb719f58585f779d0a4f9f07bd618-Abstract-Datasets_ and_Benchmarks.html).
- 511 Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. What matters when building
512 vision-language models? *arXiv preprint arXiv:2405.02246*, 2024.
- 513
514 Kimin Lee, Hao Liu, Moonkyung Ryu, Olivia Watkins, Yuqing Du, Craig Boutilier, Pieter Abbeel,
515 Mohammad Ghavamzadeh, and Shixiang Shane Gu. Aligning text-to-image models using human
516 feedback. *CoRR*, abs/2302.12192, 2023. doi: 10.48550/ARXIV.2302.12192. URL [https:
517 //doi.org/10.48550/arXiv.2302.12192](https://doi.org/10.48550/arXiv.2302.12192).
- 518 Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li.
519 Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. *CoRR*,
520 abs/2407.07895, 2024. doi: 10.48550/ARXIV.2407.07895. URL [https://doi.org/10.
521 48550/arXiv.2407.07895](https://doi.org/10.48550/arXiv.2407.07895).
- 522 Youwei Liang, Junfeng He, Gang Li, Peizhao Li, Arseniy Klimovskiy, Nicholas Carolan, Jiao Sun,
523 Jordi Pont-Tuset, Sarah Young, Feng Yang, Junjie Ke, Krishnamurthy Dj Dvijotham, Katherine M.
524 Collins, Yiwen Luo, Yang Li, Kai J. Kohlhoff, Deepak Ramachandran, and Vidhya Navalpakkam.
525 Rich human feedback for text-to-image generation. In *IEEE/CVF Conference on Computer Vision*
526 *and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*, pp. 19401–19411.
527 IEEE, 2024a. doi: 10.1109/CVPR52733.2024.01835. URL [https://doi.org/10.1109/
528 CVPR52733.2024.01835](https://doi.org/10.1109/CVPR52733.2024.01835).
- 529 Zhanhao Liang, Yuhui Yuan, Shuyang Gu, Bohan Chen, Tiankai Hang, Ji Li, and Liang Zheng.
530 Step-aware preference optimization: Aligning preference with denoising performance at each step.
531 *CoRR*, abs/2406.04314, 2024b. doi: 10.48550/ARXIV.2406.04314. URL [https://doi.org/
532 10.48550/arXiv.2406.04314](https://doi.org/10.48550/arXiv.2406.04314).
- 533
534 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction
535 tuning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024,*
536 *Seattle, WA, USA, June 16-22, 2024*, pp. 26286–26296. IEEE, 2024a. doi: 10.1109/CVPR52733.
537 2024.02484. URL <https://doi.org/10.1109/CVPR52733.2024.02484>.
- 538 Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee.
539 Llava-next: Improved reasoning, ocr, and world knowledge. January 2024b. URL [https:
//llava-vl.github.io/blog/2024-01-30-llava-next/](https://llava-vl.github.io/blog/2024-01-30-llava-next/).

- 540 J. MacQueen. Some methods for classification and analysis of multivariate observations. 1967. URL
541 <https://api.semanticscholar.org/CorpusID:6278891>.
542
- 543 Wenyi Mo, Tianyu Zhang, Yalong Bai, Bing Su, Ji-Rong Wen, and Qing Yang. Dynamic prompt
544 optimizing for text-to-image generation. In *IEEE/CVF Conference on Computer Vision and Pattern
545 Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*, pp. 26617–26626. IEEE, 2024.
546 doi: 10.1109/CVPR52733.2024.02514. URL [https://doi.org/10.1109/CVPR52733.
547 2024.02514](https://doi.org/10.1109/CVPR52733.2024.02514).
- 548 Wenyi Mo, Tianyu Zhang, Yalong Bai, Bing Su, and Ji-Rong Wen. Uniform attention maps: Boosting
549 image fidelity in reconstruction and editing. In *IEEE/CVF Winter Conference on Applications of
550 Computer Vision, WACV 2025, Tucson, AZ, USA, February 26 - March 6, 2025*, pp. 4420–4429.
551 IEEE, 2025. doi: 10.1109/WACV61041.2025.00434. URL [https://doi.org/10.1109/
552 WACV61041.2025.00434](https://doi.org/10.1109/WACV61041.2025.00434).
- 553
- 554 Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew,
555 Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with
556 text-guided diffusion models. In *ICML, 2022*.
- 557
- 558 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
559 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever.
560 Learning transferable visual models from natural language supervision. In Marina Meila and
561 Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning,
562 ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning
563 Research*, pp. 8748–8763. PMLR, 2021. URL [http://proceedings.mlr.press/v139/
564 radford21a.html](http://proceedings.mlr.press/v139/radford21a.html).
- 565
- 566 Yuxi Ren, Xin Xia, Yanzuo Lu, Jiacheng Zhang, Jie Wu, Pan Xie, Xing Wang, and Xuefeng Xiao.
567 Hyper-sd: Trajectory segmented consistency model for efficient image synthesis. *arXiv preprint
568 arXiv:2404.13686*, 2024.
- 569
- 570 Robin Rombach, A. Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution
571 image synthesis with latent diffusion models. *2022 IEEE/CVF Conference on Computer Vision
572 and Pattern Recognition (CVPR)*, pp. 10674–10685, 2022.
- 573
- 574 Peter J. Rousseeuw. Silhouettes: a graphical aid to the interpretation and validation of cluster
575 analysis. *Journal of Computational and Applied Mathematics*, 20:53–65, 1987. URL [https:
576 //api.semanticscholar.org/CorpusID:189900](https://api.semanticscholar.org/CorpusID:189900).
- 577
- 578 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aber-
579 man. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In
580 *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC,
581 Canada, June 17-24, 2023*, pp. 22500–22510. IEEE, 2023. doi: 10.1109/CVPR52729.2023.02155.
582 URL <https://doi.org/10.1109/CVPR52729.2023.02155>.
- 583
- 584 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L.
585 Denton, Seyed Kamyar Seyed Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol
586 Ayan, Tim Salimans, Jonathan Ho, David J. Fleet, and Mohammad Norouzi. Pho-
587 torealistic text-to-image diffusion models with deep language understanding. In
588 *NeurIPS, 2022*. URL [http://papers.nips.cc/paper_files/paper/2022/hash/
589 ec795aeadae0b7d230fa35cbaf04c041-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/ec795aeadae0b7d230fa35cbaf04c041-Abstract-Conference.html).
- 590
- 591 Sogand Salehi, Mahdi Shafiei, Teresa Yeo, Roman Bachmann, and Amir Zamir. Viper: Visual
592 personalization of generative models via individual preference learning. *CoRR*, abs/2407.17365,
593 2024. doi: 10.48550/ARXIV.2407.17365. URL [https://doi.org/10.48550/arXiv.
594 2407.17365](https://doi.org/10.48550/arXiv.2407.17365).
- 595
- 596 Axel Sauer, Frederic Boesel, Tim Dockhorn, Andreas Blattmann, Patrick Esser, and Robin Rombach.
597 Fast high-resolution image synthesis with latent adversarial diffusion distillation. *arXiv preprint
598 arXiv:2403.12015*, 2024a.

- 594 Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin Rombach. Adversarial diffusion
595 distillation. In Ales Leonardis, Elisa Ricci, Stefan Roth, Olga Russakovsky, Torsten Sattler,
596 and Gül Varol (eds.), *Computer Vision - ECCV 2024 - 18th European Conference, Milan, Italy,
597 September 29-October 4, 2024, Proceedings, Part LXXXVI*, volume 15144 of *Lecture Notes in
598 Computer Science*, pp. 87–103. Springer, 2024b. doi: 10.1007/978-3-031-73016-0_6. URL
599 https://doi.org/10.1007/978-3-031-73016-0_6.
- 600 Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis,
601 Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of
602 clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021.
- 603 Xiaoteng Shen, Rui Zhang, Xiaoyan Zhao, Jieming Zhu, and Xi Xiao. PMG : Personalized multimodal
604 generation with large language models. In Tat-Seng Chua, Chong-Wah Ngo, Ravi Kumar, Hady W.
605 Lauw, and Roy Ka-Wei Lee (eds.), *Proceedings of the ACM on Web Conference 2024, WWW 2024,
606 Singapore, May 13-17, 2024*, pp. 3833–3843. ACM, 2024. doi: 10.1145/3589334.3645633. URL
607 <https://doi.org/10.1145/3589334.3645633>.
- 608 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised
609 learning using nonequilibrium thermodynamics. In *ICML*, 2015.
- 610 Gowthami Somepalli, Anubhav Gupta, Kamal Gupta, Shramay Palta, Micah Goldblum, Jonas
611 Geiping, Abhinav Shrivastava, and Tom Goldstein. Measuring style similarity in diffusion models.
612 *CoRR*, abs/2404.01292, 2024. doi: 10.48550/ARXIV.2404.01292. URL [https://doi.org/
613 10.48550/arXiv.2404.01292](https://doi.org/10.48550/arXiv.2404.01292).
- 614 Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou, Aaron Lou, Senthil Purushwalkam,
615 Stefano Ermon, Caiming Xiong, Shafiq Joty, and Nikhil Naik. Diffusion model alignment using
616 direct preference optimization. In *IEEE/CVF Conference on Computer Vision and Pattern
617 Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*, pp. 8228–8238. IEEE, 2024.
618 doi: 10.1109/CVPR52733.2024.00786. URL [https://doi.org/10.1109/CVPR52733.
619 2024.00786](https://doi.org/10.1109/CVPR52733.2024.00786).
- 620 Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng
621 Li. Human preference score v2: A solid benchmark for evaluating human preferences of text-
622 to-image synthesis. *CoRR*, abs/2306.09341, 2023a. doi: 10.48550/ARXIV.2306.09341. URL
623 <https://doi.org/10.48550/arXiv.2306.09341>.
- 624 Xiaoshi Wu, Keqiang Sun, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score: Better
625 aligning text-to-image models with human preference. In *IEEE/CVF International Conference on
626 Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023*, pp. 2096–2105. IEEE, 2023b.
627 doi: 10.1109/ICCV51070.2023.00200. URL [https://doi.org/10.1109/ICCV51070.
628 2023.00200](https://doi.org/10.1109/ICCV51070.2023.00200).
- 629 Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao
630 Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation.
631 In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine
632 (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural
633 Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 -
634 16, 2023, 2023*. URL [http://papers.nips.cc/paper_files/paper/2023/hash/
635 33646ef0ed554145eab65f6250fab0c9-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/33646ef0ed554145eab65f6250fab0c9-Abstract-Conference.html).
- 636 Kai Yang, Jian Tao, Jiafei Lyu, Chunjiang Ge, Jiaxin Chen, Weihang Shen, Xiaolong Zhu, and
637 Xiu Li. Using human feedback to fine-tune diffusion models without any reward model. In
638 *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA,
639 USA, June 16-22, 2024*, pp. 8941–8951. IEEE, 2024. doi: 10.1109/CVPR52733.2024.00854. URL
640 <https://doi.org/10.1109/CVPR52733.2024.00854>.
- 641 Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li,
642 Weilin Zhao, Zhihui He, Qianyu Chen, Huarong Zhou, Zhensheng Zou, Haoye Zhang, Shengding
643 Hu, Zhi Zheng, Jie Zhou, Jie Cai, Xu Han, Guoyang Zeng, Dahai Li, Zhiyuan Liu, and Maosong
644 Sun. Minicpm-v: A GPT-4V level MLLM on your phone. *CoRR*, abs/2408.01800, 2024. doi: 10.
645 48550/ARXIV.2408.01800. URL <https://doi.org/10.48550/arXiv.2408.01800>.

648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

Guiwei Zhang, Tianyu Zhang, Mohan Zhou, Yalong Bai, and Biye Li. V2flow: Unifying visual tokenization and large language model vocabularies for autoregressive image generation. *CoRR*, abs/2503.07493, 2025. doi: 10.48550/ARXIV.2503.07493. URL <https://doi.org/10.48550/arXiv.2503.07493>.

Yujie Zhou, Jiazi Bu, Pengyang Ling, Pan Zhang, Tong Wu, Qidong Huang, Jinsong Li, Xiaoyi Dong, Yuhang Zang, Yuhang Cao, Anyi Rao, Jiaqi Wang, and Li Niu. Light-a-video: Training-free video relighting via progressive light fusion. *CoRR*, abs/2502.08590, 2025. doi: 10.48550/ARXIV.2502.08590. URL <https://doi.org/10.48550/arXiv.2502.08590>.

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

707 708 A RELATED WORK

709
710 Modeling user preferences in text-to-image generation is essential for improving alignment with
711 human aesthetics and expectations. Existing research in this area can be broadly categorized into two
712 main categories: general preference modeling, which focuses on capturing collective human judg-
713 ments to enhance overall image quality, and user-specific preference modeling, which personalizes
714 image generation based on individual tastes and behaviors.

715
716 **General Preference Modeling for Human-Aligned Image Generation.** Researchers have explored
717 various strategies to improve alignment, categorized into three approaches: (1) Filtering Training
718 Data with Preference Scores. By selecting training data based on human feedback scores or automated
719 metrics, models can benefit from high-quality examples that reflect specific user demands. For
720 instance, Liang *et al.* (Liang et al., 2024a) demonstrates how filtering data based on feedback scores
721 leads to improved model performance, as it ensures that only the most relevant examples are used
722 for fine-tuning. Similarly, HPS (Wu et al., 2023b;a) builds upon this concept by introducing a
723 scoring mechanism to prioritize image-text pairs closely aligned with user preferences, making
724 the model more responsive to varied user expectations. (2) Reward-Weighted Fine-Tuning for
725 Human-Aligned Models. In this approach, models are fine-tuned using reward signals that weigh
726 heavily on user satisfaction. Lee *et al.* (Lee et al., 2023) exemplifies this by incorporating feedback-
727 based rewards during training, which generates outputs aligned with user preferences. Furthermore,
728 ImageReward (Xu et al., 2023) provides a structured method for translating human judgments into
729 reward functions, which guides the model’s fine-tuning process. By giving greater importance to
730 rewards that capture user satisfaction, these methods tailor the model’s outputs to reflect diverse
731 and nuanced user tastes. (3) Reinforcement Learning for Preference Optimization (Mo et al., 2024;
732 Hao et al., 2023; Chen et al., 2024; Liang et al., 2024b). Recent work (Mo et al., 2024; Hao
733 et al., 2023) uses reinforcement learning to optimize the input prompts for high-quality images.
734 DiffusionDPO (Wallace et al., 2024) leverages user preferences to fine-tune the model, improving its
735 ability to generate images that reflect user choices. D3PO (Yang et al., 2024) eliminates the need to
736 train an explicit reward model by directly fine-tuning the diffusion model using human preference
737 data. Its training strategy is grounded in human preference comparisons and achieves performance
738 comparable to traditional reward-based methods.

739
740 **User-Specific Preference Modeling and Personalized Image Generation.** In recent advancements
741 in personalized image generation, several approaches have emerged to better align generative models
742 with individual needs. While customization-based methods like DreamBooth (Ruiz et al., 2023) and
743 Textual Inversion (Gal et al., 2023) focus on incorporating specific objects or styles through fine-tuning
744 with a few example images, user-preferred personalized image generation takes a different approach
745 by learning broader user preferences and aesthetic tendencies. These approaches, while effective for
746 small datasets, focus on integrating specific instances rather than broader user behaviors. To improve
747 personalization, Salehi *et al.* (Salehi et al., 2024) proposes a standardized process to collect user
748 preferences using a few query images. User feedback is then systematically incorporated to adjust
749 the preferences extracted from the user during the generation process. Additionally, Shen *et al.* (Shen
750 et al., 2024) introduces a method to integrate user-specific preferences across different modalities,
751 such as text and images, creating personalized outputs by leveraging historical interactions, such as
752 clicks and conversations. This multimodal approach significantly enhances the models’ adaptability
753 to align with user needs.

754 755 B MORE QUALITATIVE ANALYSIS RESULTS

756
757 We present a comparison between our model and ViPer (Salehi et al., 2024), supported by qualitative
758 results in Fig. 9, where target images with green borders indicate preferences aligned with the user.
759 Unlike ViPer, which primarily relies on explicit features from reference images, our method leverages

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

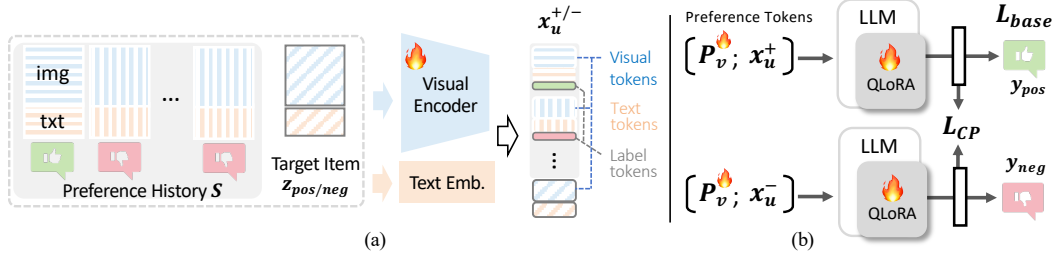


Figure 7: Overview of our MLLM-based preference learning framework. (a) The visual encoder and text embedding module extract preference representations $x_u^{+/-}$ by processing the preference history S and a target item $z_{pos/neg}$.

Multimodal Large Language Models (MLLMs) 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.

C MATHEMATICAL ANALYSIS OF AMBIGUOUS DECISION BOUNDARY CHALLENGES

The ambiguous decision boundary problem is most prominent when comparing user preference choices. As established in the main text, when $\mathcal{M}^+(\mathcal{S}, z) \approx \mathcal{M}^-(\mathcal{S}, z)$, the prediction function yields $Q(\mathcal{S}, z) \approx 0.5$, creating unstable decision boundaries. We now provide a detailed mathematical analysis of how this ambiguity leads to violations of both intra-group consistency and inter-group discrimination requirements in Assumption 1.

C.1 INTRA-GROUP DECISION INCONSISTENCY AT AMBIGUOUS BOUNDARIES

Consider a scenario where users $i, j \in \mathcal{U}_k$ from the same group evaluate an image pair (z_1, z_2) . When the model’s predictions approach the ambiguous boundary of 0.5, the following problematic situation can occur.

For the predicted scores of user i :

$$Q(\mathcal{S}_i, z_1) = 0.5 + \delta_1, \quad Q(\mathcal{S}_i, z_2) = 0.5 - \delta_1 \quad (10)$$

For the predicted scores of user j :

$$Q(\mathcal{S}_j, z_1) = 0.5 - \delta_2, \quad Q(\mathcal{S}_j, z_2) = 0.5 + \delta_2 \quad (11)$$

where δ_1 and δ_2 are small perturbations. Although these predictions satisfy the score similarity constraint from the assumption:

$$|Q(\mathcal{S}_i, z_k) - Q(\mathcal{S}_j, z_k)| \leq \epsilon_k \quad (12)$$

However, due to minute differences around the 0.5 boundary, the pairwise decisions of the two users become completely inconsistent. Specifically, the decision for user i is:

$$D_i(z_1, z_2) = \mathbf{1}[Q(\mathcal{S}_i, z_1) > Q(\mathcal{S}_i, z_2)] = 1 \quad (13)$$

whereas the decision for user j is:

$$D_j(z_1, z_2) = \mathbf{1}[Q(\mathcal{S}_j, z_1) > Q(\mathcal{S}_j, z_2)] = 0 \quad (14)$$

This inconsistency directly violates the intra-group decision consistency requirement from the assumption:

$$\mathbb{P}[D_i(z_1, z_2) = D_j(z_1, z_2)] \geq 1 - \alpha_k \quad (15)$$

810 C.2 INTER-GROUP DECISION CONSISTENCY AT AMBIGUOUS BOUNDARIES

811 Similarly, when users from different groups \mathcal{U}_k and \mathcal{U}_l both exhibit ambiguous predictions near 0.5,
812 their preference scores become unexpectedly similar, leading to undesired inter-group consistency.

813 Consider the case where both users have predictions close to the ambiguous boundary:
814

$$815 \mathcal{Q}(\mathcal{S}_i, z) = 0.5 + \delta_i, \quad \mathcal{Q}(\mathcal{S}_j, z) = 0.5 + \delta_j \quad (16)$$

816 where $|\delta_i|, |\delta_j| \ll 0.5$ are small perturbations. This results in:
817

$$818 |\mathcal{Q}(\mathcal{S}_i, z) - \mathcal{Q}(\mathcal{S}_j, z)| = |\delta_i - \delta_j| \leq |\delta_i| + |\delta_j| \quad (17)$$

819 This proximity violates the inter-group score divergence constraint, as the expected difference can be
820 arbitrarily small:
821

$$822 \mathbb{E}[|\mathcal{Q}(\mathcal{S}_i, z) - \mathcal{Q}(\mathcal{S}_j, z)|] \approx 0 < \max(\epsilon_k, \epsilon_l) \quad (18)$$

823 Furthermore, when both users' predictions hover around 0.5, their pairwise decisions for an image
824 pair (z_1, z_2) may coincidentally align:
825

$$826 D_i(z_1, z_2) = \mathbf{1}[\mathcal{Q}(\mathcal{S}_i, z_1) > \mathcal{Q}(\mathcal{S}_i, z_2)] = D_j(z_1, z_2) \quad (19)$$

827 This unexpected agreement between users from different groups violates the inter-group decision
828 divergence requirement:
829

$$830 \mathbb{P}[D_{u_i}(z_1, z_2) \neq D_{u_j}(z_1, z_2)] < 1 - \beta_{kl} \quad (20)$$

831 To resolve this fundamental issue, our contrastive preference learning method enforces clear prefer-
832 ence discrimination through the following constraints.
833

834 For a positive sample, the positive logit is forced to be significantly greater than the negative logit:

$$835 \mathcal{M}^+(\mathcal{S}, z_{\text{pos}}) \gg \mathcal{M}^-(\mathcal{S}, z_{\text{neg}}) \quad (21)$$

836 For a negative sample, the negative logit is forced to be significantly greater than the positive logit:

$$837 \mathcal{M}^-(\mathcal{S}, z_{\text{neg}}) \gg \mathcal{M}^+(\mathcal{S}, z_{\text{pos}}) \quad (22)$$

840 C.3 INTRA-GROUP CONSISTENCY ANALYSIS

841 When the contrastive constraints are satisfied, clear decision boundaries are established that ensure
842 intra-group consistency.
843

844 For a positive sample z_{pos} , we have:
845

$$846 \mathcal{Q}(\mathcal{S}, z_{\text{pos}}) = \frac{\exp(\mathcal{M}^+(\mathcal{S}, z_{\text{pos}}))}{\exp(\mathcal{M}^+(\mathcal{S}, z_{\text{pos}})) + \exp(\mathcal{M}^-(\mathcal{S}, z_{\text{pos}}))} \gg 0.5 \quad (23)$$

847 Similarly, for a negative sample z_{neg} :

$$848 \mathcal{Q}(\mathcal{S}, z_{\text{neg}}) = \frac{\exp(\mathcal{M}^-(\mathcal{S}, z_{\text{neg}}))}{\exp(\mathcal{M}^-(\mathcal{S}, z_{\text{neg}})) + \exp(\mathcal{M}^+(\mathcal{S}, z_{\text{neg}}))} \ll 0.5 \quad (24)$$

849 Let $\tau > 0$ be a confidence margin. The contrastive constraints then ensure:
850

$$851 \mathcal{Q}(\mathcal{S}, z_{\text{pos}}) \geq 0.5 + \tau, \quad \mathcal{Q}(\mathcal{S}, z_{\text{neg}}) \leq 0.5 - \tau \quad (25)$$

852 In this scenario, for users $i, j \in \mathcal{U}_k$ within the same group, even with score perturbations ϵ_k , decision
853 consistency is guaranteed when $\epsilon_k < \tau$:
854

$$855 |\mathcal{Q}(\mathcal{S}_i, z_{\text{pos}}) - \mathcal{Q}(\mathcal{S}_j, z_{\text{pos}})| \leq \epsilon_k < \tau \quad (26)$$

856 This ensures that the decisions of both users on the same image pair remain consistent:
857

$$858 D_i(z_{\text{pos}}, z_{\text{neg}}) = D_j(z_{\text{pos}}, z_{\text{neg}}) = 1 \quad (27)$$

859 This satisfies the intra-group decision consistency constraint.
860
861
862
863

C.4 INTER-GROUP DISCRIMINATION ANALYSIS

For users from different groups $i \in \mathcal{U}_k$ and $j \in \mathcal{U}_l$ where $k \neq l$, the contrastive constraints create distinct preference distributions that ensure proper inter-group discrimination.

Under the contrastive learning framework, users from different groups develop distinct preference patterns for the same images. Consider the case where user u_i from group \mathcal{U}_k has learned to prefer certain visual patterns, while user u_j from group \mathcal{U}_l has learned different preferences.

For an image z that group \mathcal{U}_k generally likes but group \mathcal{U}_l dislikes, we have:

$$\mathcal{Q}(\mathcal{S}_i, z) \geq 0.5 + \tau, \quad \mathcal{Q}(\mathcal{S}_j, z) \leq 0.5 - \tau \quad (28)$$

This leads to a significant score difference:

$$|\mathcal{Q}(\mathcal{S}_i, z) - \mathcal{Q}(\mathcal{S}_j, z)| \geq |(0.5 + \tau) - (0.5 - \tau)| = 2\tau \quad (29)$$

When $2\tau > \max(\epsilon_k, \epsilon_l)$, this satisfies the inter-group score divergence constraint:

$$\mathbb{E}[|\mathcal{Q}(\mathcal{S}_i, z) - \mathcal{Q}(\mathcal{S}_j, z)|] \geq 2\tau > \max(\epsilon_k, \epsilon_l) \quad (30)$$

Furthermore, for pairwise decisions on an image pair (z_1, z_2) where the groups have opposite preferences, we obtain:

$$D_{u_i}(z_1, z_2) = 1, \quad D_{u_j}(z_1, z_2) = 0 \quad (31)$$

This ensures the inter-group decision divergence requirement is met:

$$\mathbb{P}[D_{u_i}(z_1, z_2) \neq D_{u_j}(z_1, z_2)] = 1 > \beta_{kl} \quad (32)$$

Therefore, by enforcing $\tau > \max(\epsilon_k, \epsilon_l)/2$, our contrastive preference learning method simultaneously satisfies both intra-group consistency and inter-group discrimination constraints specified in Assumption 1.

D MORE EXPERIMENTAL DETAILS

Examples of Training Data. Our dataset, based on Pick-a-Pic v2 dataset (Kirstain et al., 2023), 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. 10 and Fig. 11 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.

Training. To conserve memory, each prompt is truncated to a maximum length of 100 tokens, and input images are resized to 512×512 pixels. Following the setup of (Salehi et al., 2024), we set the length of each user’s preference history sequence, N_{ref} , to 8. The learning rate is set to 1×10^{-5} , with a weight decay of 1×10^{-2} . The language model is fine-tuned using QLoRA (Dettmers et al., 2023), while the vision encoder is trained simultaneously. The input tokens template for the MLLM is “<image>The prompt is <prompt>. Score for this image?<label>”. We first train the MLLM using our custom loss function for 5k steps. After this phase, we continue training for 16k steps, during which both the model and the learnable preference tokens are jointly optimized. To prevent the model from overfitting to a fixed input pattern, we randomly shuffle the order of the reference history sequences during training.

User Preference Dimensions and Attribute Space in the Agent Dataset. To simulate diverse and fine-grained user preferences, we construct a dataset using the Claude-3.5-Sonnet agent. We first define a comprehensive taxonomy of aesthetic attributes spanning multiple key dimensions, as shown in Table 7. These dimensions include art styles, color palettes, compositional strategies, skill levels, visual detail, color hues, and artistic mediums. Each agent is assigned a personalized subset of liked and disliked attributes, sampled from the full attribute space. This configuration enables controllable and individualized preference simulation. The richness of the attribute space ensures that agents exhibit highly diverse and nuanced preferences, mimicking the variability observed in real-world users. Some examples of generated agents are illustrated in Fig. 8, with their corresponding attribute configurations listed in Tab. 8.



Figure 8: Some examples in Agent Dataset.

| Dimension | Example Attributes |
|-----------------|--|
| Art Styles | Surrealism, Aboriginal Art, Ukiyo-e, Romanticism, Anime/Manga, Contemporary Abstraction, Ancient Greek Art, Baroque Art, Abstract Expressionism, Art Deco, Cubism, ... |
| Color Palettes | Oceanic Tones (e.g., Turquoise, Deep Sea Blue), Neon (e.g., Laser Blue, Hot Magenta), Urban Industrial (e.g., Alloy Silver, Iron Black), Pastels (e.g., Mint Green, Peach), Muted Shades (e.g., Dusty Rose), Vibrant Colors, Earthy Palettes, ... |
| Composition | Invented vs. Real Space, Dynamic/Static Tension, Grid-Based Layouts, Pictorial vs. Installation, Foreground vs. Background Contrast, Rule of Thirds, Negative Space Use, Deep/Shallow Space, Balanced or Fragmented Structures, ... |
| Skill Level | Rigorous, Intuitive, Spontaneous, Experimental, Polished, Effortless, Graceful, Heavy-handed, Sophisticated, Controlled, Inventive, ... |
| Detail Level | Tactile, Sharp, Subtle, Elaborate, Vivid, Blurred, Simplified, Defined, Smooth, Intricate, Muted, Textured, ... |
| Hues | Turquoise, Magenta, Burgundy, Indigo, Crimson, Yellow, Slate Gray, Cerulean, Forest Green, Orange, ... |
| Artistic Medium | Mixed Media (e.g., Found Object, Assemblage), Printmaking (e.g., Lithography, Woodcut), Digital (e.g., 3D, Virtual Reality), Traditional Painting (e.g., Tempera, Watercolor), Textile Arts (e.g., Weaving, Embroidery), Ceramics, Sculpture, Drawing, ... |

Table 7: Overview of user preference attribute space across key aesthetic dimensions.

963
964
965
966
967
968
969
970
971

Evaluation Prompt of Claude-3.5-Sonnet. To provide an additional benchmark for evaluating preference prediction accuracy, we employ Claude-3.5-Sonnet, a powerful large multimodal model, as an automated annotator simulating user-level preference reasoning. For each test case, we supply the Claude agent with a set of reference images representing the user’s preferences (liked and disliked examples), along with two candidate images. The agent is instructed to infer visual preference patterns from the references and select the more preferred candidate based on visual alignment. The exact prompt used for each evaluation instance is as follows: "You are given a set of reference images indicating user preferences: the images in <image>, ..., <image> are liked, and those in <image>, ..., <image> are disliked. Based on the visual

| Agent | Dislikes | Likes |
|-------|--|--|
| 1 | Vivid Purple, Radiant Red, Social Realism, Romanticism, Contemporary Abstraction, Mesoamerican Art, Charcoal Black, Pink, Foreground, Negative Space, Closed space, Pictorial, Free-flowing, Effortless, Experimental, Polished, Unfocused, Tactile, Smooth, Ethereal, Turquoise, Burgundy, Collage, Metal, Crocheting, Drypoint | Deep Sea Blue, Jungle Green, Oceanic Art, Traditional African Art, Islamic Art, Buttercream, Alloy Silver, Rule of Thirds, Invented space, Rhythmic, Illusion of Depth, Graceful, Intuitive, Powerful, Meticulous, Elaborate, Soft, Sharp, Muted, Slate Gray, Blue, Magenta, Orange, Decoupage, Virtual Reality, Cyanotype, Found Object |
| 2 | Yellow, Land Art, Situationist Art, Performance Art, Ukiyo-e, Faded Denim, Mint Green, Turquoise, Closed space, Foreground, Invented space, Centralized, Experimental, Inventive, Graceful, Sophisticated, Smooth, Vivid, Expressive, Magenta, Gold, Emerald, Glass, 3D Modeling, Digital Collage, Ink | Naive Art, Contemporary Abstraction, Color Field Painting, Fauvism, Deep Indigo, Jet, Ocean Green, Electric Lime, Pictorial, Golden Ratio, Symmetry, Fragmented, Effortless, Classic, Sophisticated, Subtle, Fine, Abstracted, Sharp, Teal, Orange, Slate, Polaroid, Spray Paint, Watercolor, Found Object |

Table 8: Some examples in Agent Dataset.

patterns and preferences inferred from these references, classify a target image by comparing two candidates: `<image>` and `<image>`. Output only the index (0 or 1) of the image the user would prefer, with explanation."

In each query, `<image>` placeholders are replaced with actual image content using Claude’s multi-modal input interface. The agent’s selection (index 0 or 1) is parsed to compute top-1 accuracy across the test set. The performance of Claude-3.5-Sonnet is reported in Table 3 in main text, achieving a top-1 accuracy of 47.96%. This result offers a meaningful reference point for understanding model performance relative to state-of-the-art language-based multimodal reasoning capabilities.

Image Generation Guided by Our Models. Following the method outlined in (Eyring et al., 2024), 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 (in Eq. (6)) 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 THE ROLE OF LARGE LANGUAGE MODELS (LLMs)

We employed large language models (LLMs) exclusively as writing aids during the preparation of this manuscript. Their use was limited to enhancing readability through language refinement, grammar correction, and stylistic polishing. The models were not involved in formulating research questions, developing methods, conducting experiments, or interpreting findings. All conceptual contributions, technical innovations, and conclusions presented in this paper are entirely attributable to the authors.

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

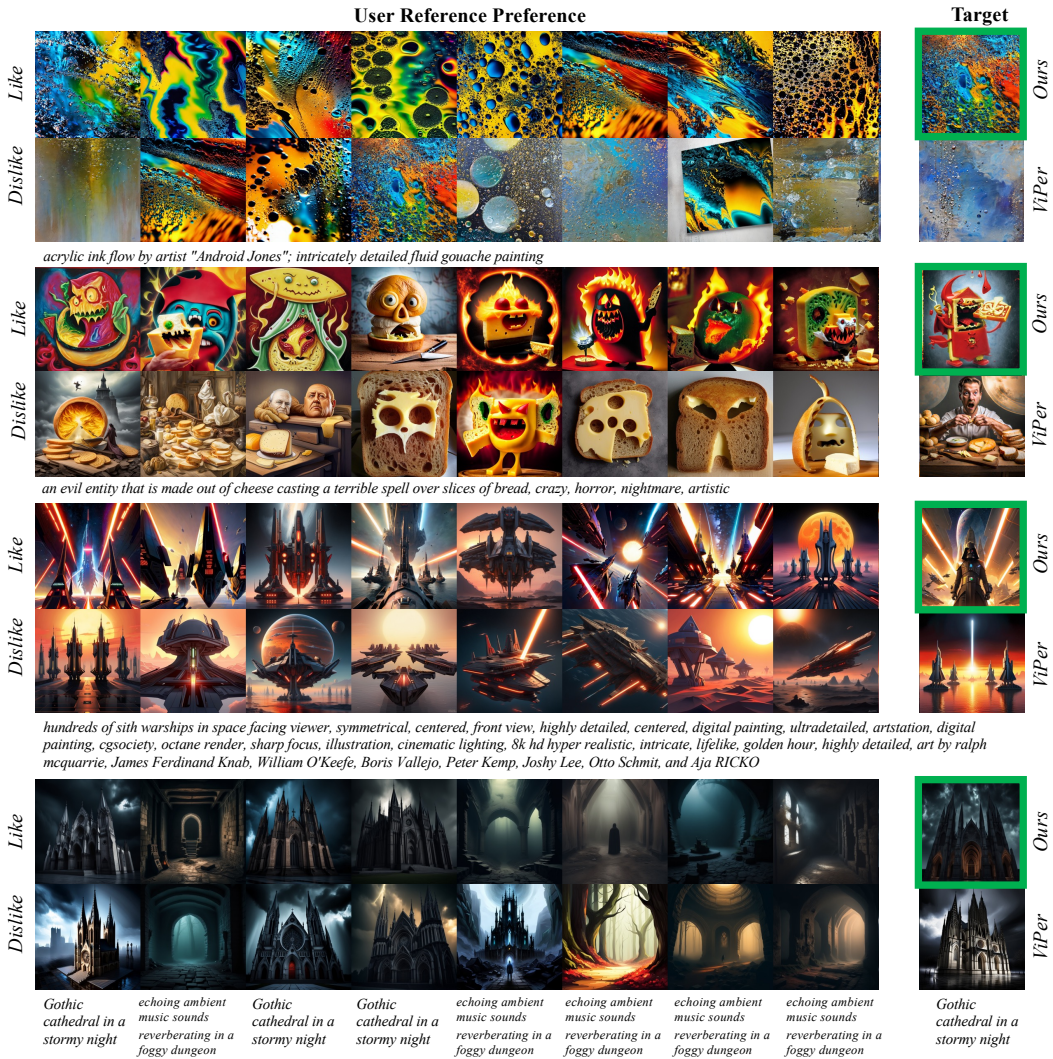


Figure 9: Visual comparison of user-specific preference alignment between our model and ViPer (Salehi et al., 2024) across varying preferences. Target images with green borders indicate preferences aligned with the user. Our method demonstrates effective capture of user-specific personalized results.

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

| | User Reference Preference | | | | | | | | Target |
|---------|---|--|---|--|---|--|---|---|---|
| Like | | | | | | | | | |
| Dislike | | | | | | | | | |
| | <i>Anime black haired tan purple dragon girl with red tracksuit manga purple horns large scaly tail red tracksuit purple eyes</i> | | | | | | | | |
| Like | | | | | | | | | |
| Dislike | | | | | | | | | |
| | <i>selfie of a surfer inside a water tornado</i> | <i>historical archive photo of a mcdonalds found inside of a prehistoric cave</i> | <i>hvn, a city street at night, a picture, by Zoltán Joó, world press photo awarded, blackout, no electricity, traversing a shadowy city, view from street angle, breathtaking, winter snow, kodak ekitar 100, Pentax G45</i> | <i>photo of a panicking surfer inside a water tornado, high in the air</i> | <i>A cute DSLR photo of bear and cow hybrid</i> | <i>hvn, closeup photo portrait of a giant Chinlun, a city street at night, a picture, by Zoltán Joó, world press photo awarded, blackout, no electricity, traversing a shadowy city, view from street angle, breathtaking, winter snow, kodak ekitar 100, Pentax G45</i> | <i>Photorealistic still from silent hill, fog, at night, 4k DSLR photo, ambient lighting</i> | <i>liminal photo of an empty plane, lights turned off, 1990s</i> | <i>focused photo of snowflakes falling in the desert</i> |
| Like | | | | | | | | | |
| Dislike | | | | | | | | | |
| | <i>Sunset reflecting on a crystal ball</i> | <i>flat illustration vector graphics style of camera, glowing gradients, noise textures, modern blue gradients</i> | <i>a flat vector illustration of a painting of a person walking through a forest, a storybook illustration by Petros Afshar, behance contest winner, fantasy art, behance hd, bioluminescence, chromatic</i> | <i>flat vector graphics style illustration of camera, glowing gradients, noise textures, modern blue gradients, behance winner</i> | <i>Generate a vector illustration of a retro car with bold colors and strong lines.</i> | <i>a glowing cube floating in space, blue gradient colors, modern flat illustration, colorful, bright colors, motion graphics</i> | <i>modern flat illustration of man standing on a platform in space, highly detailed background with abstract shapes, blue gradient colors, vector graphics, flat illustration style, 2d art</i> | <i>flat vector graphics style illustration of apple, glowing gradients, noise textures, modern blue gradients, behance winner</i> | <i>modern flat illustration of man standing on a platform in space, highly detailed background with abstract shapes, blue gradient colors, vector graphics, flat illustration style, 2d art</i> |
| Like | | | | | | | | | |
| Dislike | | | | | | | | | |
| | <i>Rob Zombie holding ba pentagram</i> | <i>John 5 holding a pentagram</i> | <i>Andy Warhol holding a sign that says \$</i> | <i>wizard playing electric guitar album cover</i> | <i>Alice Cooper holding a pentagram</i> | <i>A woman holding a sign with a pentacle on it</i> | <i>Alice Cooper sticking tongue out wearing sunglasses holding a sign that says Famous</i> | <i>Jesus holding a pentacle</i> | <i>Alice Cooper sticking tongue out wearing sunglasses holding a sign that says Famous</i> |

Figure 10: Some examples of the training data.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

| User Reference Preference | | | | | | | | Target |
|---------------------------|---|--|--|--|--|--|--|--------|
| Like | | | | | | | | |
| Dislike | | | | | | | | |
| | <i>She wears lilacs in her hair, and picks roses and picks daisys by artist Ralph Horsley</i> | | | | | | | |
| Like | | | | | | | | |
| Dislike | | | | | | | | |
| | <i>motion blur, heavy rain, street photo of an 30 yo Asian woman with short hair, she is laughing</i> <i>a candid shot of lan Mckellen as Gandalf eating soft icecream cone</i> <i>Photo of a blonde girl, intricate cyberpunk respirator and armor</i> <i>a tall woman with purple hair in leather, alcohol, bar, tatoeod, neon light</i> <i>Beautiful woman standing in armour, Futuristic Cyberpunk city</i> <i>still shot from a cyberpunk western, girl fedora firing a handgun</i> <i>Movie still of starwars princess leah cworking as a waitress in a dinner, extremely detailed, intricate, high resolution, hdr, trending on artstation</i> <i>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, avarad winning photography, intricate details</i> <i>Movie still of starwars princess leah cworking as a waitress in a dinner, extremely detailed, intricate, high resolution, hdr, trending on artstation</i> | | | | | | | |
| Like | | | | | | | | |
| Dislike | | | | | | | | |
| | <i>A logo of ai laptop and surveillance camera</i> <i>A logo for a computer vision lady developer</i> <i>A logo for a computer vision lady developer</i> <i>A logo of a laptop and cameras</i> <i>A logo for a computer vision lady developer</i> <i>A logo for a computer vision lady developer</i> <i>A logo of laptop with woman</i> <i>A logo of a woman vollyball playe.</i> <i>A logo of laptop camera and woman</i> | | | | | | | |
| Like | | | | | | | | |
| Dislike | | | | | | | | |
| | <i>logo truck, vector, simply, modern, black and white</i> <i>simply, modern triangle, mountain, geometric design vector, watercolor, travel, tree, blue and pink, white background, primary colors</i> <i>DIY splash right side art colorful detail mountain, travel, river vector, vintage, blue, pink, white background</i> <i>circle border, kids, play, school, sport, fun, comic style</i> <i>DIY promaster Camper VAN art colorful camper detail mountain, tree, splash</i> <i>DIY promaster Camper VAN art colorful camper detail</i> <i>circle border, kids, play, school, sport, fun, comic style</i> <i>play, sport, horse, round, border background, color, comic style</i> <i>red deer stag roaring side view vector logo dark comic style black and white</i> | | | | | | | |

Figure 11: Some examples of the training data.