

---

# Enhancing Intent Understanding for Ambiguous Prompts: A Human-Machine Co-Adaption Strategy

---

Yangfan He<sup>1,2</sup> Yuxuan Bai<sup>3</sup> Tianyu Shi<sup>4</sup>

## Abstract

Today’s image generation systems are capable of producing realistic and high-quality images. However, user prompts often contain ambiguities, making it difficult for these systems to interpret users’ actual intentions. Consequently, many users must modify their prompts several times to ensure the generated images meet their expectations. While some methods focus on enhancing prompts to make the generated images fit user needs, the model is still hard to understand users’ real needs, especially for non-expert users. In this research, we aim to enhance the visual parameter-tuning process, making the model user-friendly for individuals without specialized knowledge and better understand user needs. We propose a human-machine co-adaption strategy using mutual information between the user’s prompts and the pictures under modification as the optimizing target to make the system better adapt to user needs. We find that an improved model can reduce the necessity for multiple rounds of adjustments. We also collect multi-round dialogue datasets with prompts and images pairs and user intent. Various experiments demonstrate the effectiveness of the proposed method in our proposed dataset.

## 1 Introduction

Generative AI has immense potential to boost economic development by optimizing creative and non-creative tasks. Models like DALL·E 2, IMAGEN, Stable Diffusion, and

Muse can produce unique, convincing, and lifelike images from textual descriptions (Gozalo-Brizuela & Garrido-Merchan, 2023). Despite significant progress, there’s still room for improvement, especially in generating higher-resolution images that better reflect the semantics of input text and in creating more user-friendly interfaces (Frolov et al., 2021). Many models struggle to understand nuanced human instructions, often resulting in a mismatch between user expectations and outputs. Additionally, the impact of variable adjustments on the final image is not always clear, posing challenges for non-expert users who haven’t systematically studied prompt engineering. This complexity hinders those without technical backgrounds from fully utilizing advanced AI models. To address these challenges, we introduce an innovative approach to enhance the user experience for non-professional users. Unlike traditional models that require a deep understanding of underlying mechanisms and control elements, our approach enables users to adjust and optimize image generation with minimal technical knowledge. Inspired by human-in-the-loop co-adaptation (Reddy et al., 2022), our model evolves with user feedback to better meet user expectations. Figure 3 illustrates the operational flow as interacted by users. Our main contributions are:

- **Adaptive Prompt Engineering and Personalized Image Generation:** We propose visual co-adaptation (VCA), an adaptive framework that fine-tunes user prompts using a pre-trained language model enhanced through reinforcement learning, aligning image outputs more closely with user preferences and creating images that truly reflect individual styles and intentions.
- **Human-in-the-Loop Feedback Integration:** Our work considers incorporating human feedback within the training loops of diffusion models. By assessing its impact, we demonstrate how human-in-the-loop methods can surpass traditional reinforcement learning in enhancing model performance and output quality.
- **Comparative Analysis and Tool Development for Non-Experts:** Through comparative analysis, we explore the superiority of mutual information maximization over conventional reinforcement learning in tuning model outputs to user preferences. Additionally, we

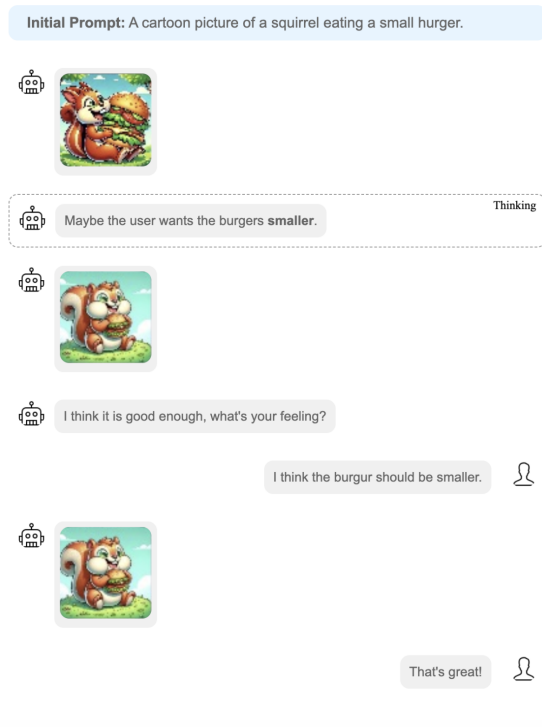
---

<sup>1</sup>University of Minnesota - Twin Cities, Department of Computer Science, Minneapolis, USA <sup>2</sup>Henan RunTai Digital Technology Group Co., Ltd., Zhengzhou, China <sup>3</sup>Southwest Jiaotong University, Department of Computer Science, Chengdu, China <sup>4</sup>University of Toronto, Department of Computer Science, Toronto, Canada. Correspondence to: Yangfan He <he000577@umn.edu>, Yuxuan Bai <18007146970@163.com>, Tianyu Shi <tianyushi3@mail.mcgill.ca>.

*ICML 2024 Workshop on Models of Human Feedback for AI Alignment*, Vienna, Austria. PMLR 235, 2024. Copyright 2024 by the author(s).

# Enhancing Intent Understanding for Ambiguous Prompts: A Human-Machine Co-Adaption Strategy

## Single-round Dialogue - Self Correction



## Multi-round Dialogue - User Correction

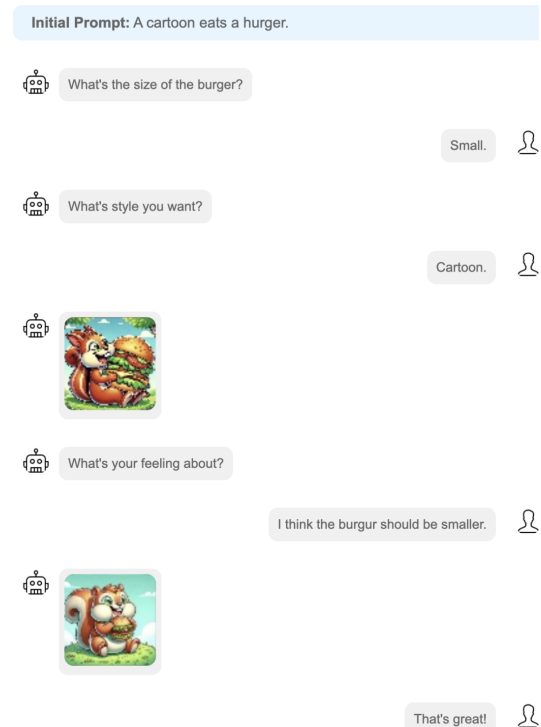


Figure 1. Users have the choice between single-round dialogue, where they provide detailed inputs for the model to generate and self-adjust an image on the left, and multi-round dialogue on the right, where the model engages in iterative refinement based on user feedback, asking questions to clarify any unclear requirements. This allows for either model-driven optimization through self-reflection or user-driven customization to meet specific needs. Our proposed visual co-adaption system can successfully handle both scenarios.

introduce an interactive tool that grants non-experts easy access to advanced generative models, enabling the creation of personalized, high-quality images, thus broadening the applicability of text-to-image technologies in creative domains.

## 2 Related Work

### 2.1 Text-Driven Image Editing Framework:

Image editing is fundamental in computer graphics, with textual prompts providing an intuitive way for users to edit images. Recent advancements in text-to-image generation focus on aligning models with human preferences, using feedback for image refinement. Studies like Hertz et al. (Hertz et al., 2022)'s framework leverage diffusion models' cross-attention layers for high-quality, prompt-driven modifications. Methods like ImageReward (Xu et al., 2024) develop reward models based on human preferences, collecting rich feedback (Wu et al., 2023; Liang et al., 2023) and training models for better image-text alignment and adaptability (Lee et al., 2023) to diverse preferences. However, these methods often lack convenient and efficient editing capabilities. For instance, Hertz et al.'s framework requires

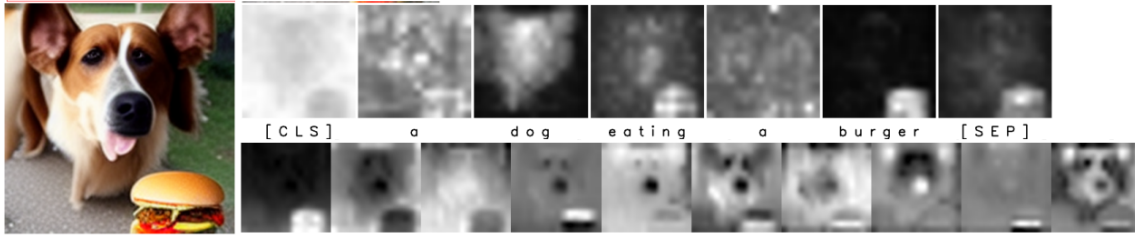
users to adjust complex parameters such as cross-attention and attention weights, demanding high professional knowledge. Unlike traditional image editing, ImageReward (Xu et al., 2024) generates a new image instead of editing the existing one, failing to preserve previous information and risking invalidating prior modifications.

### 2.2 Human Preference-Driven Optimization for Text-to-Image Generation Models:

Zhong et al. (Zhong et al., 2024) significantly advance the adaptability of large language models (LLMs) to human preferences through their innovative approach. Their method leverages advanced mathematical techniques for nuanced, preference-sensitive model adjustments, eliminating the need for exhaustive model retraining. Xu et al. (Xu et al., 2024) adopt a distinctive strategy by harnessing extensive expert insights to develop their ImageReward system, setting a new benchmark for creating images that resonate deeply with human desires. Together, these advancements represent a pivotal shift towards more intuitive, user-centric LLM technologies, heralding a future where AI seamlessly aligns with the intricate mosaic of individual human expectations.

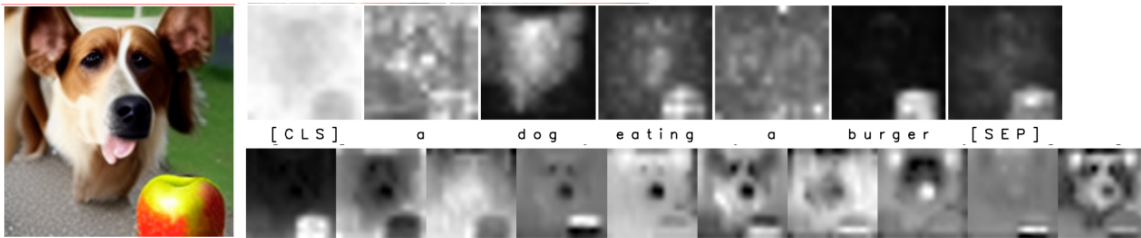
## Enhancing Intent Understanding for Ambiguous Prompts: A Human-Machine Co-Adaption Strategy

Round 1 User: a dog eating a burger



Round 2 User: I change my mind. I want the dog eating an apple instead of burger

Prompt: a dog eating an apple



Round 3 User: I want the photo taken in autumn now

Prompt: a dog eating an apple at autumn

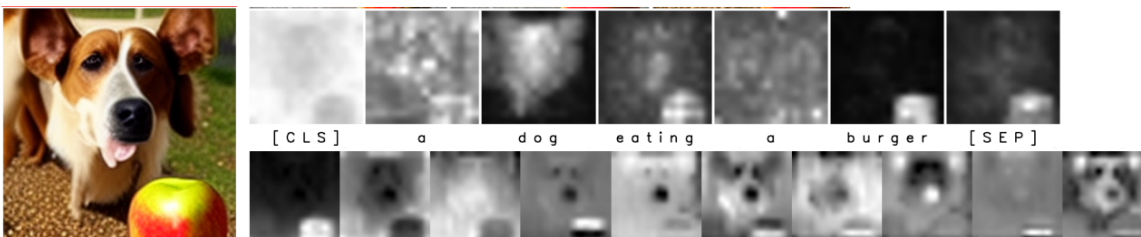


Figure 2. The diagram shows our model’s architecture with cross attention in the first row and self attention in the second. It incorporates an improved cross attention mechanism that maintains shape consistency and aligns well with prompt tokens, enabling effective multi-round modifications based on user feedback. The model captures intricate cross attention details, optimizing parameters for progressively better single-generation performance, demonstrating few-shot learning adaptation with minimal dialogue iterations.

### 2.3 Exploration of Self-Correction Strategies:

Advances in large language models (LLMs) self-correction such as Pan et al (Pan et al., 2023), Shinn et al. (Shinn et al., 2023), Madaan et al (Madaan et al., 2024), improving language understanding and production. Huang et al (Huang et al., 2022) showcased self-debugging and zero-shot learning for reasoning evaluation, underscoring the potential and limits of self-correction. These contributions collectively highlight the progress and future challenges in enhancing LLMs’ self-corrective capabilities (Hertz et al., 2022; Rosenman et al., 2023; Mehrabi et al., 2022; Xu et al., 2024). Meanwhile, we can find that multi-modal self-correction is less investigated. It is also very important to teach the vision model to think it step by step. We explore the integration of self-correction strategies into image generation to produce

images that more closely align with user intentions.

### 2.4 Ambiguity Resolution in Text-to-Image Generation:

Natural dialogue often contains ambiguity due to grammar, polysemy, and vagueness. Humans manage this ambiguity with clarifying questions and contextual cues, but machines find it challenging. To address this, text-to-image generation employs various strategies. For example, masked transformers (Chang et al., 2023) and visual annotations (Endo, 2023) help clarify prompts, while model evaluation benchmarks (Lee et al., 2024) and auto-regressive models (Yu et al., 2022) improve image alignment. Frameworks for abstract (Liao et al., 2023) and inclusive imagery (Zhang et al., 2023), as well as layout guidance (Qu et al., 2023) and

## Enhancing Intent Understanding for Ambiguous Prompts: A Human-Machine Co-Adaption Strategy

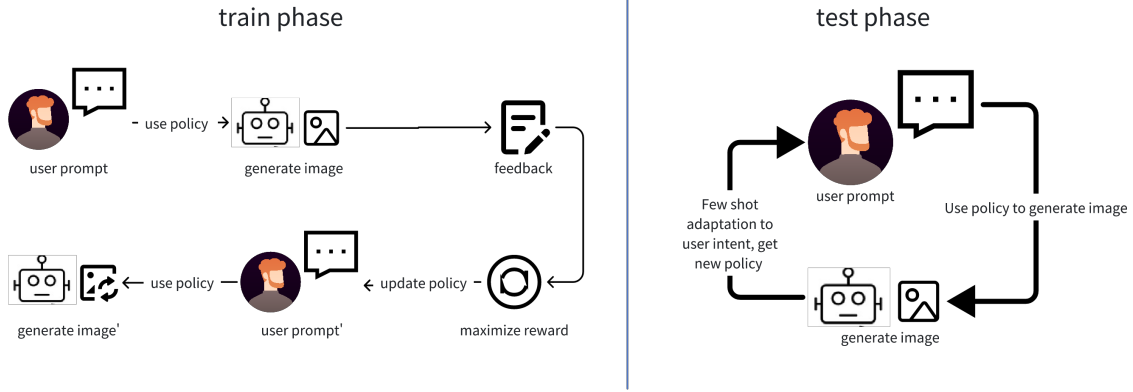


Figure 3. This figure illustrates a reinforcement learning framework with training and testing phases. In training, the policy (three editing operations with trainable parameters, more details in section 3.1.1 and A.3) updates based on human feedback (environment), where the state is the prompt and the action is the generated image. In testing, few-shot adaptation refines the policy ( $\pi_{\text{new}}$ ) to generate images, allowing efficient model adaptation with minimal dialogue interactions.

feedback mechanisms (Liang et al., 2023), further enhance quality. The TIED framework and TAB dataset (Mehrabi et al., 2023) use user interaction to refine prompt clarity. Our model integrates these techniques across multiple dialogue rounds to elicit users’ true intentions, effectively reducing prompt ambiguity and generating results that align with user expectations, thus enhancing image generation quality.

### 3 Method

#### 3.1 Policy Model: Controlling Cross-Attention in a Reinforcement Learning Framework

In our framework, the Imagen text-guided synthesis model (Saharia et al., 2022) constructs the basic composition and geometric layout of images at a  $64 \times 64$  resolution. The model uses a U-shaped network during each diffusion step  $t$  to predict the noise component  $\epsilon$  based on the text embedding  $\psi(P)$  and the noise-added image  $z_t$ . Crucial to shaping the image’s final appearance  $I = z_0$ , the attention maps  $M = \text{Softmax}\left(\frac{QK^T}{\sqrt{d}}\right)$  influence its spatial and geometric properties. Here,  $Q$  and  $K$  are the query and key matrices formed from image and text features, respectively. We define the diffusion step function  $\text{DM}(z_t, P, t, s)$  that computes a single step of the diffusion process, outputting the noisy image  $z_{t-1}$  and the attention map  $M_t$ , if utilized. Overriding the attention map with an additional map  $M_c$  while maintaining the values  $V$  from the prompt is indicated as  $\text{DM}(z_t, P, t, s)\{M \leftarrow M_c\}$ . The modified prompt  $P^*$  generates a new attention map  $M_t^*$ , and the general edit function  $\text{Edit}(M_t, M_t^*, t)$  manages the attention maps at any step  $t$  for both the original and modified images.

#### 3.1.1 EDITING OPERATIONS

In our framework, we employ three strategic editing operations—Word Swap, Adding a New Phrase, and Attention Re-weighting—each optimized through reinforcement learning (RL) as the policy model to enhance the reward function, which is based on the interaction results between the action output in a specific state and the environment (human feedback), using gradient ascent. This approach learns parameters that are highly aligned with human preferences. For more details about the RL training framework, refer to Appendix A.2.

In the **Word Swap** method, users replace tokens in the prompt (e.g., ”a big red bicycle” to ”a big red car”), and we control attention map injection steps to manage compositional freedom:

$$\text{Edit}(M_t, M_t^*, t) := \begin{cases} M_t^* & \text{if } t < \tau \\ M_t & \text{otherwise} \end{cases} \quad (1)$$

The attention map  $M_t^*$  is updated as follows:

$$M_t^* = M_t^* + \eta \nabla_{M_t^*} \mathcal{R}(M_t^*) \quad (2)$$

In the **Adding a New Phrase** method, new tokens are added to the prompt (e.g., ”a castle next to a river” to ”children drawing of a castle next to a river”), targeting shared tokens with an alignment function  $A$ :

$$(\text{Edit}(M_t, M_t^*, t))_{i,j} := \begin{cases} (M_t)_{i,A(j)} & \text{if } A(j) \neq \text{None} \\ (M_t)_{i,j} & \text{otherwise} \end{cases} \quad (3)$$

The alignment function  $A_t$  is updated as follows:

$$A_t = A_t + \eta \nabla_{A_t} \mathcal{R}(A_t) \quad (4)$$

In the **Attention Re-weighting** method, token influence is adjusted to enhance or diminish features (e.g., scaling the attention map of "fluffy red ball" for token  $j^*$  with a parameter  $c \in [-2, 2]$ ):

$$(\text{Edit}(M_t, M_t^*, t))_{i,j} := \begin{cases} c \cdot (M_t)_{i,j} & \text{if } j = j^* \\ (M_t)_{i,j} & \text{otherwise} \end{cases} \quad (5)$$

This parameter  $c$  provides intuitive control over the induced effect. The scaling parameter  $c_t$  is updated as follows:

$$c_t = c_t + \eta \nabla_{c_t} \mathcal{R}(c_t) \quad (6)$$

Each operation refines text-image interactions through cross-attention layers, aligning outputs with human preferences. The RL framework optimizes these strategies by updating  $M_t$ ,  $A_t$ , and  $c_t$  through gradient ascent. For detailed optimization processes of the three editing operations, see Appendix A.3.

### 3.2 Human-Machine Co-Adaptation with Mutual Information

In this section, we explain how our model can adapt to human intent. Let  $X$  denote the user inputs and  $Y$  the images generated by the model. The adaptation mechanism seeks to maximize the mutual information  $I(X; Y)$ , which quantifies the amount of information shared between  $X$  and  $Y$ . The mutual information is given by:

$$I(X; Y) = \int_{x \in X} \int_{y \in Y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)} dy dx, \quad (7)$$

where  $p(x, y)$  is the joint probability distribution of  $X$  and  $Y$ , and  $p(x)$  and  $p(y)$  are the marginal distributions of  $X$  and  $Y$ , respectively.

#### Adaptive Feedback Loop

The adaptive feedback loop updates the model parameters  $\theta$  to better align with human intent, utilizing the gradient of mutual information that is now conditioned on user feedback  $f$ . This feedback directly represents human preferences and intents, guiding the model towards desired outcomes:

$$\theta_{\text{new}} = \theta_{\text{old}} + \eta \nabla_{\theta} I(X; Y | f), \quad (8)$$

where  $\eta$  is the learning rate and  $f$  encapsulates the feedback signals from users. This adaptive approach measures effectiveness through an increase in conditional mutual information, reflecting improved alignment with user expectations, and higher user satisfaction scores in image generation tasks.

---

#### Algorithm 1 Prompt-to-Prompt Image Editing with Human-Machine Co-Adaptation (Training)

---

**Input:** Original prompt  $P_0$ , Edited prompt  $P_1$ , Initial image  $I_0$   
**Output:** Edited image  $I_1$

- 1: Initialize interface  $\pi$  with parameters  $\theta$
- 2: Generate initial attention maps  $A_0$  for  $I_0$  using  $\pi(P_0)$
- 3: Set  $I_t \leftarrow I_0$
- 4: Initialize user feedback loop
- 5: **for**  $t = 1$  **to** Convergence **do**
- 6:   Collect user feedback on image  $I_t$  and prompt  $P_t$
- 7:   Adapt  $\pi$  (Using editing operation in Section 3.1.1) to maximize mutual information  $I(A; I|P)$  incorporating feedback
- 8:   Apply  $P_1$  to generate new attention maps  $A_1$
- 9:   Generate  $I_1$  by applying  $A_1$  in diffusion step
- 10:   Evaluate  $I(A; I|P)$  between  $(P_0, P_1)$  and  $(I_0, I_1)$
- 11:   Update  $\theta$  to align more closely with user preferences
- 12: **end for**
- 13: Conduct final evaluation of  $I_1$  with user
- 13: **return**  $I_1 = 0$

---



---

#### Algorithm 2 Evaluation of Adaptation to New User Preferences

---

**Input:** Trained interface  $\pi$  with parameters  $\theta$ , New user initial prompt  $P_{\text{new}}$   
**Output:** Adapted image  $I_{\text{adapted}}$  aligns with new user preferences

- 1: Initialize new user interaction session
- 2: **for**  $i = 1$  **to** few-shot rounds **do**
- 3:   Present  $I_{\text{current}}$  generated from  $P_{\text{new}}$  using  $\pi$
- 4:   Collect new user feedback on  $I_{\text{current}}$
- 5:   Update  $P_{\text{new}}$  based on user feedback
- 6:   Adapt pre-trained  $\theta$  minimally to reflect new user preferences
- 7:   Generate new  $I_{\text{current}}$  using updated  $\pi(P_{\text{new}})$
- 8:   **if** user feedback is positive **then**
- 9:     Break the loop and finalize  $I_{\text{adapted}}$
- 10:   **end if**
- 11: **end for**
- 12: Evaluate user satisfaction with  $I_{\text{adapted}}$
- 12: **return**  $I_{\text{adapted}} = 0$

---

## 4 Experiments

### 4.1 Settings

The experiments are conducted using 4 NVIDIA 4090 GPUs, This setup allows us to utilize complex algorithms such as diverse beam search with a beam size of 8 and a diversity penalty of 1.0, ensuring thorough exploration and diversity in the generated responses. The model parameters are initialized from a fine-tuned baseline, which provides a robust starting point for further optimization. Over three days of training session, which encompass 12,000 episodes, with four PPO epochs per batch and a batch size of 256. The learning rate is set at  $5 \times 10^{-5}$ , and the value and KL reward coefficients are meticulously calibrated to 2.2 and 0.3, respectively, to balance the learning dynamics. For additional details due to page constraints, see Appendix A.1.

4.2 Dataset

We have developed a Q&A software that annotates prompts on our platform, automatically generating JSON files that capture detailed multi-turn dialogue information. An example of user interface annotations is showcased in the Appendix A.4. Our training set includes 1673 meticulously crafted JSON files, each annotated with prompts, detailed Q&A sequences, image paths, unique identifiers, and ratings for image alignment and fidelity. This dataset instructs our model on user expectations and artistic intentions, analyzing subjects, emotions, settings, styles, perspectives, and extra elements. Feedback is synthesized into refined prompts, enabling the model to grasp complex artistic directions crucial for user resonance. We use 95% of the data for training and 5% for validation, supporting efficient few-shot learning to enhance both performance and user satisfaction.

4.3 Evaluation Metrics

The experimental framework of this study is meticulously designed to evaluate our text-to-image generation model across four key dimensions.

**Lpips** (Zhang et al., 2018): is a deep learning metric that evaluates how image modifications preserve the original structure, with lower scores indicating minimal visual differences and alignment with human perception. It measures the consistency and perceptual coherence of images generated in successive dialogue rounds.

**Clip Score** (Radford et al., 2021): Based on the CLIP model, the system evaluates image-text alignment, assigning scores from 0 (no similarity) to 1 (perfect alignment). In dialogues, the LLM subtly adjusts prompts and selects one of three strategies following user feedback. The text-to-image model, using reinforcement learning and CLIPScore, iteratively refines images until reaching a satisfactory score. For detailed information on how the ChatGPT-4 modifies prompts based on human input, refer to the Appendix A.8.

**Human Evaluation:** In a study with 100 diverse users, we utilize a randomized control trial with stratified sampling based on age, gender, and technical proficiency. Using a blind design, participants are unaware of the models or components being tested to prevent biases. Detailed feedback is collected through electronic surveys post-interaction, utilizing standardized forms with scaled and open-ended questions. A cross-over design ensures that each user experiences all model variations in a randomized order, maximizing exposure. Statistical power analysis confirms that 100 participants provide sufficient power to detect significant results.

4.4 Comparison Study

4.4.1 TRENDS ACROSS BASELINES OVER ITERATIVE ROUNDS

Figure 4.4.1 showcases our model’s superior performance on a validation prompt describing ”A serene ancient fantasy sanctuary constructed of stone, with white birds flying in the

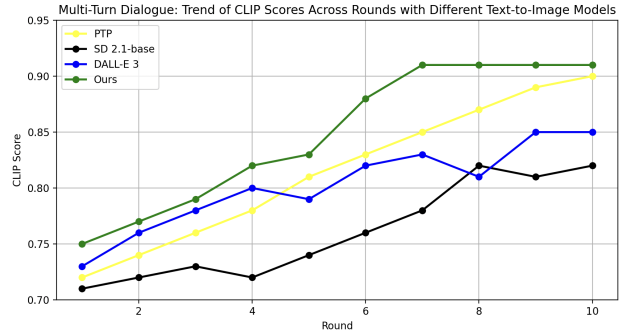


Figure 4. This graph shows CLIP score trends over 10 rounds for various text-to-image models (PTP (Hertz et al., 2022), SD 2.1-base, DALL-E 3, and ours)

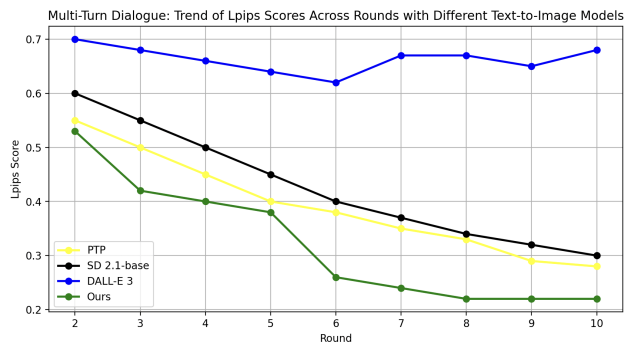


Figure 5. Illustrated in the graph are the trends of LPIPS scores for several text-to-image models (PTP, SD 2.1-base, DALL-E 3, and ours) over 10 rounds.

distance.” and achieves high CLIP scores early, our model reaches 0.78 by round 3 and peaks at 0.91 by round 7, surpassing competitors. It also excels in Lpips, as is shown in Figure 5 recording a score of 0.42 by round 3 and stabilizing at 0.22 by round 8. This rapid stabilization highlights our model’s adaptability and efficiency, maintaining high consistency and user satisfaction across fewer dialogue rounds. Each round incrementally builds on the last, refining details without altering the prompt’s core structure.

4.4.2 PROMPT REFINEMENT

Table 1 provides a detailed comparison between **self-reflection prompt refinement** and **multi-round dialogue prompt refinement**. Self-reflection is notably quicker (3.4s vs. 12s), yet multi-round dialogue more effectively captures user preferences, resulting in higher satisfaction ratings (4.7 vs. 3.0). Additionally, it demonstrates a significant improvement in Purpose Adaptability (4.8 vs. 3.3) along with modest enhancements in Clarity (4.7 vs. 4.2) and Detail Level (4.2 vs. 4.1). For a deeper exploration of the algorithms behind these refinement methods, please refer to Appendix A.7.

## Enhancing Intent Understanding for Ambiguous Prompts: A Human-Machine Co-Adaption Strategy

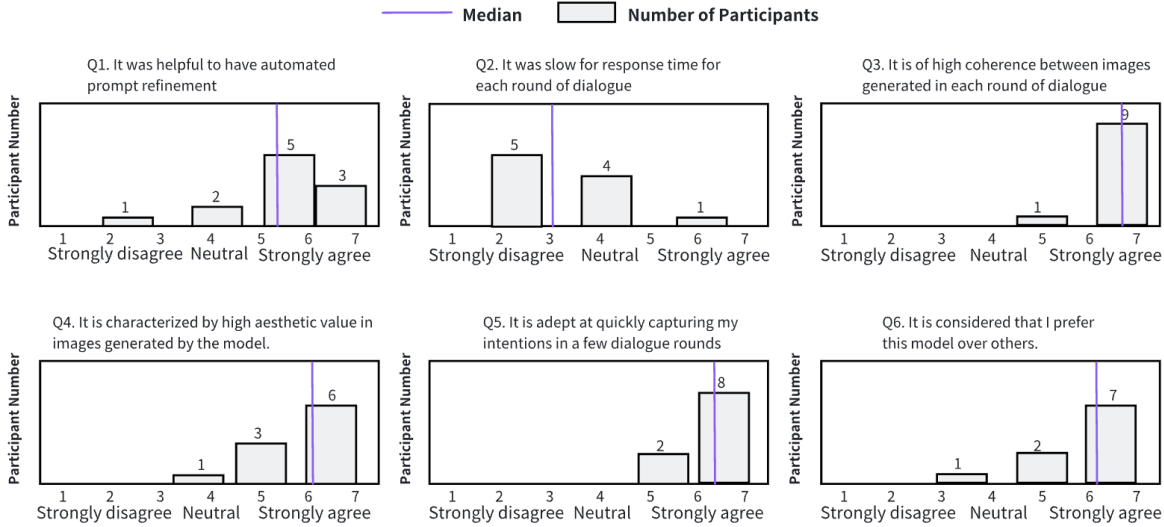


Figure 6. The chart shows user feedback on a model, highlighting mixed responses with positive feedback on image coherence and capturing intentions, but concerns over response time.

Table 1. Comparative Analysis of Prompt Refinement from 100 users, averaged and rounded to one decimal. Metrics are scored on a 0-5 scale. Response Time indicates average duration for self-reflection and multi-dialogue processes.

Metric & Category	Refine Type	
	Self-reflection	Multi-dialogue
<b>Prompt Quality</b>		
Clarity	4.2/5	<b>4.7/5</b>
Detail Level	4.1/5	<b>4.2/5</b>
Purpose Adaptability	3.3/5	<b>4.8/5</b>
<b>Image Reception</b>		
User Satisfaction	3.0/5	<b>4.7/5</b>
Clip Value	0.8/1	<b>0.9/1</b>
<b>Response Time</b>	<b>3.4s</b>	12s

### 4.5 Ablation Study: Reinforcement Learning for Parameter Tuning

Table 2 demonstrates the substantial impact of Reinforcement Learning (RL) tuning on dialogue system performance. Systems equipped with RL require significantly fewer dialogue rounds, averaging 4.3 compared to 6.9 for those without RL, highlighting enhanced efficiency in responding to user inputs. Additionally, RL tuning improves the CLIP score from 0.83 to 0.92, indicating better alignment of generated images with textual prompts. User satisfaction also increases markedly with RL, from 4.14 to 4.73 out of 5, reflecting a more pleasing user experience. While both systems perform similarly in aesthetic quality (4.89 vs. 4.88), the primary benefits of RL tuning are seen in functionality and user satisfaction. Users, unaware of the tuning status during tests, noted lower consistency in image quality from

Table 2. Compares RL effects using data averaged from random 10 of 100 users, with final interaction CLIP and Aesthetic Scores.

Metrics	With RL	Without RL
<b>Rounds</b>	<b>4.3</b>	6.9
<b>CLIP Score</b>	<b>0.92/1.0</b>	0.83/1.0
<b>User Satisfaction</b>	<b>4.73/5</b>	4.14/5
<b>Aesthetic Score</b>	<b>4.89/5</b>	4.88/5

Table 3. Assesses cross attention(CA)’s impact, averaging data from random 10 of 100 users, with CLIP and Aesthetic Scores from the final interaction.

Metrics	Edited CA	Normal CA
<b>Rounds</b>	<b>3.7</b>	6.1
<b>CLIP Score</b>	<b>0.88/1.0</b>	0.81/1.0
<b>User Satisfaction</b>	<b>4.82/5</b>	3.94/5
<b>Aesthetic Score</b>	<b>4.71/5</b>	4.48/5

the non-RL-tuned model. This underscores the effectiveness of RL in adapting dynamically to user feedback, leading to quicker, more relevant, and satisfying interactions. For a detailed discussion on the parameter updates facilitated by RL tuning, refer to Appendix A.5.

### 4.6 Ablation Study: Comparing Edited Cross Attention with Normal Cross Attention.

Table 3 highlights the superior performance of edited cross attention (CA) over normal CA in dialogue systems, showcasing their distinct approaches to adaptability. Normal CA computes attention weights based on initial inputs and maintains them statically throughout the interaction, whereas edited CA dynamically adjusts these weights in response

to changes in dialogue context and user feedback. This adaptability significantly reduces dialogue rounds, averaging 3.7 compared to 6.1 for normal CA, and leads to notable enhancements in system performance. For instance, edited CA achieves a higher CLIP score of 0.88 versus 0.81 and increases user satisfaction from 3.94 to 4.82 out of 5. The aesthetic quality of images also improves with edited CA, scoring 4.71 compared to 4.48 for normal CA. These results underscore the effectiveness of integrating reinforcement learning with edited CA to refine the tuning process and improve the consistency and relevance of outputs in denoising tasks. For an in-depth exploration of how edited cross attention mechanisms function within the system, refer to Appendix A.6.

4.7 Visualization Results

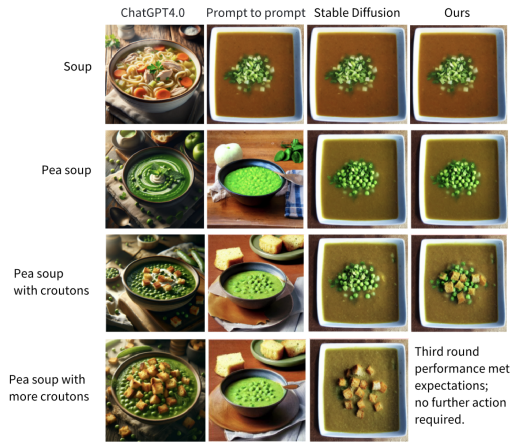


Figure 7. The comparison demonstrates our model’s few-shot learning capability, effectively adapting to user preferences with minimal dialogue.

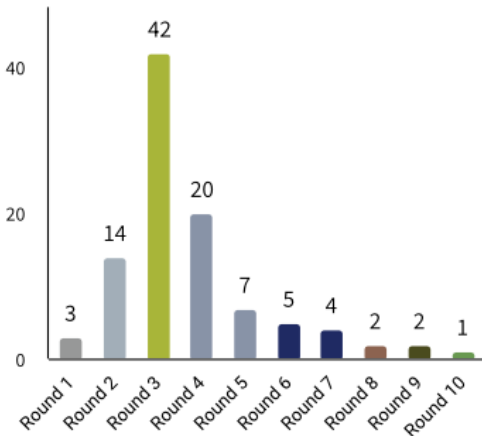


Figure 8. The chart shows the rapid decline in user interaction rounds needed for satisfaction, peaking by Round 5, demonstrating the model’s efficient few-shot learning.

Dialogue Rounds Across Different Models

Figure 2 compares dialogue rounds across different models: ChatGPT, Stable Diffusion v2.1, Prompt-to-Prompt (Hertz et al., 2022), and our model. Initially, images from Stable Diffusion, Prompt-to-Prompt, and our model are similar due to the lack of feedback. By the second round, ”pea soup” preferences cause significant changes in ChatGPT-4 and Stable Diffusion, affecting consistency. In the third round, with croutons added, our model excels by fine-tuning parameters via reinforcement learning, maintaining balance, while Prompt-to-Prompt struggles, and ChatGPT-4 shows inconsistencies. By the fourth round, our model achieves satisfactory results and opts out, while the others continue ineffective adjustments. This highlights our model’s superior ability to understand and respond to user feedback, achieving optimal results by the third round and demonstrating effective multi-round dialogue learning. Despite ChatGPT-4’s realistic visuals, it struggles with consistency and adapting to human preferences. Our model, preferred by 89% of users, effectively adapts with minimal dialogue.

User Satisfaction Distribution for Our Model Over Multiple Rounds

Figure 8 illustrates our model’s efficiency in adapting to user feedback. Initially, the satisfaction rate increases rapidly, with 59 users satisfied by Round 3, demonstrating the model’s quick alignment with user preferences. By Round 5, satisfaction peaks at 99 out of 100 users, underscoring the model’s effectiveness in achieving high user satisfaction swiftly.

Users’ Overall Evaluation of Our Model

Figure 6 presents user evaluations across various model aspects. The majority found the automated prompt refinement to be helpful, indicating approval. In contrast to typical concerns about speed in models with complex computations, most users disagreed with the notion that the model’s response time per dialogue round was slow, suggesting that the integration of reinforcement learning for fine-tuning did not significantly impact perceived efficiency. The model was highly praised for its coherence across images generated in each dialogue round and received commendations for aesthetic quality. It was also recognized for adeptly capturing user intentions within just a few rounds of dialogue. Overall, the participants showed a strong preference for this model over others, reflecting its effectiveness and user satisfaction.

5 Conclusion and Future Work

In this study, we introduced a new image generation method using a human-in-the-loop approach that enhances user interaction and responsiveness to ambiguous prompts. Our findings highlight the model’s ability to closely match user



expectations through adaptive prompt engineering and mutual information optimization. Looking ahead, we plan to release our training dataset, improving transparency and enabling broader testing. Additionally, we aim to refine the model’s interpretive skills, expand its applications across different domains, and conduct comprehensive benchmarks to gauge the alignment between user intentions and generated images. These initiatives will advance personalized and intuitive image generation technologies, making advanced modeling tools more accessible without requiring deep technical expertise.

### References

- Chang, H., Zhang, H., Barber, J., Maschinot, A., Lezama, J., Jiang, L., Yang, M.-H., Murphy, K., Freeman, W. T., Rubinstein, M., et al. Muse: Text-to-image generation via masked generative transformers. *arXiv preprint arXiv:2301.00704*, 2023.
- Endo, Y. Masked-attention diffusion guidance for spatially controlling text-to-image generation. *The Visual Computer*, pp. 1–13, 2023.
- Frolov, S., Hinz, T., Raue, F., Hees, J., and Dengel, A. Adversarial text-to-image synthesis: A review. *Neural Networks*, 144:187–209, 2021.
- Gozalo-Brizuela, R. and Garrido-Merchan, E. C. Chatgpt is not all you need. a state of the art review of large generative ai models. *arXiv preprint arXiv:2301.04655*, 2023.
- Hertz, A., Mokady, R., Tenenbaum, J., Aberman, K., Pritch, Y., and Cohen-Or, D. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022.
- Huang, J., Gu, S. S., Hou, L., Wu, Y., Wang, X., Yu, H., and Han, J. Large language models can self-improve. *arXiv preprint arXiv:2210.11610*, 2022.
- Lee, K., Liu, H., Ryu, M., Watkins, O., Du, Y., Boutilier, C., Abbeel, P., Ghavamzadeh, M., and Gu, S. S. Aligning text-to-image models using human feedback. *arXiv preprint arXiv:2302.12192*, 2023.
- Lee, T., Yasunaga, M., Meng, C., Mai, Y., Park, J. S., Gupta, A., Zhang, Y., Narayanan, D., Teufel, H., Bellagente, M., et al. Holistic evaluation of text-to-image models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Liang, Y., He, J., Li, G., Li, P., Klimovskiy, A., Carolan, N., Sun, J., Pont-Tuset, J., Young, S., Yang, F., et al. Rich human feedback for text-to-image generation. *arXiv preprint arXiv:2312.10240*, 2023.
- Liao, J., Chen, X., Fu, Q., Du, L., He, X., Wang, X., Han, S., and Zhang, D. Text-to-image generation for abstract concepts. *arXiv preprint arXiv:2309.14623*, 2023.
- Madaan, A., Tandon, N., Gupta, P., Hallinan, S., Gao, L., Wiegrefe, S., Alon, U., Dziri, N., Prabhunoye, S., Yang, Y., et al. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36, 2024.
- Mehrabi, N., Goyal, P., Verma, A., Dhamala, J., Kumar, V., Hu, Q., Chang, K.-W., Zemel, R., Galstyan, A., and Gupta, R. Is the elephant flying? resolving ambiguities in text-to-image generative models. *arXiv preprint arXiv:2211.12503*, 2022.
- Mehrabi, N., Goyal, P., Verma, A., Dhamala, J., Kumar, V., Hu, Q., Chang, K.-W., Zemel, R., Galstyan, A., and Gupta, R. Resolving ambiguities in text-to-image generative models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 14367–14388, 2023.
- Pan, L., Saxon, M., Xu, W., Nathani, D., Wang, X., and Wang, W. Y. Automatically correcting large language models: Surveying the landscape of diverse self-correction strategies. *arXiv preprint arXiv:2308.03188*, 2023.
- Qu, L., Wu, S., Fei, H., Nie, L., and Chua, T.-S. Layoutllm-t2i: Eliciting layout guidance from llm for text-to-image generation. In *Proceedings of the 31st ACM International Conference on Multimedia*, pp. 643–654, 2023.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Reddy, S., Levine, S., and Dragan, A. First contact: Unsupervised human-machine co-adaptation via mutual information maximization. *Advances in Neural Information Processing Systems*, 35:31542–31556, 2022.
- Rosenman, S., Lal, V., and Howard, P. Neuroprompts: An adaptive framework to optimize prompts for text-to-image generation. *arXiv preprint arXiv:2311.12229*, 2023.
- Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E. L., Ghasemipour, K., Gontijo Lopes, R., Karagol Ayan, B., Salimans, T., et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35: 36479–36494, 2022.

Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

Shinn, N., Cassano, F., Labash, B., Gopinath, A., Narasimhan, K., and Yao, S. Reflexion: Language agents with verbal reinforcement learning.(2023). *arXiv preprint cs.AI/2303.11366*, 2023.

Vijayakumar, A. K., Cogswell, M., Selvaraju, R. R., Sun, Q., Lee, S., Crandall, D., and Batra, D. Diverse beam search: Decoding diverse solutions from neural sequence models. *arXiv preprint arXiv:1610.02424*, 2016.

Wu, X., Sun, K., Zhu, F., Zhao, R., and Li, H. Better aligning text-to-image models with human preference. *arXiv preprint arXiv:2303.14420*, 2023.

Xu, J., Liu, X., Wu, Y., Tong, Y., Li, Q., Ding, M., Tang, J., and Dong, Y. Imagereward: Learning and evaluating human preferences for text-to-image generation. *Advances in Neural Information Processing Systems*, 36, 2024.

Yu, J., Xu, Y., Koh, J. Y., Luong, T., Baid, G., Wang, Z., Vasudevan, V., Ku, A., Yang, Y., Ayan, B. K., et al. Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2(3):5, 2022.

Zhang, C., Chen, X., Chai, S., Wu, C. H., Lagun, D., Beeler, T., and De la Torre, F. Iti-gen: Inclusive text-to-image generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3969–3980, 2023.

Zhang, R., Isola, P., Efros, A. A., Shechtman, E., and Wang, O. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 586–595, 2018.

Zhong, Y., Ma, C., Zhang, X., Yang, Z., Zhang, Q., Qi, S., and Yang, Y. Panacea: Pareto alignment via preference adaptation for llms. *arXiv preprint arXiv:2402.02030*, 2024.

## A Appendix

### A.1 Reinforcement Learning configuration

To train our policy model, we employ Proximal Policy Optimization (PPO) (Schulman et al., 2017), initializing the value and policy networks from a supervised fine-tuned model. We use diverse beam search (Vijayakumar et al., 2016) with a beam size of 8 and a diversity penalty of 1.0 to ensure exploration quality and diversity. The maximum generation length is randomly set between 15 to 75 at each step, and one completion is randomly selected to update the

policy. Each prompt generates one image, computing the clip score as the reward function to reduce variance. Training involves 12,000 episodes, four PPO epochs per batch, a batch size of 256, and a learning rate of  $5e-5$ , with value and KL reward coefficients set at 2.2 and 0.3, respectively. Based on human fragmented language feedback, ChatGPT provides new prompts with minimal structural changes but reflects human intent very well.

### A.2 Reinforcement Learning Framework

The reinforcement learning framework for our human-machine co-adaptation system in image editing involves the following elements:

#### STATE (S)

The state in our framework represents the current situation of the system, which includes:

- The current image  $I_t$  being edited.
- The current prompt  $P_t$  describing desired modifications or features in the image.
- Optionally, it can also include historical user interactions and feedback to provide context to the state, enabling the model to better understand and predict user preferences.

#### ACTION (A)

Actions in this context refer to the modifications applied to the image based on the input prompt and model’s interpretation:

- Adjustments or transformations applied to the image  $I_t$  to generate a new image  $I_{t+1}$ .
- These actions are driven by the interpretation of the user’s prompt, potentially influenced by machine learning algorithms that predict optimal changes.

#### REWARD (R)

The reward function is crucial as it guides the training of the RL model by quantifying the success of actions taken based on the state:

- It could be defined using objective metrics such as the similarity between the generated image and user’s expected outcome, measured by tools like CLIP score.
- Feedback from users after viewing the modified image can also be used as part of the reward, where positive feedback increases the reward and negative feedback decreases it.

## Enhancing Intent Understanding for Ambiguous Prompts: A Human-Machine Co-Adaption Strategy

- The reward aims to maximize the alignment between the user’s intent and the image output, effectively training the model to interpret and act upon ambiguous prompts accurately.

This reinforcement learning setup enables our system to iteratively learn and adapt from each user interaction, improving its ability to decode ambiguous prompts and align image outputs with user expectations.

### A.3 Optimization Details

To optimize image generation, the model dynamically selects among three strategies (adding phrases, word swapping, re-weighting) using the CLIP score as the reward function to update all parameters of the chosen strategy. This feedback-driven approach optimizes parameters within one strategy per iteration, yielding three well-adjusted parameter sets that adapt image generation to human preferences. The strategies correspond to three controllers: Attention-Replace, Attention-Refine, and Attention-Reweight. Our text-to-image model uses controllers to adjust cross-attention during generation, with each controller utilizing cross-attention information between images and prompts in each dialogue round. The controllers correspond to three strategies with trainable parameters, including the dynamic proportion of self-attention during the sampling process, the proportion of attention injection steps, and adaptive updates to cross-attention maps based on dialogue feedback. The optimization process for parameter updates can be mathematically represented as follows:

#### Reward function:

This is computational framework for the reward function  $\mathcal{R}(\theta)$  in a reinforcement learning context, where the CLIP score assesses the similarity between generated images and textual prompts. Specifically:

$$\mathcal{R}(\theta) = \text{CLIPScore}(I_{\text{gen}}, P_{\text{prev}}) + \lambda \cdot \text{CLIPScore}(I_{\text{gen}}, P_{\text{new}}) \quad (9)$$

This formula ensures that the parameters are finely tuned, with  $\lambda$  serving as a balancing factor between aligning the generated image with the previous prompt and the new prompt, fostering both continuity and responsiveness to new requirements. Extensive experimentation has determined that setting  $\lambda = 0.2$  is optimal, as it allows the CLIP score to converge more rapidly to its maximum value. When incrementally increasing  $\lambda$  from 0.1 to 1, the performance peaks at 0.2. However, increasing  $\lambda$  beyond 1 leads to a significant decline in performance, falling even below the levels observed at  $\lambda = 0.1$ . Further, to underscore the iterative update mechanism integral to the reinforcement

learning cycle:

$$I_{\text{gen}}^{(k+1)} = \text{Update}(I_{\text{gen}}^{(k)}, \theta^{(k)})$$

Here,  $I_{\text{gen}}^{(k)}$  signifies the image generated at iteration  $k$ , and  $\theta^{(k)}$  indicates the parameters at that iteration. The update function modifies the image based on the current parameters, capturing the dynamic nature of the learning process across successive rounds.

#### Attention-Replace Strategies:

Update method directly adjusts the mapping matrix  $M$  using gradient ascent and then multiplies it with the cross-attention matrix  $M_{\text{cross.attention}}$  called `mapper` to alter the attention distribution, impacting the generated image’s features and quality.

$$M_{\text{new}} = (M + \eta \cdot \Delta M) \cdot M_{\text{cross.attention}} \quad (10)$$

#### Attention-Refine Strategies:

Update the attention weights by combining the original and new attention maps derived from the modified prompt. In the `Attention-Refine` class, the `mapper` aligns base attention weights with the new prompt structure while `alphas` blend original and modified weights, ensuring the final output accurately reflects user modifications and maintains consistency. The `mapper` tensor aligns tokens between prompts, enabling correct transfer of attention weights; updated as

$$\theta'_m = \theta_m + \eta \nabla_{\theta_m} \mathbb{E}[R]$$

to maximize the expected reward ( $\mathbb{E}[R]$ ) using gradient ascent with learning rate  $\eta$ . The `alphas` weights control the blending of original and modified attention weights, determining each token’s influence; updated as

$$\theta'_\alpha = \theta_\alpha + \eta \nabla_{\theta_\alpha} \mathbb{E}[R]$$

to maximize the expected reward ( $\mathbb{E}[R]$ ) using gradient ascent with learning rate  $\eta$ .

The attention weights are updated by combining the original and new attention maps derived from the modified prompt. The original attention is processed using the `mapper`, which aligns the attention weights by permuting dimensions based on the mapped indices:

$$\begin{aligned} \text{attn\_base\_replace}_{ijk} &= \text{attn\_base}_{ijk} \cdot \text{mapper}_{kj} \\ \implies (\text{attn\_base\_replace})_{\text{permute}(2,0,1,3)} \end{aligned}$$

Here, `mapperkj` indicates the mapping from index  $k$  in the original prompt to index  $j$  in the new prompt. The operation

## Enhancing Intent Understanding for Ambiguous Prompts: A Human-Machine Co-Adaption Strategy

$(\text{attn\_base\_replace})_{\text{permute}(2,0,1,3)}$  permutes the dimensions of the resulting tensor to align with the expected structure for further processing.

The updated attention weights are then calculated as:

$$M_{\text{update}}^{(t)} = \beta_t \cdot M_{\text{orig}}^{(t)} + (1 - \beta_t) \cdot M_{\text{new}}^{(t)}$$

### Attention-Reweight Strategies:

Modifies the distribution of attention by first blending the original and new attention maps, and then scaling the weights according to user preferences. The blending of attention maps is given by:

$$M_{\text{refine}}^{(t)} = \beta_t \cdot M_{\text{orig}}^{(t)} + (1 - \beta_t) \cdot M_{\text{new}}^{(t)}, \quad \beta_t = \beta_{t-1} + \gamma \cdot \nabla_{\beta_t} \mathcal{R}(\theta) \quad (11)$$

with  $\beta_t$  adjusting the blending ratio dynamically based on feedback, and  $\gamma$  is the learning rate for  $\beta_t$ . After blending, the attention distribution is further modified by scaling the weights:

$$M_{\text{reweight}}^{(t)} = \sum_i \gamma_{t,i} \cdot M_{\text{refine}}^{(t,i)}, \quad \gamma_{t,i} = \gamma_{t-1,i} + \kappa \cdot \nabla_{\gamma_{t,i}} \mathcal{R}(\theta) \quad (12)$$

where  $\gamma_{t,i}$  are the weight multipliers that adapt the emphasis on specific features, and  $\kappa$  is the learning rate for  $\gamma_{t,i}$ .

In addition to these, we also update the proportions related to specific attention mechanisms:

$$\alpha_{t+1} = \alpha_t + \eta \nabla_{\alpha_t} \mathcal{R}(\theta) \quad (13)$$

$$\zeta_{t+1} = \zeta_t + \gamma \nabla_{\zeta_t} \mathcal{R}(\theta) \quad (14)$$

$$\delta_{t+1} = \delta_t + \kappa \nabla_{\delta_t} \mathcal{R}(\theta) \quad (15)$$

Here,  $\alpha$  represents the proportion of self-attention features injected at different stages of the sampling process,  $\zeta$  represents the replacement proportion of the cross-attention map, and  $\delta$  represents the overall number of sampling steps.

### A.4 Q&A Software Annotation Interface

### A.5 Ablation of RL tuning

The RL tuning process and static parameter configuration are mathematically represented as:

$$\theta^{\text{RL}} = \theta_0 + \sum_{t=1}^T \eta \nabla_{\theta} \mathcal{R}(\theta_t), \quad \theta^{\text{Fixed}} = \theta_0 \quad (16)$$

Here,  $\theta^{\text{RL}}$  are the parameters iteratively updated with RL,  $\theta_0$  is the initial parameter setting,  $\eta$  is the learning rate, and  $\nabla_{\theta} \mathcal{R}(\theta_t)$  is the gradient of the reward function at iteration  $t$ . This setup without RL results in more dialogue rounds and less optimal outcomes.

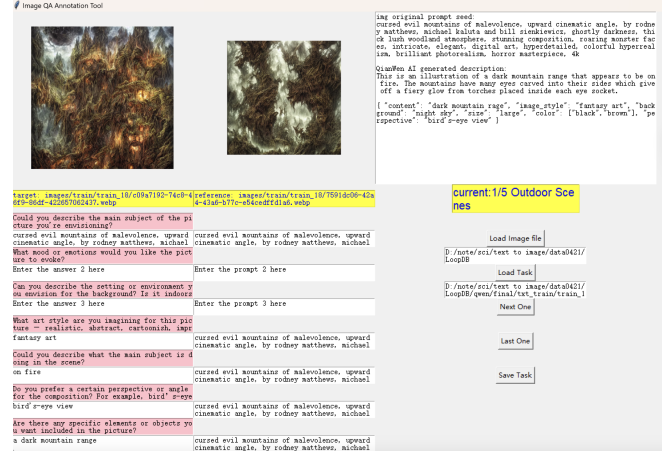


Figure 9. Screenshot of the Q&A software annotation interface.

### A.6 Ablation of cross attention control

$$\theta_{\text{Weighted}}^{(t+1)} = \theta_{\text{Weighted}}^{(t)} + \eta \nabla_{\theta} \mathcal{L}(I_t, \text{Feedback}_t, M) \quad (17)$$

$$\theta_{\text{Empty}}^{(t+1)} = \theta_{\text{Empty}}^{(t)} + \eta \nabla_{\theta} \mathcal{L}(I_t, \text{Feedback}_t, M_{\text{new}}) \quad (18)$$

This setup employs only new attention without blending it with the base cross attention. Each strategy involves a distinct function to modify the cross attention map, directed by its corresponding controller. For standard cross attention, the controller is set to 'empty control' within the code.

### A.7 LLM Prompt Refinement

#### Algorithm 3 Multi-dialogue Prompt Refine Process for ChatGPT-4

```

0: Input: Initial prompt  $p_0$ 
0: Output: Refined prompt  $p_i$  that meets conditions and is ambiguity-free
0: Define  $C(p)$ : Checks if prompt  $p$  meets all predefined conditions.
0: Define  $A(p)$ : Checks if prompt  $p$  is free of ambiguities.
0:  $i \leftarrow 0$ 
0: while  $\neg C(p_i) \vee \neg A(p_i)$  do
0:   if  $\neg A(p_i)$  then
0:      $p_{i+1} \leftarrow \text{ResolveAmbiguities}(p_i)$  {Clarify prompt, ensuring clarity.}
0:   else if  $\neg C(p_i)$  then
0:      $p_{i+1} \leftarrow \text{ModifyToMeetConditions}(p_i)$  {Adjust prompt to meet conditions.}
0:   end if
0:    $i \leftarrow i + 1$ 
0: end while
0: return  $p_i = 0$ 

```

The Multi-dialogue Refine process in ChatGPT-4 iteratively refines prompts until they meet predefined conditions and

## Enhancing Intent Understanding for Ambiguous Prompts: A Human-Machine Co-Adaption Strategy

are ambiguity-free. Initially, the model assesses if the prompt  $p_0$  meets specific criteria and lacks ambiguities. If issues are identified, the process loops to rectify them. The model evolves with each iteration, described mathematically as:

$$y_{t+1} = M(p_{\text{refine}} \parallel x \parallel y_0 \parallel \text{fb}_0 \parallel \dots \parallel y_t \parallel \text{fb}_t),$$

where  $y_t$  is the output at iteration  $t$ ,  $M$  represents the model,  $p_{\text{refine}}$  is the refined prompt,  $x$  is the input data, and  $\text{fb}_t$  is the feedback at iteration  $t$ . The model refines prompts by engaging in multi-turn dialogue, asking clarifying questions until the prompts are comprehensive and unambiguous. This self-reflection mechanism allows the model to produce initial responses and evaluate them for retrieval, relevance, support, and utility. Necessary modifications are made based on feedback to enhance accuracy and usefulness, represented as:

$$y_{t+1} = M(x \parallel y_t \parallel \text{fb}_t).$$

### A.8 The Processing of Setup Prompts from Human Feedback Using Large Language Models (LLM)

Table 4. ChatGPT-4 Prompt Rewriting and Type Judgment Process

---

#### Process Description

---

Given the current prompt: ‘{current\_prompt}’, the user has requested changes described as: ‘{user\_input}’.

Please generate a new prompt by incorporating these changes. The alterations should be subtle and maintain the structural integrity of the original prompt.

Modify the original prompt using one of the following methods: ‘word swapping’, ‘adding phrases’, or ‘attention reweighting’.

Ensure that the modifications align closely with the user’s request, and specify which method you used to alter the prompt. The final output format should be ‘{new\_prompt, type}’.

---