

Personalized Image Generation with Large Multimodal Models

Anonymous Author(s)

ABSTRACT

Personalized content filtering, such as recommender systems, has become a critical infrastructure to alleviate information overload. However, these systems merely filter existing content and are constrained by its limited diversity, making it difficult to meet users' varied content needs. To address this limitation, personalized content generation has emerged as a promising direction with broad applications. Nevertheless, most existing research focuses on personalized text generation, with relatively little attention given to personalized image generation. The limited work in personalized image generation faces challenges in accurately capturing users' visual preferences and needs from noisy user-interacted images and complex multimodal instructions. Worse still, there is a lack of supervised data for training personalized image generation models.

To overcome the challenges, we propose a *Personalized Image Generation Framework* named Pigeon, which adopts exceptional large multimodal models with three dedicated modules to capture users' visual preferences and needs from noisy user history and multimodal instructions. To alleviate the data scarcity, we introduce a two-stage preference alignment scheme, comprising masked preference reconstruction and pairwise preference alignment, to align Pigeon with the personalized image generation task. We apply Pigeon to personalized sticker and movie poster generation, where extensive quantitative results and human evaluation highlight its superiority over various generative baselines.

CCS CONCEPTS

• Information systems → Recommender systems.

KEYWORDS

Personalized Image Generation, Large Multimodal Models, Preference Alignment

ACM Reference Format:

Anonymous Author(s). 2018. Personalized Image Generation with Large Multimodal Models. In *Proceedings of ACM Conference (Conference 'XX)*. ACM, New York, NY, USA, 10 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

In the era of information overload, individuals are overwhelmed with vast amounts of multimodal content on the Web, underscoring the importance of personalized content delivery [46, 47, 55]. The predominant approach, personalized content filtering like recommender systems [7, 43, 44], relies on user interaction history

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
Conference 'XX, June 03–05, 2018, Woodstock, NY

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-XXXX-X/18/06
<https://doi.org/XXXXXXX.XXXXXXX>

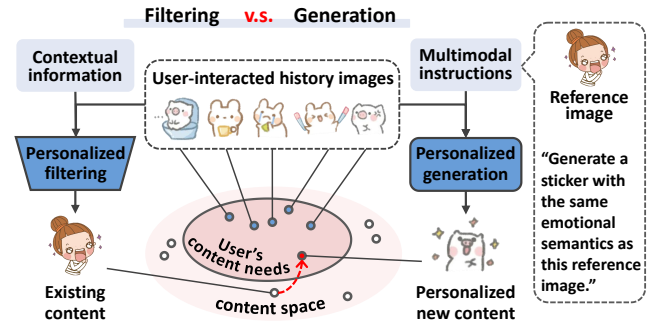


Figure 1: Personalized filtering selects the most relevant existing content while personalized generation creates new and customized ones, more precisely satisfying users' diverse content needs.

and contextual information to infer user preferences and filter existing content. However, this approach is constrained by the limited diversity of available content, rendering it inadequate to fully meet users' varied content needs (see an example in Figure 1). To address this limitation, generating personalized new content is becoming increasingly important across various domains, including personalized movie posters [34], advertisements [40, 50], music [4, 25], and fashion designs [49, 52].

Previous works on personalized content generation primarily focus on personalized text generation [14, 29, 32, 33] while personalized image generation receives little attention. Technically, personalized image generation aims to capture implicit user preferences from user-interacted history images and then integrate users' explicit needs from multimodal instructions to generate personalized target images, as illustrated in Figure 1. Existing methods mainly rely on Diffusion Models (DMs) or Large Language Models (LLMs) for personalized image generation:

- **DM-based methods** [3, 5, 31, 49, 50] might learn the representations of implicit user preferences from user-interacted history images and combine these representations with explicit user instructions for target images to guide the generation of DMs. However, these methods struggle to accurately capture user preferences from noisy history images, which typically cover diverse and complex user interests.
- **LLM-based Personalized Multimodal Generation (PMG)** [34] converts history images and multimodal instructions into textual descriptions, and then utilizes pre-trained LLMs to encode textual descriptions for guiding image generation. However, the discrete nature of text makes it difficult to convey complex visual information in history images and instructions, leading to imprecise representations.

In this light, the key to personalized image generation lies in accurately inferring implicit user preferences from noisy history images while adhering to explicit multimodal instructions for image generation. This necessitates robust multimodal understanding, reasoning, and instruction-following capabilities, driving the adoption of Large Multimodal Models (LMMs) [6, 12] for personalized

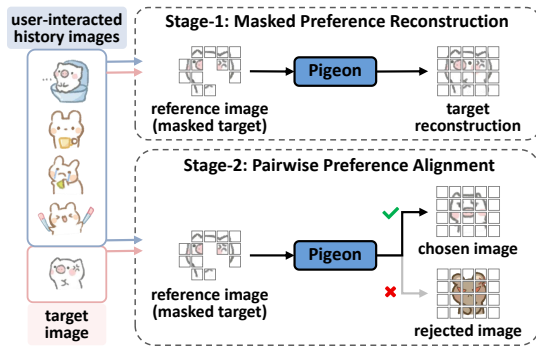


Figure 2: Two-stage preference alignments for Pigeon: given user-interacted images, the last image is treated as the target, with the preceding ones as user history.

image generation. An intuitive approach is to transform history images and multimodal instructions into visual and textual tokens as the input of LMMs for cross-modal understanding and image generation. However, this approach faces critical challenges:

- User-preferred and disliked features (e.g., characters and colors) are often entangled within user-interacted history images, producing fine-grained noise at the feature level. This significantly challenges LMMs to infer implicit user preferences.
- The multimodal instructions may include a reference image alongside textual instructions, e.g., “generate a sticker with the same emotional semantics as this reference image”, requiring LMMs to generate the target image with high-level semantic alignment with the reference image.
- Worse still, existing LMMs are not specifically trained for personalized image generation, making it challenging to infer user preferences and align with multimodal instructions. Furthermore, there is a lack of supervised data containing triplets of \langle user-interacted history images, multimodal instructions, a personalized target image \rangle for LMM training.

To address the challenges, we propose a *Personalized Image Generation Framework* (shorted as Pigeon) for LMMs, comprising three key modules: 1) *Mask generation module* incorporates a mask generator to create token-level masks for reference-aware history filtering, effectively removing noisy signals from the history images at the feature level (cf. Section 2.2.1). 2) *Personalization module* integrates masked history tokens and encodes multimodal instructions with the transformed semantic features of the reference image to generate personalized tokens (cf. Section 2.2.2). 3) *Image generation module* employs a DM to convert the generated personalized tokens into the personalized target image.

Due to the lack of supervised data, Pigeon adopts a two-stage preference alignment scheme to adapt LMMs to the personalized image generation task. As shown in Figure 2, the first stage assumes that user-interacted history images, despite some noise, still partially reflect implicit user preferences. Given a sequence of user-interacted images, Pigeon treats the last one as the target image and the preceding images as the history images. We then mask the target image as a reference image to construct the user’s multimodal instructions and fine-tune Pigeon to reconstruct the target image based on this user’s history images and multimodal instructions, regulating Pigeon to infer user preference from history.

In the second stage, Pigeon generates multiple target images based on the first-stage alignment and ranks them using a preference reward strategy, thus forming pseudo-labeled preference data pairs of “chosen” and “rejected” images. Pigeon is then optimized with the preference data pairs via Direct Preference Optimization (DPO) [28] to generate more personalized target images, enhancing personalization capabilities.

We validate the effectiveness of Pigeon in two popular scenarios: personalized sticker and movie poster generation. Extensive quantitative evaluation demonstrates that Pigeon outperforms the best baseline in personalization, achieving improvements of 7%–31% while maintaining comparable semantic alignment with the reference image. Notably, human evaluation on Amazon MTurk¹ reveals that, on average, 71% participants rate Pigeon-generated images with superior personalization and semantic alignment. Furthermore, we discuss the versatility of Pigeon extending to more domains such as personalized product images, advertisement, and fashion images in Section 2.3, highlighting Pigeon’s broad applicability and significant economic value. Our code and data are available at <https://anonymous.4open.science/r/Pigeon>.

In summary, the key contributions of this work are as follows:

- We empower LMMs with the capability of personalized image generation by the Pigeon framework, which can infer user preferences from noisy history images and integrate multimodal instructions for personalized image generation. Pigeon offers a wide range of applications, catering to diverse user demands and driving the evolution of content delivery paradigms.
- We introduce a two-stage preference alignment scheme to effectively adapt LMMs for the personalized image generation task, eliminating the need for supervised data.
- We propose multiple quantitative evaluation metrics for personalized image generation and conduct extensive experiments across two scenarios. Both quantitative results and human evaluation validate that Pigeon significantly surpasses all the baselines, effectively aligning with personalized user preferences.

2 PERSONALIZED IMAGE GENERATION

In this section, we first formulate the personalized image generation task, followed by the elaboration of our proposed Pigeon framework and its potential applications across various domains.

2.1 Task Formulation

Personalized image generation aims to synthesize personalized images tailored to implicit user preferences and explicit multimodal instructions. Formally, given a set of use-interacted history images $\mathcal{H} = \{x_i\}_{i=1}^N$ and multimodal instructions $\mathcal{R} = \{x_0, txt\}$, where x_0 and txt represent the reference image and textual instruction, respectively, this goal is to generate a personalized target image x_{N+1} that not only meets user visual preferences but also adheres to multimodal instructions by high-level semantic alignment with the reference image. This task has broad applications in enhancing user experience across various domains, such as generating personalized product images in e-commerce or creating personalized movie posters and video thumbnails on platforms like Netflix and YouTube.

¹<https://www.mturk.com/>.

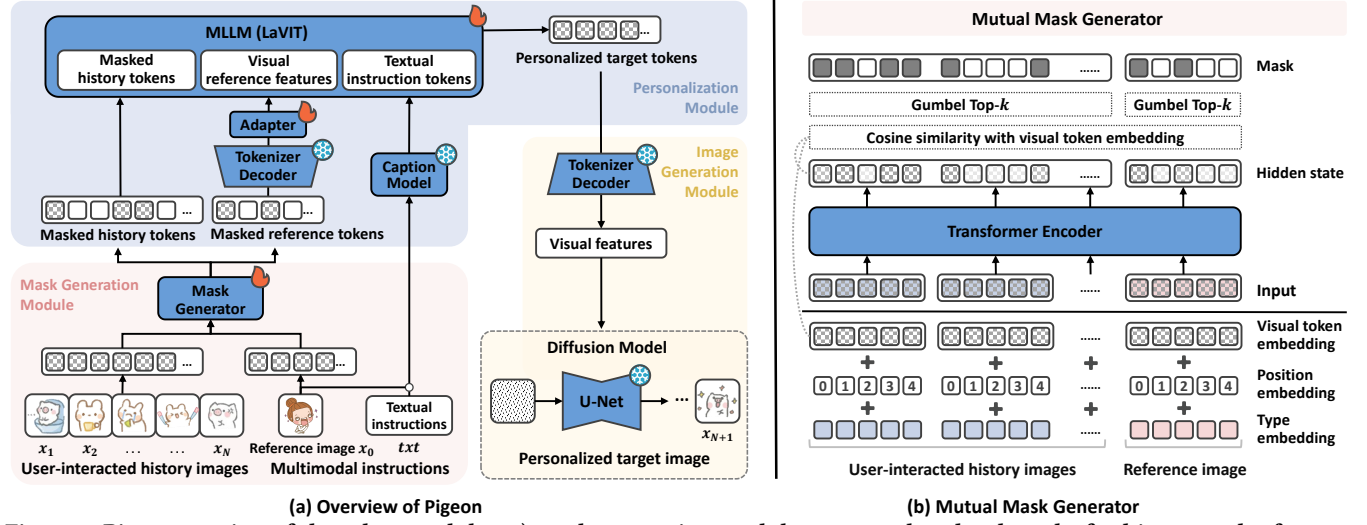


Figure 3: Pigeon consists of three key modules: 1) mask generation module creates token-level masks for history and reference images, 2) personalized module encodes multimodal instructions and integrates them with masked history to generate personalized tokens, and 3) image generation module utilizes these tokens to produce personalized images.

2.2 Pigeon

To achieve personalized image generation, Pigeon leverages a representative LMM named LaVIT [12] for instantiation². Specifically, LaVIT includes a visual tokenizer that translates images into visual tokens for multimodal understanding, and a tokenizer decoder that transforms generated visual tokens into dense visual features to guide image generation. Built upon LaVIT, as depicted in Figure 3(a), Pigeon comprises three key modules: 1) *mask generation module* employs a mask generator to create token-level masks for both history and reference images. 2) *personalization module* extracts high-level semantic features of multimodal instructions and combines them with the masked history tokens to guide LaVIT to generate personalized tokens that reflect users’ content needs. 3) *image generation module* converts these tokens into visual features to generate personalized target images via a DM.

2.2.1 Mask Generation Module. To discard the noise from user-interacted history images, we introduce a mask generator based on a Transformer encoder [41]. It leverages attention mechanisms to encode both history and reference images, and identifies key history tokens that are more relevant to the reference image and contain more personalized information, producing a history mask to filter out noisy tokens. Besides, the mask generator can also create a token-level mask for the reference image to support the two-stage preference alignments, which will be illustrated in Section 2.2.4.

• **Identification of important visual tokens.** Given a set of user-interacted history images \mathcal{H} and a reference image x_0 , we first tokenize these images into visual token sequences:

$$E_i = \text{Visual_Tokenizer}(x_i), i = 0, \dots, N, \quad (1)$$

where $E_i = [e_{i1}, \dots, e_{iL_i}]$ represents the visual token embedding sequence of each image x_i with length L_i , and **Visual_Tokenizer**(\cdot) refers to the visual tokenizer with a visual embedding layer from the pre-trained LaVIT. This process is omitted in Figure 3(a) for brevity.

²Pigeon can also be applied to more LMMs, which is left for future exploration.

The mask generator, as shown in Figure 3(b), combines position and type embeddings with the visual token embeddings via element-wise addition to form the input, which allows the Transformer encoder to distinguish between history and reference tokens and capture the token sequence order within each image. The encoding process is formulated as follows:

$$Z_1, \dots, Z_N, Z_0 = \text{Encoder}(E_1, \dots, E_N, E_0), \quad (2)$$

where $Z_i = [z_{i1}, \dots, z_{iL_i}]$ represents the hidden states of each token sequence E_i , and **Encoder**(\cdot) encapsulates both the element-wise addition and the encoding process. During the encoding process, the attention mechanism allows the visual tokens from both history and reference images to attend to each other, prioritizing important information while reducing the impact of outlier noise. To quantify the importance of each token, we compute the cosine similarity between the hidden states and the original visual token embeddings:

$$s_{ij} = \text{cosine}(z_{ij}, e_{ij}), j = 1, \dots, L_i, \quad (3)$$

where s_{ij} denotes the importance score of the j -th token in each visual token sequence E_i . Intuitively, a higher score indicates more key information is retained in the token.

• **Reference-aware history filtering.** We create a multi-hot binary mask m_h to mask the low-score tokens according to the history mask ratio $\alpha_h \in [0, 1]$. This mask filters out noisy or reference-irrelevant history tokens, yielding the filtered token embeddings for each history image:

$$[\tilde{E}_1, \dots, \tilde{E}_N] = m_h \odot [E_1, \dots, E_N], \quad (4)$$

where \tilde{E}_i denotes masked history token embeddings. For gradient backpropagation in this discrete sampling process, the Gumbel-Softmax trick [22] is applied to the non-differentiable binary mask.

2.2.2 Personalization Module. To effectively handle multimodal instructions, this module first encodes them to extract essential high-level semantic features, then combines these features with masked history tokens into a hybrid prompt, which serves as the input to LMM, enabling the generation of personalized tokens.

• **Multimodal instructions encoding.** When directly utilizing the reference image to guide target image generation, LMMs often duplicate the reference image, failing to effectively incorporate personalized information (see empirical results in Section 3.4.2). This highlights the necessity to extract high-level semantics from the reference image for image generation. To enrich the semantics of the reference image x_0 and enhance the comprehension of multimodal instructions in LMMs, we utilize a caption model (e.g., BLIP-2 [17] and LLaVA [20]) to generate a textual description of the reference image, which is then tokenized into textual tokens:

$$r_t = \text{Text_Tokenizer}(\text{Caption}(x_0)), \quad (5)$$

where r_t refers to the high-level textual semantic features extracted from the reference image, and $\text{Text_Tokenizer}(\cdot)$ denotes the text tokenizer with the word embedding layer from LaVIT.

For visual semantics, we transform the low-level reference token embedding sequence E_0 into high-level dense visual features. Here, we utilize the pre-trained tokenizer decoder of LaVIT for the transform to avoid introducing extra parameters, followed by average pooling to aggregate the multiple feature vectors from the tokenizer decoder:

$$v_0 = \text{AvgPooling}(\text{Tokenizer_Decoder}(E_0)). \quad (6)$$

Next, an adapter layer is introduced to align the feature dimension of v_0 with the LaVIT embeddings, i.e., $r_v = \text{Adapter}(v_0)$, where r_v denotes the extracted high-level visual semantic features.

• **Hybrid prompt for LMM.** To integrate these encoded semantic features with filtered history into prompts for LMM instruction tuning, we propose a hybrid prompt that is structured as follows:

$$p = \text{Prompt}(\tilde{E}_1, \dots, \tilde{E}_N, r_t, r_v). \quad (7)$$

Instruction: You are a helpful personalized assistant. You will receive a list of user-liked images that reflect the user’s visual preferences. By analyzing user preferences, please generate a personalized image that aligns with the user’s aesthetic taste and the semantics in a specified reference image.

Input: The user likes the following images: $\tilde{E}_1, \dots, \tilde{E}_N$. The reference image: r_t, r_v .

Response: <Personalized Target Tokens E_{N+1} >

By using a hybrid prompt similar to the above one, LMMs can adapt to various scenarios to generate personalized target tokens.

2.2.3 Image Generation Module. With personalized target tokens E_{N+1} , the pre-trained tokenizer decoder of LaVIT converts these discrete tokens into dense visual features, which can guide the generation of the personalized target image x_{N+1} in DM.

2.2.4 Two-stage Preference Alignments. To optimize Pigeon for personalized image generation, an intuitive strategy is maximizing the generation likelihood of the target tokens E_{N+1} , based on the prompt p in Eq. (7). However, since there is no supervised dataset containing triplets of <user-interacted history images, multimodal instructions, personalized target image>, we propose a two-stage preference alignment process for effective instruction tuning.

• **Stage-1: Masked Preference Reconstruction.** We assume that user-interacted history images, despite containing some noise,

still reflect user visual preferences. Based on this, as shown in Figure 2, given a sequence of user-interacted images $\{x_i\}_{i=1}^{N+1}$, the last one x_{N+1} is considered the personalized target image, while the preceding images are treated as history images $\mathcal{H} = \{x_i\}_{i=1}^N$.

Supervised dataset construction. Considering the lack of multimodal instructions, we adopt the target image as the reference to construct multimodal instructions $\mathcal{R} = \{x_{N+1}, \text{txt}\}$. A token-level reference mask is then applied to corrupt the reference image, encouraging the model to extract user preferences from history images for target reconstruction. Specifically, we utilize the importance score defined in Eq. (3) to rank all the reference tokens and create the token-level mask for the reference image.

Unlike the history mask, which filters out noise by discarding low-score tokens, we introduce a dual-phase mask scheme for the reference image. During training, we mask high-score reference tokens, which contain more personalized information (as discussed in Section 2.2.1), forcing the model to rely on history images to recover the target. During inference, low-score tokens are masked instead, utilizing the preference reconstruction capability to generate more personalized content. Formally, the dual-phase mask m_r with a reference mask ratio $\alpha_r \in [0, 1]$ is applied to the reference tokens by $\tilde{E}_0 = m_r \odot E_0$. We then replace E_0 in Eq. (6) with \tilde{E}_0 to derive the modified visual features for the hybrid prompt p in Eq. (7), optimizing the model to reconstruct the target token sequence E_{N+1} . In this way, we could construct a supervised prompt-response dataset $\mathcal{D} = \{(p, E_{N+1})^k\}_k$ from the available interaction sequences for masked preference reconstruction.

Supervised fine-tuning. For parameter-efficient fine-tuning, we introduce a LoRA [9] module into the pre-trained LaVIT, which keeps the LaVIT parameters frozen and imports trainable low-rank decomposition matrices for updates. As shown in Figure 3(a), we only fine-tune specific components of Pigeon, namely the mask generator, adapter, and LoRA for LaVIT, while freezing all the other parameters. During training, we randomly sample the reference mask ratio $\alpha_r \in [0, 1]$ and fine-tune Pigeon for target reconstruction, aiming to capture more robust user preferences. Formally, the loss function is defined as the negative likelihood of the target token sequence via an auto-regressive manner:

$$\mathcal{L}_{sft} = - \sum_{(p, E_{N+1}) \in \mathcal{D}} \sum_{j=1}^{L_{N+1}} \log (P_{\Theta}(e_{N+1,j} | p(\alpha_r), e_{N+1, <j})), \quad (8)$$

where $e_{N+1,j}$ is the j -th token in the sequence E_{N+1} of length L_{N+1} , $p(\alpha_r)$ is the hybrid prompt with a uniformly sampled reference mask ratio, and Θ includes all the learnable parameters of Pigeon.

• **Stage-2: Pairwise Preference Alignment.** After the first-stage fine-tuning, Pigeon is capable of following the instructions for personalized image generation. To further enhance its personalization capability, we adopt DPO [28] for pairwise preference alignment, which utilizes preference pairs of chosen and rejected responses to optimize the model to produce the chosen one.

Preference dataset construction. To construct the preference data pairs for DPO, we first generate multiple personalized target token sequences for each prompt $p(\alpha_r)$ with varying reference mask ratios $\alpha_r \in \{0.0, 0.1, \dots, 1.0\}$ based on the first-stage alignment. These tokens are then transformed into images x via the image generation module. To identify the best and worst personalized

images, we introduce a preference reward strategy to rank all generated images. Following [34], we compute the CLIP similarity between each generated image and the history images:

$$s(\alpha_r) = \frac{1}{N} \sum_{i=1}^N \text{CLIPSim}(x(\alpha_r), x_i), \quad (9)$$

where $s(\alpha_r)$ is the preference score of image $x(\alpha_r)$. We rank the generated images based on these scores to form the pseudo-labeled preference dataset $\tilde{\mathcal{D}} = \{(\mathbf{p}, E', E'')^k\}_k$, where E' and E'' denote the chosen and rejected token sequences for DPO, corresponding to images with the highest and lowest preference scores.

Preference optimization. In this stage, we continue updating the LoRA weights while keeping all the other parameters frozen. With the preference dataset, the loss function can be formulated as:

$$\mathcal{L}_{DPO} = -\mathbb{E}_{(\mathbf{p}, E', E'') \sim \tilde{\mathcal{D}}}_{\alpha_r \sim \mathcal{U}(0,1)} \left[\log \sigma \left(\beta \frac{P_{\Theta_l}(E' | \mathbf{p}(\alpha_r))}{P_{\Theta_l}(E' | \mathbf{p}(\alpha_r))} - \beta \frac{P_{\Theta_l}(E'' | \mathbf{p}(\alpha_r))}{P_{\Theta_l}(E'' | \mathbf{p}(\alpha_r))} \right) \right], \quad (10)$$

where Θ_l denotes the learnable parameters of the LoRA module, and β is a parameter controlling the deviation from the reference model $\hat{\Theta}_l$ obtained in the first-stage alignment.

2.2.5 Inference. To manage the trade-off between personalization and semantic alignment with the reference image, users could adjust the reference mask ratio to control how much reference information is retained in the generated images. During inference, given history images $\mathcal{H} = \{x_i\}_{i=1}^N$, multimodal instructions $\mathcal{R} = \{\mathbf{x}_0, \mathbf{txt}\}$ and a user-specified reference mask ratio α_r , Pigeon can mask the low-score reference tokens accordingly to generate an image x_{N+1} that aligns with the user’s visual preferences and multimodal instructions.

2.3 Domain Applications of Pigeon

Pigeon empowers LMMs with the capability to generate personalized images, which is applicable in various scenarios such as personalized stickers on social media platforms like Twitter and personalized movie posters on platforms like Netflix (see demonstration in Section 3). Beyond these, we showcase the potential of Pigeon in other representative domains.

• **E-commerce: personalized product images.** In e-commerce, compelling product images are crucial for drawing attention and driving purchase decisions. Pigeon can analyze user visual preferences from their behaviors to generate personalized product images that match individual tastes in personalized display style and background, delivering a more customized shopping experience.

• **Advertising: personalized advertisements.** Pigeon can assist advertisers in creating highly customized and context-aware multimodal advertisements based on user behaviors, which are more likely to improve user engagement and conversion rates.

• **Fashion: personalized fashion designs.** Pigeon can infer users’ fashion preferences to generate personalized designs for fashion products like clothing, shoes, and jewelry. Besides, both fashion designers and users can provide their preferred fashion images with explicit multimodal instructions for Pigeon to customize designs, fostering an interactive and collaborate design experience.

3 EXPERIMENTS

We evaluate Pigeon in sticker and movie poster scenarios to validate its superiority by answering the following research questions:

- **RQ1:** How does Pigeon perform compared with DM-based, LLM-based, and LMM-based personalized image generation methods, based on quantitative evaluation?
- **RQ2:** Can Pigeon surpass the baselines in human evaluation?
- **RQ3:** How do the special designs of Pigeon (e.g., history mask, multimodal instruction encoding strategy, and two-stage preference alignment process) affect the performance?

3.1 Experimental Settings

3.1.1 Datasets. We conduct experiments on two publicly available datasets, focusing on sticker and movie poster scenarios: 1) **SER30K**³ is a large-scale dataset of stickers, each categorized by theme and annotated with an associated emotion label; and 2) **ML-Latest**⁴, a benchmark dataset containing user ratings on movies.

For the sticker scenario, we exclude low-quality themes or those with fewer than six stickers, constructing user interaction sequences where each user interacts with a single theme. For the movie scenario, we adopt the small version of the dataset, retaining user interactions with ratings of four or higher, sorted by the timestamps. We apply a sliding window of six interactions, moving one step at a time to create data samples for each user in both scenarios. Each sample treats the first five interactions as the user history images and the last as the target image. We split the samples into training, validation, and testing sets with a ratio of 8:1:1. In the sticker testing set, we randomly select one sticker from a different theme than the user history as the reference image, while in the movie poster scenario, the target image is used as the reference. Dataset statistics are summarized in Table 4 in Appendix.

3.1.2 Baselines. We compare Pigeon with various generative baselines, including methods based on DMs, LLMs, and LMMs: 1) **Textual Inversion (TI)** [5] introduces a word embedding to learn user preference representation, which is then combined with textual instructions to guide the text-to-image generation process in DMs. 2) **PMG** [34] transforms user-interacted and reference images into textual descriptions, using pre-trained LLMs to extract user preferences through keywords and implicit embeddings to condition the image generator. 3) **LlVA** [20] is an LMM designed to extract dense image features for visual reasoning, generating text by default but capable of producing images when integrated with an external text-to-image generator. 4) **LaVIT** [12] is another LMM that converts images into discrete visual tokens for reasoning and generates visual tokens to guide the image generation process.

Additionally, we include two results for reference: 5) **Recon**, which utilizes the visual tokenizer, tokenizer decoder, and DM of the pre-trained LaVIT for image reconstruction without personalization; and 6) **Grd**, representing the evaluation results of the reference images. The performance gap between Recon and Grd reflects the difference between generated and real-world images.

3.1.3 Evaluation Metrics. We employ various quantitative evaluation metrics for performance comparison. Following [34, 36], we mainly focus on **personalization** and **semantic alignment** with the reference image by measuring the semantic and perceptual similarity between generated and history/reference images.

³<https://github.com/nku-shengzheliu/SER30K>.

⁴<https://grouplens.org/datasets/movielens>.

Table 1: Quantitative performance comparison between Pigeon and the baselines in both scenarios. Baselines labeled with “*” indicate the pre-trained models. The best results are highlighted in bold, while the second-best results are underlined.

| #Sticker | | Overall | Personalization | | | | | Semantic Alignment | | | Fidelity |
|---------------|--------|--------------|-----------------|----------------|----------------|--------------------|--------------------|--------------------|----------------|----------------|------------------|
| Methods | | | CS \uparrow | CIS \uparrow | DIS \uparrow | LPIPS \downarrow | MS-SSIM \uparrow | CS \uparrow | CIS \uparrow | DIS \uparrow | FID \downarrow |
| DM-based | TI | <u>36.91</u> | 18.67 | 40.90 | 36.58 | 0.7654 | 0.0887 | 32.91 | <u>53.67</u> | <u>48.50</u> | 105.48 |
| LLM-based | PMG | 32.83 | <u>19.16</u> | 47.34 | 39.15 | 0.7383 | 0.0827 | 18.31 | 45.45 | 37.80 | <u>84.91</u> |
| LMM-based | LLaVA* | 32.40 | 17.88 | 47.26 | 42.59 | 0.7575 | 0.0966 | 17.54 | 42.65 | 39.25 | 93.23 |
| | LLaVA | 32.23 | 18.72 | 37.44 | 33.19 | 0.7552 | 0.0851 | <u>27.02</u> | 49.15 | 43.88 | 95.19 |
| | LaVIT* | 34.56 | 18.77 | <u>53.63</u> | <u>50.96</u> | <u>0.6855</u> | <u>0.1376</u> | 15.49 | 40.76 | 39.09 | 107.53 |
| | LaVIT | 33.15 | 16.39 | 40.56 | 40.84 | 0.7377 | 0.1128 | 25.74 | 70.80 | 69.93 | 83.39 |
| Reference | Recon | 33.22 | 16.30 | 40.60 | 40.76 | 0.7370 | 0.1126 | 25.84 | 71.09 | 70.14 | 83.57 |
| | Grd | 36.98 | 16.93 | 45.00 | 43.71 | 0.6443 | 0.1349 | 28.95 | 100.00 | 100.00 | - |
| #Movie poster | | Overall | Personalization | | | | | Semantic Alignment | | | Fidelity |
| Methods | | | CS \uparrow | CIS \uparrow | DIS \uparrow | LPIPS \downarrow | MS-SSIM \uparrow | CS \uparrow | CIS \uparrow | DIS \uparrow | FID \downarrow |
| DM-based | TI | <u>31.07</u> | 12.41 | 28.29 | 19.18 | 0.7721 | 0.0399 | 33.84 | 43.53 | 39.81 | 79.77 |
| LLM-based | PMG | 20.36 | 13.61 | 25.11 | 22.73 | 0.7692 | 0.0261 | 15.60 | 27.29 | 25.15 | 77.25 |
| LMM-based | LLaVA* | 22.08 | 12.24 | 29.60 | 19.73 | 0.7607 | 0.0373 | 14.55 | 31.76 | 21.99 | 73.77 |
| | LLaVA | 30.59 | 12.62 | <u>30.64</u> | 19.33 | 0.7690 | 0.0370 | <u>30.53</u> | <u>48.50</u> | 41.45 | 54.55 |
| | LaVIT* | 23.81 | 12.64 | 28.23 | 17.50 | <u>0.7546</u> | <u>0.0458</u> | 19.39 | 36.93 | 37.71 | 50.08 |
| | LaVIT | 27.82 | <u>13.86</u> | 30.49 | 19.95 | 0.7548 | 0.0370 | 25.15 | 46.02 | 60.07 | 33.53 |
| Reference | Recon | 27.81 | 13.85 | 30.33 | 19.95 | 0.7548 | 0.0367 | 25.29 | 46.08 | 60.52 | 33.74 |
| | Grd | 41.58 | 10.94 | 51.34 | 20.75 | 0.7502 | 0.0402 | 31.81 | 100.00 | 100.00 | - |

- **Semantic similarity.** We adopt CLIP [27] and DINO [24], two popular visual feature extraction models, to extract image features from generated and history/reference images, and compute the cosine similarity between them to obtain the CLIP Image Score (CIS) and DINO Image Score (DIS). Additionally, the CLIP Score (CS) measures the similarity between generated images and textual descriptions of the history/reference images. To assess the overall performance, we also calculate a unified F1-score, combining the history CIS and reference CS.
- **Perceptual similarity.** To evaluate finer-grained visual personalization, we apply LPIPS [53] and MS-SSIM [45] to quantify the perceptual similarity between generated and history images.
- **Fidelity.** We also employ the widely-used FID metric to assess the fidelity of the generated images.

3.1.4 Implementation Details. All the baselines are tuned with a fixed learning rate of $1e^{-5}$. We implement PMG following its default model designs, while other baselines are implemented with Stable Diffusion XL [26] as the image generator for fair comparisons.

In Pigeon, the learning rate is also set to $1e^{-5}$, and the history mask ratio α_h is fixed at 0.2. During inference, we select the optimal reference mask ratio $\alpha_r \in \{0.0, 0.1, \dots, 1.0\}$ for each reference image by averaging the history CIS and reference CS. All experiments are conducted using a single NVIDIA-A100 GPU.

3.2 Quantitative Evaluation (RQ1)

The comparison between Pigeon and the baselines is shown in Table 1. The observations are summarized as follows:

- DM-based TI outperforms most baselines in semantic alignment by directly using the textual description of the reference image for text-to-image generation. However, noisy signals in interaction history hinder its ability to precisely capture user preferences, resulting in inferior personalization.

- PMG converts images into textual descriptions and uses LLMs to infer user preferences for guiding image generation. The image-to-text conversion may overlook critical visual details, leading to inaccurate preference modeling and multimodal instruction understanding. As a result, PMG presents moderate performances in both personalization and semantic alignment.
- The decent performance of the pre-trained LLaVA and LaVIT in personalization validates the strength of advanced instruction-following and visual understanding capabilities in LMMs for personalized image generation. Among them, LLaVA relies on personalized text to guide image generation, which can cause misalignments between expressed textual preferences and actual visual preferences, resulting in relatively lower performance.
- After fine-tuning in each scenario, both LLaVA and LaVIT tend to reconstruct reference images rather than generate personalized ones, as evidenced by significant improvements in semantic alignment alongside a decline in personalization. This is mainly due to the lack of supervised data for model training.
- Pigeon exhibits superior performance in most personalization metrics across two scenarios, while maintaining comparable semantic alignment and fidelity. These results underscore the effectiveness of Pigeon in capturing user visual preferences from noisy history images and accurately understanding multimodal instructions to produce personalized images.

3.3 Human Evaluation (RQ2)

To assess the qualitative performance of Pigeon in personalization and semantic alignment, we conduct a human evaluation on Amazon MTurk⁵, comparing it against Grd and two representative baselines: 1) TI, which exhibits the second-best overall performance

⁵<https://www.mturk.com/>.

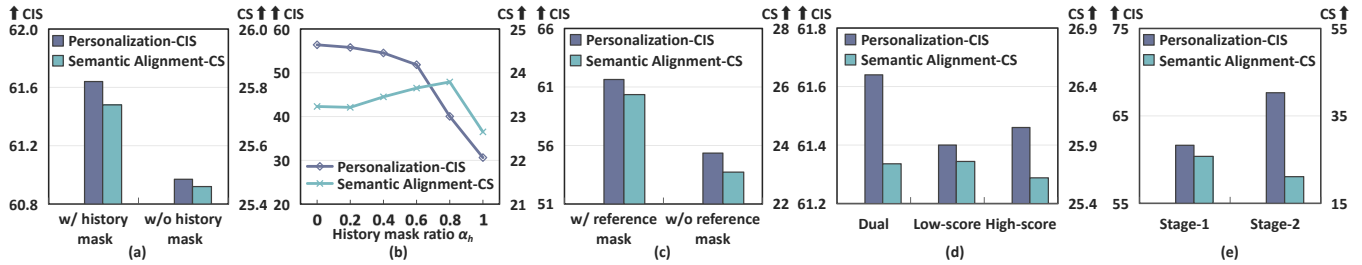


Figure 4: In-depth analysis of the history mask and the two-stage preference alignment process.

Table 2: The human evaluation results, where “±” denotes 95% confidence interval. Pigeon is consistently preferred (≥ 50%) over the baselines across sticker and movie poster scenarios.

| Pigeon | | Grd | TI | PMG |
|--------------------|---------|------------|------------|------------|
| Personalization | Sticker | 0.91±2.19% | 0.91±2.19% | 0.89±1.79% |
| | Movie | 0.62±2.85% | 0.66±2.16% | 0.57±2.51% |
| Semantic Alignment | Sticker | - | 0.54±2.67% | 0.67±3.65% |
| | Movie | - | 0.58±2.58% | 0.73±2.22% |

in Table 1, and 2) PMG, designed for personalized image generation. The evaluation adopts binary-choice tests across sticker and movie poster scenarios, each with 50 cases. For personalization, we present five user-interacted history images and the generated images, with the question: “When provided with someone’s five previously liked stickers (movies), please select the next sticker (movie poster) that is more attractive to her/him.” For semantic alignment, we display the reference and generated images with the question: “Which image aligns more closely with the semantics of the reference image?” As shown in Table 2, Pigeon consistently surpasses (≥ 50%) the baselines, even the Grd, in personalization and maintains decent results in semantic alignment with reference images. These findings emphasize its superiority in capturing user preferences from noisy history images and effectively integrating multimodal instructions for image generation, which aligns with the quantitative analysis.

3.4 In-depth Analysis (RQ3)

In this section, we conduct additional experiments in the sticker scenario to further investigate the effects of various Pigeon designs, including the history mask, multimodal instruction encoding strategy, and the two-stage preference alignment process. To reduce resource costs, we mainly focus on the results after first-stage preference alignment for fair comparisons.

3.4.1 Effect of history mask. To assess the effectiveness of the history mask in managing noisy history images, we exclude it during training and present the results on two key metrics in Figure 4(a). The findings show that: 1) noise in the history images prevents the model from accurately capturing user preferences and even disrupts the semantic alignment with the reference image. 2) The history mask could effectively filter out the noisy signals, thereby enhancing model performance.

Additionally, we vary the history mask ratio α_h during inference, with the reference mask ratio fixed at 0.5. The results in Figure 4(b) reveal that increasing α_h discards both noise and useful personalized information in history images, causing Pigeon to rely more on the reference image, thus slightly improving the semantic

Table 3: Effects of multimodal instruction encoding.

| | Personalization CIS↑ | LPIPS↓ | Semantic Alignment CS↑ |
|----------------------|-------------------------|--------|---------------------------|
| Pigeon | 61.64 | 0.6800 | 25.74 |
| - w/o visual feature | 55.46 | 0.6828 | 23.53 |
| - w/o textual tokens | 65.73 | 0.6731 | 20.37 |
| - w/o encoding | 55.35 | 0.6976 | 24.72 |

alignment. However, this also makes it harder for Pigeon to extract user preferences, reducing the performance in personalization.

3.4.2 Effect of multimodal instruction encoding. To validate the necessity to extract high-level semantics via multimodal instruction encoding, we perform three ablation studies during the training phase. Specifically, we remove the encoded visual features and textual instruction tokens, referred to as “w/o visual features” and “w/o textual tokens”, respectively. We also disable the encoding process by directly inputting the masked reference tokens into LaViT, denoted as “w/o encoding”. Results on three key metrics, reported in Table 3, reveal the following insights: 1) removing the visual features reduces the performance, highlighting the importance of high-level visual semantics for understanding the reference image and enhancing personalization. 2) Excluding textual tokens improves personalization while significantly reducing semantic alignment, indicating that the model over-prioritizes user preferences when textual semantics are absent. 3) Disabling the encoding process leads to simple duplication of the reference image rather than true personalization, as evidenced by a notable drop in personalization and an increase in semantic alignment.

3.4.3 Effect of two-stage preference alignments.

• **Stage-1: masked preference reconstruction.** To evaluate the impact of the first-stage masked preference reconstruction, we perform additional experiments that analyze the effect of the reference mask and explore alternative masking schemes: 1) removing the reference mask, as shown in Figure 4(c), leads to a notable performance decline, underscoring the importance of masked preference reconstruction, which allows Pigeon to effectively integrate user preferences with reference semantics for personalization. 2) Exploring alternative masking schemes for the reference tokens: “Low-score” refers to masking low-score tokens during both training and inference, while “High-score” masks high-score tokens in both phases. These schemes are compared to the dual mask scheme of Pigeon, with results presented in Figure 4(d). The significant decline in personalization suggests that masking either high-score or low-score tokens during both phases causes the model to over-focus on preference reconstruction, limiting its ability to generalize this reconstruction for broader personalization.



Figure 5: Examples of generated images in sticker and movie poster scenarios, along with four user-interacted history images and one reference image.

• **Stage-2: pairwise preference alignment.** We evaluate the effect of the second-stage pairwise preference alignment by comparing the performance after the first and second stages of alignments, as shown in Figure 4(e). Despite a slight decline in semantic alignment, the second-stage preference alignment further enhances personalization. This demonstrates the effectiveness of DPO in aligning the generation process more closely with user preferences, ultimately resulting in more personalized image generation.

3.5 Case Study

In this section, we present two examples of Pigeon-generated images in sticker and movie poster scenarios, along with four user-interacted history images and one reference image. We compare Pigeon with two competitive baselines, TI and PMG, as shown in Figure 5. In the sticker scenario, Pigeon effectively captures the user’s visual preference for Yoda and integrates it with the high-level semantics of the reference sticker, such as “drinking coffee”, achieving impressive personalization and semantic alignment with the reference image. In the movie poster scenario, Pigeon-generated poster for the movie “Rise of the Planet of the Apes” showcases high semantic alignment with the reference poster by emphasizing an intense central ape figure, evoking a similar sense of power and conflict. Meanwhile, it matches the user’s preference for character-centered movie posters with a dark and dramatic color palette. More examples are provided in Figure 6 and Figure 7 in Appendix.

4 RELATED WORK

• **Personalized Content Filtering.** Traditional filtering-based personalized content delivery approaches, such as recommender systems [2, 19, 43, 54], typically rank existing content based on user interaction history and contextual information, delivering the top-ranked content. However, constrained by the limited diversity of available content, they often fall short of meeting users’ diverse needs [42, 49, 52], motivating the emergence of personalized content generation across various domains.

• **Personalized Content Generation.** The rise of powerful generative models, such as DMs [26, 30], LLMs [39], and LMMs [12, 20], has sparked increasing interest in their potential for personalized content generation. Most previous work focuses on personalized text generation [14, 29, 32, 37]. For example, the LaMP benchmark [33] is developed to train and evaluate LLMs in various personalized text scenarios like personalized news headline generation and

tweet paraphrasing. Further work, such as RSPG [32], studies the retrieval-augmented solutions to personalize LLM outputs, while PER-PCS [37] introduces a parameter-sharing framework to enable more efficient and fine-grained personalization.

In contrast, personalized image generation has received relatively less attention. Current research mainly adopts DMs and LLMs for this task: 1) DM-based methods [8, 15, 18, 35], such as TI [5] and DreamBooth [31], focus on aligning image generation with users’ explicit multimodal instructions, without consideration of user implicit visual preferences. Other approaches like DiFashion [49], CG4CTR [50], and AdBooster [36], integrate user data (*e.g.*, interaction history and user features) with multimodal instructions to guide personalized fashion and product image generation. However, these methods often struggle with the noisy signals in user-interacted history images, leading to inaccurate preference modeling. 2) LLM-based PMG [34] translates images into texts for the LLM to extract user visual preferences, while the limitations of text in conveying complex visual details hinder its effectiveness. In this work, we leverage the notable visual understanding and reasoning capabilities of LMMs, along with dedicated modules, to develop the Pigeon framework that effectively handles noisy history images for accurate, tailored image generation.

• **Multimodal Content Generation.** A lot of prior studies utilize pre-trained generative models for content generation across various modalities, including image [10, 13], text [17, 51], video [11, 21], and audio [16, 23, 48]. From these works, we have witnessed the impressive capabilities of the pre-trained LMMs, such as GPT-4 [1], LLaVA [20], LaVIT [12], and Gemini [38] in instruction-following, multimodal content understanding, and generation. However, despite their success, these models primarily generate content conditioning on text prompts or other given modalities, without incorporating users’ personalized information. When directly applied to personalized content generation, these models often exhibit suboptimal performance (*cf.* the empirical results of pre-trained LaVIT and LLaVA in Table 1) due to their limited understanding of user preferences. Therefore, we propose the Pigeon framework, empowering the pre-trained LMMs with personalization capabilities.

5 CONCLUSION AND FUTURE WORK

In this work, we proposed a novel framework named Pigeon, which integrates a pre-trained LMM with specialized modules to infer implicit user preferences from noisy user history and incorporate explicit multimodal instructions for personalized image generation. To alleviate data scarcity, Pigeon adopts a two-stage preference alignment scheme with masked preference reconstruction and pairwise preference alignment, enhancing the personalization capabilities of LMMs. Both quantitative results and human evaluation highlight Pigeon’s effectiveness in generating personalized images.

This work marks an initial attempt to align pre-trained LMMs with implicit user visual preferences, paving the way for several promising directions: 1) adapting Pigeon to consider evolving user preferences; 2) developing efficient strategies to manage lifelong user history for superior personalization; 3) integrating personalized content generation and filtering to construct more powerful personalized content delivery systems.

[50] Hao Yang, Jianxin Yuan, Shuai Yang, Linhe Xu, Shuo Yuan, and Yifan Zeng. 2024. A New Creative Generation Pipeline for Click-Through Rate with Stable Diffusion Model. In *Companion WWW*. ACM, 180–189.

[51] Yukang Yang, Dongnan Gui, Yuhui Yuan, Weicong Liang, Haisong Ding, Han Hu, and Kai Chen. 2024. Glyphcontrol: Glyph conditional control for visual text generation. *NeurIPS* 36 (2024).

[52] Cong Yu, Yang Hu, Yan Chen, and Bing Zeng. 2019. Personalized fashion design. In *ICCV*. IEEE, 9046–9055.

[53] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. 2018. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*. IEEE, 586–595.

[54] Jujia Zhao, Wenjie Wang, Yiyang Xu, Teng Sun, Fuli Feng, and Tat-Seng Chua. 2024. Denoising diffusion recommender model. In *SIGIR*. ACM, 1370–1379.

[55] Zhengbang Zhu, Rongjun Qin, Junjie Huang, Xinyi Dai, Yang Yu, Yong Yu, and Weinan Zhang. 2024. Understanding or Manipulation: Rethinking Online Performance Gains of Modern Recommender Systems. *TOIS* 42, 4 (2024), 1–32.

A APPENDIX

We present five additional Pigeon-generated examples for both sticker and movie poster scenarios, respectively, along with several user-interacted history images and one reference image.

- Sticker scenario.** As illustrated in Figure 6, Pigeon effectively captures users’ visual preferences for character figures and styles in stickers, and combines these preferences with the high-level semantics of the reference sticker to generate personalized stickers. The generated stickers exhibit high semantic alignment with the reference image, including the conveyed emotions, facial expressions, character actions, and elements like hearts.
- Movie poster scenario.** As depicted in Figure 7, each user shows a distinct set of visual preferences, ranging from action and sci-fi to historical drama and crime thrillers. Pigeon-generated posters effectively mirror these preferences through character-centered designs, dynamic compositions, and color palettes that align with each user’s unique taste. By tailoring its designs to the emotional tone, genre, and thematic focus of the reference posters, Pigeon creates personalized posters that strongly resonate with individual users’ past interactions and preferences. For instance, the user in the first row shows a strong preference for action-heavy, explosive films with a focus on dramatic visuals and blockbuster-style presentations. Pigeon matches the user’s love for explosive visuals, with characters taking center stage and environments filled with dynamic elements like fire, destruction, and warfare.

Table 4: Dataset of sticker and movie poster scenarios, where each “sample” consists of user-interacted history images and one reference image.

| | #Users | #Items | #Samples |
|---------------|--------|--------|----------|
| Stickers | 725 | 14,345 | 10,719 |
| Movie posters | 594 | 6,961 | 31,058 |



Figure 6: Examples of generated stickers, along with user-interacted history stickers and one reference sticker.

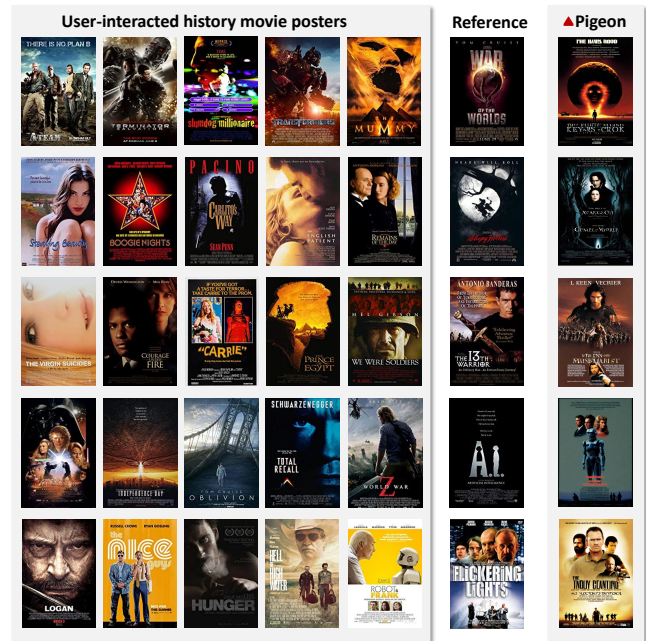


Figure 7: Examples of generated movie posters, along with user-interacted history posters and one reference poster.