

# Parrot: Pareto-optimal Multi-Reward Reinforcement Learning Framework for Text-to-Image Generation

Seung Hyun Lee<sup>1,6\*</sup>, Yinxiao Li<sup>1</sup>, Junjie Ke<sup>1</sup>, Innfarn Yoo<sup>2</sup>,  
Han Zhang<sup>3</sup>, Jiahui Yu<sup>4†</sup>, Qifei Wang<sup>2</sup>, Fei Deng<sup>2,5\*</sup>, Glenn Entis<sup>1</sup>,  
Junfeng He<sup>1</sup>, Gang Li<sup>1</sup>, Sangpil Kim<sup>6</sup>, Irfan Essa<sup>1</sup>, Feng Yang<sup>1</sup>

Google Research<sup>1</sup>, Google<sup>2</sup>, Google DeepMind<sup>3</sup>, OpenAI<sup>4</sup>,  
Rutgers University<sup>5</sup>, Korea University<sup>6</sup>



**Fig. 1: Parrot visual examples.** Parrot consistently improves the quality of generated images across multiple criteria: aesthetics, human preference, text-image alignment, and image sentiment. Each column shows generated images using the same seed.

**Abstract.** Recent works have demonstrated that using reinforcement learning (RL) with multiple quality rewards can improve the quality of generated images in text-to-image (T2I) generation. However, manually adjusting reward weights poses challenges and may cause over-optimization in certain metrics. To solve this, we propose Parrot, which addresses the issue through multi-objective optimization and introduces an effective multi-reward optimization strategy to approximate Pareto optimal. Utilizing batch-wise Pareto optimal selection, Parrot automatically identifies the optimal trade-off among different rewards. We use the novel multi-reward optimization algorithm to jointly optimize the T2I model and a prompt expansion network, resulting in significant improvement of image quality and also allow to control the trade-off of

\* This work was done during an internship at Google.

† This work was done during working at Google.

different rewards using a reward related prompt during inference. Furthermore, we introduce original prompt-centered guidance at inference time, ensuring fidelity to user input after prompt expansion. Extensive experiments and a user study validate the superiority of Parrot over several baselines across various quality criteria, including aesthetics, human preference, text-image alignment, and image sentiment.

## 1 Introduction

Despite significant advancements in text-to-image (T2I) generation [17, 29, 38, 45, 55], recent work like Imagen [42] and Stable Diffusion [39], still struggle to produce high quality images. Images in the first row in Fig. 1 illustrates such quality issues in Stable Diffusion [39], including poor composition (*e.g.* bad cropping), misalignment with input prompts (*e.g.* missing objects), or overall lack of aesthetic appeal. Assessing the quality of generated images can involve various metrics such as aesthetics [23, 24, 33], human preference [26], text-image alignment [35], and emotional appeal [44]. Enhancing T2I generation across multiple quality metrics remains a challenging task.

Recent works [4, 12, 15, 27] have demonstrated that incorporating quality signals as reward functions in fine-tuning T2I with reinforcement learning (RL) can improve image quality. For instance, Promptist [15] fine-tunes prompt expansion model using RL with the sum of aesthetics and text-image alignment scores. However, the simple weighted sum approach may not effectively handle trade-offs among multiple quality metrics. As the number of rewards increases, manually adjusting reward weights becomes impractical. Moreover, optimizing one quality metric may inadvertently compromise others, as the model might prioritize aesthetics over relevance to the input prompt. Additionally, the trade-off for different reward is not controllable after training.

To address these challenges, we propose **Parrot**, a **Pareto-optimal multi-reward reinforcement learning** algorithm to improve text-to-image generation. Unlike previous approaches that treat T2I reward optimization as a single objective optimization problem, Parrot tackles this challenge through multi-objective optimization and introduces an effective multi-reward optimization strategy to achieve Pareto optimal approximation. Intuitively, each generated sample in a batch embodies a distinctive trade-off among various quality rewards, with some samples exhibiting superior trade-offs compared to others. Instead of updating gradients using all batch samples, Parrot uses non-dominated points [32] which have better trade-offs. Consequently, Parrot automatically learns from the optimal trade-off among different rewards. Moreover, Parrot learns reward-specific preference prompts, which can be utilized individually or in combination to control the trade-off among different rewards during inference time. Unlike prior work, which either solely fine-tunes the T2I model [12] or only tunes the prompt expansion network while freezing the T2I model [15], we employ the Parrot multi-reward optimization algorithm to jointly optimize both the T2I model and the prompt expansion network. This collaborative optimization unlocks the

full potential of Parrot by encouraging both more details from added context from the prompt expansion model, and the overall quality improvement on the T2I generation. During inference, we further introduce *original prompt-centered guidance* to ensure the output image is relevant to input prompts after prompt expansion.

In summary, our contributions can be listed as follows:

- We propose Parrot, a novel multi-reward optimization algorithm for T2I RL fine-tuning. Leveraging batch-wise Pareto-optimal selection, it effectively optimizes multiple T2I rewards, enabling collaborative improvement in aesthetics, human preference, image sentiment, and text-image alignment and also allowing to control the trade-off of different rewards using a reward related prompt during inference.
- We show the advantage of jointly optimizing both the prompt expansion network and the T2I model, which has never been explored before.
- We introduce original prompt-centered guidance during inference time after prompt expansion, ensuring better alignment with the original prompt while enriching image details.
- Extensive results and a user study validate that Parrot outperforms several baseline methods across various quality criteria.

## 2 Related Work

**T2I Generation:** The goal of T2I generation is to create an image given an input text prompt. Several T2I generative models have been proposed and have demonstrated promising results [5, 7, 14, 21, 22, 28, 37, 41, 42, 54]. Stable Diffusion [39] shows impressive generation performance in T2I generation, leveraging latent text representations from LLMs. Despite substantial progress, the images generated by those models still exhibit quality issues, such as bad cropping or misalignment with the input texts.

**RL for T2I Fine-tuning:** Starting by Fan *et al.* [11] to explore RL fine-tuning for T2I models, following works [4, 8, 9, 12, 16] have explored RL fine-tuning technique for T2I diffusion model, showcasing superior performance for human preference learning. DPOK [12] improves quality through RL using ImageReward [51] score as a reward with a few prompts. In addition to fine-tuning the T2I model directly using RL, Promptist [15] fine-tunes the prompt expansion model by using a simple sum of aesthetic and text-image alignment scores as reward. DRaFT [6] proposed not only differentiable rewards for efficient fine-tuning but also effectiveness of using linear summation of multi-rewards. These methods treat T2I RL as a single-objective optimization problem, while Parrot employs multi-objective optimization. Additionally, prior approaches either fine-tune the T2I model or the prompt expansion model while freezing the other. In contrast, Parrot proposes joint optimization of the prompt expansion model and the T2I model using multi-reward RL to foster better collaboration.

**Multi-objective Optimization:** Multi-objective optimization problem involves optimizing multiple objective functions simultaneously. The scalarization technique [31, 47] formulates multi-objective problem into single-objective problems with the weighted sum of each score, which requires pre-defined weights for each objective. Rame *et al.* [36] proposed weighted averaging method to find Pareto frontier, leveraging multiple fine-tuned models. Lin *et al.* [30] proposes to learn a set model to map trade-off preference vectors to their corresponding Pareto solutions. Inspired by this, Parrot introduces a language-based preference vector constructed from task identifiers for each reward, then encoded by the text-encoder. In the context of multi-reward RL for T2I diffusion models, Promptist [15] uses a simple weighted sum of two reward scores. This approach requires manual tuning of the weights, which makes it time-consuming and hard to scale when the number of rewards increases.

**Generated Image Quality:** The quality assessment of images generated by T2I models involves multiple dimensions, and various metrics have been proposed. In this paper, we consider using four types of quality metrics as rewards: aesthetics, human preference, text-image alignment, and image sentiment. Aesthetics captures the overall visual appealingness of the image, and it is learned using human ratings for aesthetics in real images [13, 19, 23, 24, 33, 48, 52]. Human preferences, rooted the concept of learning from human feedback [3, 34, 50], involves gathering preferences at scale by having raters to compare generated images [26, 51]. Text-image alignment measures the extent to which the generated image aligns with the input prompt, CLIP [35] score is often employed, measuring the cosine distance of between contrastive image embedding and text embedding. Image sentiment is important for ensuring the generated image evokes positive emotions in the viewer. Serra *et al.* [44] predict average polarity of sentiments an image elicits and learn estimates for positive, neutral, and negative scores. In Parrot, we use its positive score as a reward for positive emotions.

### 3 Preliminary

**Diffusion Probabilistic Models:** Diffusion probabilistic models [17] generate the image by gradually denoising a noisy image. Specifically, given a real image  $\mathbf{x}_0$  from the data distribution  $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ , the forward process  $q(\mathbf{x}_t|\mathbf{x}_0, c)$  of diffusion probabilistic models produce a noisy image  $\mathbf{x}_t$ , which induces a distribution  $p(\mathbf{x}_0, c)$  conditioned on text prompt  $c$ . In classifier-free guidance [18], denoising model predicts noise  $\bar{\epsilon}_\theta$  with a linear combination of the unconditional score estimates  $\epsilon_\theta(\mathbf{x}_t, t)$  and the conditional score estimates  $\epsilon_\theta(\mathbf{x}_t, t, c)$  as follows:

$$\bar{\epsilon}_\theta = w \cdot \epsilon_\theta(\mathbf{x}_t, t, c) + (1 - w) \cdot \epsilon_\theta(\mathbf{x}_t, t, \text{null}), \quad (1)$$

where  $t$  denotes diffusion time step, the null indicates a null text and  $w$  represents the guidance scale of classifier-free guidance where  $w \geq 1$ . Note that  $\epsilon_\theta$  is typically parameterized by the UNet [40].

**RL-based T2I Diffusion Model Fine-tuning:** Given a reward signal from generated images, the goal of RL-tuning for T2I diffusion models is to optimize the policy defined as one denoising step of T2I diffusion models. In particular, Black *et al.* [4] apply policy gradient algorithm, which regards the denoising process of diffusion models as a Markov decision process (MDP) by performing multiple denoising steps iteratively. Subsequently, a black box reward model  $r(\cdot, \cdot)$  predicts a single scalar value from sampled image  $\mathbf{x}_0$ . Given text condition  $c \sim p(c)$  and image  $\mathbf{x}_0$ , objective function  $\mathcal{J}$  can be defined to maximize the expected reward as follows:

$$\mathcal{J}_\theta = \mathbb{E}_{p(c)} \mathbb{E}_{p_\theta(\mathbf{x}_0|c)} [r(\mathbf{x}_0, c)], \quad (2)$$

where the pre-trained diffusion model  $p_\theta$  produces a sample distribution  $p_\theta(\mathbf{x}_0|c)$  using text condition  $c$ . Modifying this equation, Fan *et al.* [12] demonstrate that the gradient of objective function  $\nabla \mathcal{J}_\theta$  can be calculated through gradient ascent algorithm without using the gradient of reward model as follows:

$$\nabla \mathcal{J}_\theta = \mathbb{E}[r(\mathbf{x}_0, c) \sum_{t=1}^T \nabla_\theta \log p_\theta(\mathbf{x}_{t-1}|c, t, \mathbf{x}_t)], \quad (3)$$

where  $T$  denotes the total time step of the diffusion sampling. With parameters  $\theta$ , the expectation value can be taken over the trajectories of diffusion sampling.

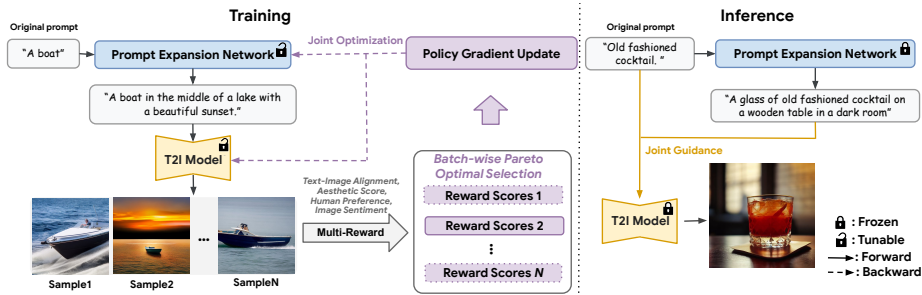
## 4 Method

### 4.1 Parrot Overview

Fig. 2 shows the overview of Parrot, which consists of the prompt expansion network (PEN)  $p_\phi$  and the T2I diffusion model  $p_\theta$ . The PEN is first initialized from a supervised fine-tuning checkpoint on demonstrations of prompt expansion pairs, and the T2I model is initialized from pretrained diffusion model. Given the original prompt  $c$ , the PEN generates an expanded prompt  $\hat{c}$ , and the T2I model generates images based on this expanded prompt. During the multi-reward RL fine-tuning, a batch of  $N$  images is sampled, and multiple quality rewards are calculated for each image, encompassing aspects like text-image alignment, aesthetics, human preference, and image sentiment. Based on these reward scores, Parrot identifies the batch-wise Pareto-optimal set using a non-dominated sorting algorithm. This optimal set of images is then used for joint optimization of the PEN and T2I model parameters through RL policy gradient update. During inference, Parrot leverages both the original prompt and its expansion, striking a balance between maintaining faithfulness to the original prompt and incorporating additional details for higher quality.

### 4.2 Batch-wise Pareto-optimal Selection

Li *et al.* [30] has demonstrated that using batchwise Pareto-set learning and selecting good samples in a batch can approximate Pareto-optimality across multiple objectives. Backed up by this theory, we propose to select non-dominated



**Fig. 2:** Overview of **Parrot**. During the training,  $N$  images are sampled from the T2I model using the expanded prompt from the prompt expansion network. Multiple quality rewards are calculated for each image, and the Pareto-optimal set is identified using the non-dominated sorting algorithm. These optimal images are then used to perform policy gradient update of the parameters of T2I model and prompt expansion network jointly. During the inference, both the original prompt and the expanded prompt are provided to the T2I model, enabling better faithfulness while adding detail.

points in a batch for policy gradient update in RL to achieve Pareto-optimality. Algorithm 1 outlines the procedure of Parrot. Rather than updating the gradients using all images, Parrot focuses on high-quality samples, considering multiple quality rewards in each mini-batch. In the multi-reward RL, each sample generated by the T2I model presents distinct trade-offs for different rewards. Among these samples, a subset with varied optimal trade-offs across multiple objectives, also known as the Pareto set, exists. For a Pareto-optimal sample, none of its objective values can be further improved without damaging others. In other words, the Pareto-optimal set is not dominated by any data points, also known as *the non-dominated set*. To achieve a Pareto-optimal solution with text-to-image generation diffusion model, Parrot selectively uses data points from the non-dominated set using non-dominated sorting algorithm. This naturally encourages the T2I model to produce Pareto-optimal samples with respect to the multi-reward objectives.

**Reward-specific Preference:** Inspired by the use of preference information in multi-objective optimization [30], Parrot incorporates the preference information through reward-specific identifiers. This enables Parrot to automatically determine the importance for each reward objective. Concretely, we enrich the expanded prompt  $\hat{c}$  by prepending reward-specific identifier “<reward  $k$ >” for  $k$ -th reward. Based on this reward-specific prompt,  $N$  images are generated and are used for maximizing the corresponding  $k$ -th reward model during gradient update. At inference time, a concatenation of all the reward identifiers “<reward 1>, ..., <reward  $K$ >” is used for image generation.

**Non-dominated Sorting:** Parrot constructs Pareto set with non-dominated points based on trade-offs among multiple rewards. These non-dominated points are superior to the remaining solutions and are not dominated by each other. Formally, the dominance relationship is defined as follows: the image  $\mathbf{x}_0^a$  domi-

---

**Algorithm** Parrot: Pareto-optimal Multi-Reward RL

**Input:** Prompt  $c$ , Batch size  $N$ , Total iteration  $E$ , the number of rewards:  $K$ , Prompt expansion network  $p_\phi$ , T2I diffusion model:  $p_\theta$ , Total diffusion time step  $T$ , Non-dominated set:  $\mathcal{P}$

```

for  $e = 1$  to  $E$  do
    Sample text prompt  $c \sim p(c)$ 
    for  $k = 1$  to  $K$  do
        Expand text prompt  $\hat{c} \sim p_\phi(\hat{c}|c)$ 
        Prepend reward-specific tokens " $\langle \text{reward } k \rangle$ " to  $\hat{c}$ 
        Sample a set of images  $\{\mathbf{x}_0^1, \dots, \mathbf{x}_0^N\} \sim p_\theta(x_0|\hat{c})$ 
        A set of reward vector  $\mathcal{R} = \{R_1, \dots, R_N\}$ 
         $\mathcal{P} \leftarrow \text{NDSSET}(\{\mathbf{x}_0^1, \dots, \mathbf{x}_0^N\})$ 
         $\nabla \mathcal{J}_\phi += -r_k(\mathbf{x}_0^j, \hat{c}) \times \nabla \log p_\phi(\hat{c}|c)$ 
    Update the gradient  $p_\theta$  from Eq. 4
    Update the gradient  $p_\phi$ 
    function NDSSET( $\{\mathbf{x}_0^1, \dots, \mathbf{x}_0^N\}$ )
         $\mathcal{P} \leftarrow \emptyset$ 
        for  $i = 1$  to  $N$  do
            dominance  $\leftarrow$  True
            for  $j = 1$  to  $N$  do
                if  $\mathbf{x}_0^j$  dominates  $\mathbf{x}_0^i$  then
                    dominance  $\leftarrow$  False
            if dominance is True then
                Add  $i$  to  $\mathcal{P}$ 
        return  $\mathcal{P}$ 

```

**Output:** Fine-tuned diffusion model  $p_\theta$ , prompt expansion network  $p_\phi$

---

nates the image  $\mathbf{x}_0^b$ , denoted as  $\mathbf{x}_0^b < \mathbf{x}_0^a$ , if and only if  $R_i(\mathbf{x}_0^b) \leq R_i(\mathbf{x}_0^a)$  for all  $i \in 1, 2, \dots, m$ , and there exists  $j \in 1, 2, \dots, m$  such that  $R_j(\mathbf{x}_0^b) < R_j(\mathbf{x}_0^a)$ . For example, given the  $i$ -th generated image  $\mathbf{x}_0^i$  in a mini-batch, when no point in the mini-batch dominates  $\mathbf{x}_0^i$ , it is referred to as a non-dominated point.

**Policy Gradient Update:** We assign a reward value of zero to the data points not included in non-dominated sets and only update the gradient of these non-dominated data points as follows:

$$\nabla \mathcal{J}_\theta = \sum_{k=1}^K \frac{1}{n(\mathcal{P})} \sum_{i=1, \mathbf{x}_0^i \in \mathcal{P}}^N \sum_{t=1}^T r_k(\mathbf{x}_0^i, c_k) \times \nabla_\theta \log p_\theta(\mathbf{x}_{t-1}^i | c_k, t, \mathbf{x}_t^i), \quad (4)$$

where  $i$  indicates the index of images in mini-batches, and  $\mathcal{P}$  denotes batch-wise a set of non-dominated points.  $K$  and  $T$  are the total number of reward models and total diffusion time steps, respectively. The same text prompt is used when updating the diffusion model in each batch.

### 4.3 Original Prompt Centered Guidance

While prompt expansion enhances details and often improves generation quality, there is a concern that the added context may dilute the main content of the original input. To mitigate this during the inference, we introduce original prompt-centered guidance. When sampling conditioned on the original prompt, the diffusion model  $\epsilon_\theta$  typically predicts noises by combining the unconditioned score estimate and the prompt-conditioned estimate. Instead of relying solely on the expanded prompt from PEN, we propose using a linear combination of two guidances for T2I generation: one from the user input and the other from the expanded prompt. The strength of the original prompt is controlled by guidance scales  $w_1$  and  $w_2$ . The noise  $\bar{\epsilon}_\theta$  is estimated, derived from Eq. 1, as follows:

$$\bar{\epsilon}_\theta = w_1 \cdot \epsilon_\theta(\mathbf{x}_t, t, c) + (1 - w_1 - w_2) \cdot \epsilon_\theta(\mathbf{x}_t, t, \text{null}) + w_2 \cdot \epsilon_\theta(\mathbf{x}_t, t, \hat{c}), \quad (5)$$

where null denotes a null text.

## 5 Experiments

### 5.1 Experiment Setting

**Dataset:** The PEN is first supervised fine-tuned on a large-scale text dataset named the Promptist [15], which has 360K constructed prompt pairs for original prompt and prompt expansion demonstration. The original instruction “Rephrase” is included per pair in Promptist. We modify the instruction into “Input: <original prompt>. This is a text input for image generation. Expand prompt for improving image quality. Output: ”. Subsequently, we use the RL tuning prompts (1200K) from Promptist for RL training of the PEN and T2I model.

**T2I Model:** Our T2I model is based on the JAX version of Stable Diffusion 1.5 [39] pre-trained with the LAION-5B [43] dataset. We conduct experiments on a machine equipped with 16 NVIDIA RTX A100 GPUs. DDIM [46] with 50 denoising steps is used, and the classifier-free guidance weight is set to 5.0 with the resolution  $512 \times 512$ . Instead of updating all layers, we specifically update the cross-attention layer in the Denoising U-Net. For optimization, we employ the Adam [25] optimizer with a learning rate of  $1 \times 10^{-5}$ .

**Prompt Expansion Network:** For prompt expansion, we use PaLM 2-L-IT [2], one of the PaLM2 variations, which is a multi-layer Transformer [49] decoder with casual language modeling. We optimize LoRA [20] weights for RL-based fine-tuning. The output token length of the PEN is set to 77 to match the maximum number of token length for Stable Diffusion. For original prompt-centered guidance, we set both  $w_1$  and  $w_2$  to 5 in Eq. 5.

**Reward Models:** We incorporate four quality signals as rewards: Aesthetics, Human preference, Text-Image Alignment, Image Sentiment. For aesthetics, we use the VILA-R [24] pre-trained with the AVA [33] dataset. For human preference, we train a ViT-B/16 [10] using the Pick-a-Pic [26] dataset, which contains 500K examples for human feedback in T2I generation. The ViT-B/16 image encoder consists of 12 transformer layers, and the image resolution is  $224 \times 224$  with a patch size of  $16 \times 16$ . For text-image alignment, we use CLIP [35] with the image encoder ViT-B/32. For image sentiment, we use the pre-trained model from [44], which outputs three labels: positive, neutral, negative. We use the positive score ranging from 0 to 1 as the sentiment reward.

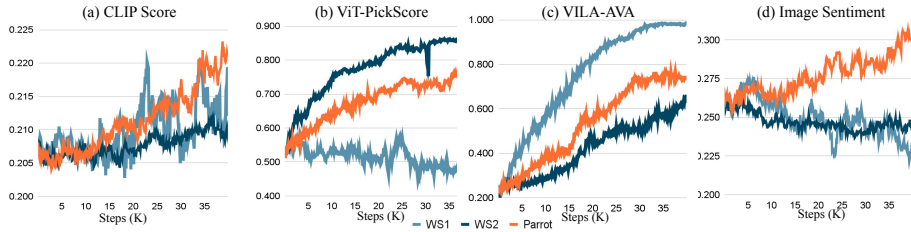
### 5.2 Qualitative Analysis

**Comparison with Baselines:** Fig. 3 shows the visual comparison of Parrot and multiple baselines. We include results from Stable Diffusion 1.5, DPOK [12] with a weighted sum of rewards, Promptist [15], and Parrot. DPOK exclusively fine-tunes the T2I model, while Promptist focuses on fine-tuning only the prompt expansion network. Parrot shows visually better images, particularly in aspects like color combination, cropping, perspective, and fine details in the image. This improvement can be attributed to Parrot’s T2I model being fine-tuned together





**Fig. 3: Comparison of Parrot and diffusion-based RL baselines.** From left to right, we provide results of Stable diffusion 1.5 [39] (1st column), DPOK [12] (2nd column) with the weighted sum, Promptist [15] (3rd column), and Parrot (4th column).



**Fig. 4: Training curve for fine-tuning on weighted sum and Parrot.** For weighted sum, WS1 denotes  $\{0.7, 0.1, 0.1, 0.1\}$  and WS2 denotes  $\{0.25, 0.25, 0.25, 0.25\}$  for aesthetics, human preference, text-image alignment and image sentiment. Using weighted sum leads to decrease in human preference score and image sentiment score despite an improvement in the aesthetic score. In contrast, Parrot exhibits stable increases across all metrics.

Model	Quality Metrics				
	TIA ( $\uparrow$ )	Aesth. ( $\uparrow$ )	HP ( $\uparrow$ )	Sent. ( $\uparrow$ )	Average ( $\uparrow$ )
SD 1.5 [39]	0.2322	0.5755	0.1930	0.3010	0.3254
DPOK [12] (WS)	0.2337	0.5813	0.1932	0.3013	0.3273 (+0.58%)
Parrot w/o PE	0.2355	0.6034	0.2009	0.3018	0.3354 (+3.07%)
Parrot T2I Model Tuning Only	<b>0.2509</b>	<b>0.7073</b>	<b>0.3337</b>	<b>0.3052</b>	<b>0.3992 (+22.6%)</b>
Promptist [15]	0.1449	0.6783	0.2759	0.2518	0.3377 (+3.77 %)
Parrot with HP Only	0.1543	0.5961	<b>0.3528</b>	0.2562	0.3398 (+4.42 %)
Parrot PEN Tuning Only	0.1659	0.6492	0.2617	0.3131	0.3474 (+6.76 %)
Parrot w/o Joint Optimization	0.1661	0.6308	0.2566	0.3084	0.3404 (+4.60 %)
Parrot w/o ori prompt guidance	0.1623	0.7156	0.3425	0.3130	0.3833 (+17.8 %)
Parrot	<b>0.1667</b>	<b>0.7396</b>	0.3411	<b>0.3132</b>	<b>0.3901 (+19.8 %)</b>

**Table 1:** Quantitative comparison between Parrot and alternatives on the Parti dataset [53]. Abbreviations: WS - Weighted Sum; PE - Prompt Expansion; TIA - Text-Image Alignment; Aesth. - Aesthetics; HP - Human Preference; Sent. - Image Sentiment. TIA score is measured against the original prompt without expansion.

with the prompt expansion model that incorporates aesthetic keywords during training. Parrot generates results that are more closely aligned with the input prompt, as well as more visually pleasing.

**Weighted sum vs. Parrot:** Fig. 4 shows the training curve comparison of Parrot and using a linear combination of the reward scores. Each subgraph represents a reward. WS1 and WS2 denote two different weights with multiple reward scores. WS1 places greater emphasis on the aesthetic score, while WS2 adopts balanced weights across aesthetics, human preference, text-image alignment, and image sentiment. Employing the weighted sum of multiple rewards leads to a decrease in the image sentiment score, despite notable enhancements in aesthetics and human preference. In contrast, Parrot consistently exhibits improvement across all metrics.

### 5.3 Quantitative Evaluation

**Comparison with Baselines:** Table 1 presents our results of the quality score across four quality rewards: text-image alignment score, aesthetic score, human preference score, and emotion score. The first group shows methods without prompt expansion, and the second group compares methods with expansion. The prompt expansion and T2I generation are performed on the PartiPrompts [53]. Using a set of 1632 prompts, we generate 32 images for each text input and calculate the average for each metric.

*w/o PE* indicates the generation of images solely based on original prompt without expansion. Note that *w/o PE* is not the same as *Parrot T2I Model Tuning Only*. The former takes a model trained with PEN and removes it during inference, which results in a notable disparity between prompts used in training and testing. The latter trains with the multi-objective RL using only the T2I model, and it shows substantial improvement upon DPOK [12] (weighted sum) with balanced weights of  $\{0.25, 0.25, 0.25, 0.25\}$ . This shows that the proposed RL method is indeed effective in multi-objective optimization.

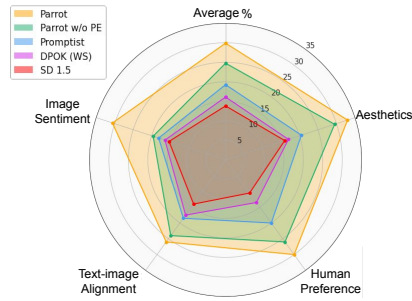
For Promptist [15], we generate prompt expansion from their model. *Parrot with HP Only* shows solely using a human preference reward leads to decline in others. *Parrot PEN Tuning Only* shows suboptimal aesthetics and human preference. *Parrot w/o Joint Optimization* shows suboptimal results than Parrot which demonstrates the necessity of jointly optimizing PEN and T2I models.

Our method outperforms both compared methods in aesthetics, human preference and sentiment scores. The text-image alignment score is measured with the original prompt before expansion for fair comparison. As a result, the group without prompt expansion generally shows a higher text-image alignment score. Parrot shows better text-image alignment in each subgroup.

Area	Question
Aesthetics	“Which image shows better aesthetics without blurry texture, unnatural focusing, and poor color combination?”
Human Preference	“Which generated image do you prefer?”
Text-Image Alignment	“Which image is well aligned with the text?”
Image Sentiment	“Which image is closer to amusement, excitement, and contentment?”

**Table 2:** Questions for user study. For performing user study, we carefully design questions suitable to each quality.

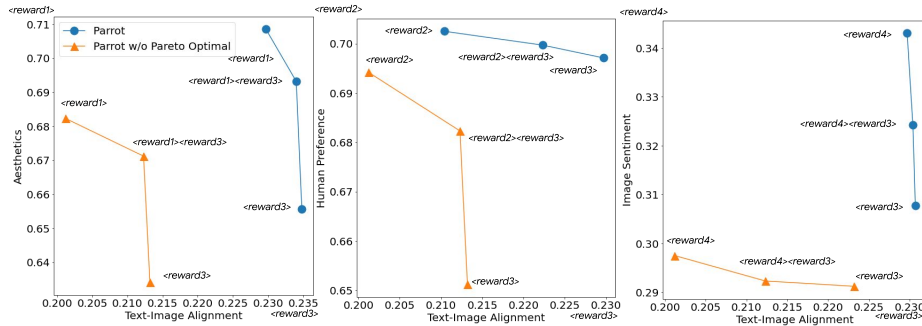
**User Study:** We conduct a user study using MTurk [1] with generated images from 100 random prompts in the PartiPrompts [53]. Five models are compared: Stable Diffusion v1.5, DPOK [12] with an equal weighted sum, Promptist [15], Parrot without prompt expansion, and Parrot. Each rater is presented with the original prompt (before expansion) and a set of five generated images, with the image order being randomized. Raters are then tasked with selecting the best image from the group, guided by questions outlined in Table 2. Each question



**Fig. 5:** User study results on ParrotPrompts [53]. Parrot outperforms baselines across all metrics.



**Fig. 6:** Ablation study. We perform an ablation study by removing one of quality signals. We observe that each quality signal affects their improvement of (a) aesthetics, (b) human preference score, (c) text-image alignment score, (d) image sentiment score.



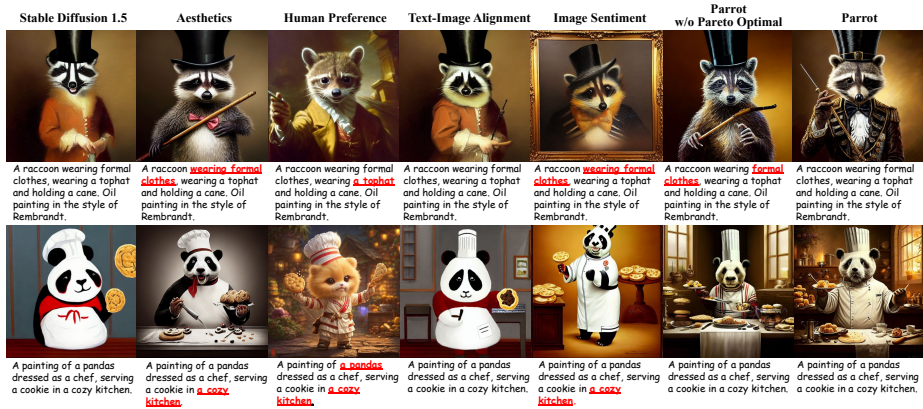
**Fig. 7:** Comparing Parrot with/without Pareto Optimal. By prepending different reward-specific preferences, Parrot can change the tradeoff between rewards. Specifically, “<reward 1>”, “<reward 2>”, “<reward 3>”, “<reward 4>” indicates aesthetics, human preference, text-image alignment, and image sentiment respectively. With Pareto-optimal selection, each reward increases based on reward-specific preference.

pertains to a specific quality aspect: aesthetics, human preference, text-image alignment, and image sentiment. For each prompt, 20 rounds of random sampling are conducted and sent to different raters. The user study results, illustrated in Fig. 5, show that Parrot outperforms other baselines across all dimensions.

**Proportion of non-dominated points:** Using batch size of 256, we observe that in a batch around 20% to 30% are non-dominated points and the proportion of non-dominated points in single batch slightly increase as training proceeds.

## 5.4 Ablations

**Effect of Pareto-optimal Multi-reward RL:** To show the efficacy of Pareto-optimal Multi-reward RL, we conduct an ablation study by removing one reward model at a time. Fig. 6 shows quantitative results using one hundred random text



**Fig. 8: The comparisons of the diffusion fine-tuning between Pareto-optimal multi-reward RL and single reward RL.** We show results with same seed from various methods: Stable Diffusion 1.5 [39] (1st column), T2I model fine-tuned with the aesthetic model (2nd column), the human preference model (3rd column), text-image alignment (4th column), image sentiment (5th column), Parrot without Pareto-optimal selection (6th column) and Parrot (7th column). Parrot is effective to generate acceptable images without sacrificing one of quality signals. For example, T2I model fine-tuned with a single quality signal such as aesthetics, human preference and image sentiment results in text-image misalignment, while our method achieves a balanced visual outcome across multiple criteria.

prompts from the Promptist [15]. We observe that our training scheme improves multiple target objectives.

To verify whether Parrot achieved better trade-off for different reward scores, we generate 1000 images from common animal dataset [4], where each text prompt consists of the name of a common animal. As shown in Fig. 7, using only text-image alignment with reward-specific preference “<reward 3>” generates images with higher text-image alignment score, while using only aesthetic model with reward-specific preference “<reward 1>” yields images with higher aesthetic score. In the case of using two reward-specific preferences “<reward 1>, <reward 3>”, we observe that scores are balanced and show that better results across multiple rewards than Parrot without Pareto optimal selection.

Fig 8 shows the visual comparison between Parrot, Parrot with a single reward, and Parrot without selecting the batch-wise Pareto-optimal solution. Using a single reward model tends to result in degradation of other rewards, especially text-image alignment. In the third column, results of the first row miss the text *a tophat* in input prompt, even though the Stable Diffusion result includes that attribute. On the other hand, Parrot results capture all prompts, improving other quality signals, such as aesthetics, image sentiment and human preference.

**Effect of Original Prompt Centered Guidance:** Fig 9 shows the effect of the proposed original prompt-centered guidance. As evident from the figure, using only the expanded prompt as input often results in the main content be-

ing overwhelmed by the added context. For instance, given the original prompt “A *shiba inu*”, the result from the expanded prompt shows a zoomed-out image and the intended main subject (*shiba inu*) becomes small. The proposed original prompt-centered guidance effectively addresses this issue, generating an image that faithfully captures the input prompt while incorporating visually more pleasing details.



**Fig. 9: Results of original prompt centered guidance.** As we expand the prompt, the content in the generated image often fades away. This guidance is helpful for keeping the main content of the original prompt.

## 6 Conclusion and Limitation

We propose Parrot, a novel multi-reward optimization algorithm aimed to improve text-to-image generation by effectively optimizing multiple quality rewards using RL. With batch-wise Pareto-optimal selection, Parrot adaptively balance the optimization of multiple quality rewards. By applying Parrot to jointly fine-tune both the T2I model and the prompt expansion model, we achieve the generation of higher-quality images with richer details. Additionally, our original prompt centered guidance technique ensures that the generated image maintains fidelity to the user prompt after prompt expansion during inference. Results from the user study indicate that Parrot significantly improves the quality of generated images across multiple criteria, including text-image alignment, human preference, aesthetics, and image sentiment. While Parrot has shown effectiveness in enhancing generated image quality, its efficacy is limited by the quality metrics it relies on. Therefore, advancements of the generated image quality metrics will directly enhance the capabilities of Parrot. Additionally, Parrot is adaptable to a broader range of rewards that quantify generated image quality.

**Societal Impact:** Parrot could potentially raise ethical concerns related to the generation of immoral content. This concern stems from the user’s ability to influence T2I generation, allowing for the creation of visual content that may be deemed inappropriate. The risk may be tied to the potential biases in reward models inherited from various datasets.

## References

1. Amazon mechanical turk. <https://www.mturk.com/> (2005)
2. Anil, R., Dai, A.M., Firat, O., Johnson, M., Lepikhin, D., Passos, A., Shakeri, S., Taropa, E., Bailey, P., Chen, Z., et al.: Palm 2 technical report. arXiv preprint arXiv:2305.10403 (2023)
3. Bai, Y., Jones, A., Ndousse, K., Askell, A., Chen, A., DasSarma, N., Drain, D., Fort, S., Ganguli, D., Henighan, T., et al.: Training a helpful and harmless assistant with reinforcement learning from human feedback. arXiv preprint arXiv:2204.05862 (2022)
4. Black, K., Janner, M., Du, Y., Kostrikov, I., Levine, S.: Training diffusion models with reinforcement learning. arXiv preprint arXiv:2305.13301 (2023)
5. 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)
6. Clark, K., Vicol, P., Swersky, K., Fleet, D.J.: Directly fine-tuning diffusion models on differentiable rewards. ICLR (2024)
7. Dai, X., Hou, J., Ma, C.Y., Tsai, S., Wang, J., Wang, R., Zhang, P., Vandenhende, S., Wang, X., Dubey, A., et al.: Emu: Enhancing image generation models using photogenic needles in a haystack. arXiv preprint arXiv:2309.15807 (2023)
8. Deng, F., Wang, Q., Wei, W., Grundmann, M., Hou, T.: Prdp: Proximal reward difference prediction for large-scale reward finetuning of diffusion models. CVPR (2024)
9. Dong, H., Xiong, W., Goyal, D., Pan, R., Diao, S., Zhang, J., Shum, K., Zhang, T.: Raft: Reward ranked finetuning for generative foundation model alignment. arXiv preprint arXiv:2304.06767 (2023)
10. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020)
11. Fan, Y., Lee, K.: Optimizing ddp sampling with shortcut fine-tuning. ICML (2023)
12. Fan, Y., Watkins, O., Du, Y., Liu, H., Ryu, M., Boutilier, C., Abbeel, P., Ghavamzadeh, M., Lee, K., Lee, K.: Dpok: Reinforcement learning for fine-tuning text-to-image diffusion models. NeurIPS (2023)
13. Fang, Y., Zhu, H., Zeng, Y., Ma, K., Wang, Z.: Perceptual quality assessment of smartphone photography. In: CVPR (2020)
14. Han, L., Li, Y., Zhang, H., Milanfar, P., Metaxas, D., Yang, F.: Svdif: Compact parameter space for diffusion fine-tuning. CVPR (2023)
15. Hao, Y., Chi, Z., Dong, L., Wei, F.: Optimizing prompts for text-to-image generation. arXiv preprint arXiv:2212.09611 (2022)
16. He, H., Wang, T., Yang, H., Fu, J., Yuan, N.J., Yin, J., Chao, H., Zhang, Q.: Learning profitable nft image diffusions via multiple visual-policy guided reinforcement learning. arXiv preprint arXiv:2306.11731 (2023)
17. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. NeurIPS (2020)
18. Ho, J., Salimans, T.: Classifier-free diffusion guidance. arXiv preprint arXiv:2207.12598 (2022)

19. Hosu, V., Lin, H., Sziranyi, T., Saupe, D.: Koniq-10k: An ecologically valid database for deep learning of blind image quality assessment. *TIP* **29**, 4041–4056 (2020)
20. Hu, E.J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W.: Lora: Low-rank adaptation of large language models. arXiv preprint arXiv:2106.09685 (2021)
21. Jeong, Y., Ryoo, W., Lee, S., Seo, D., Byeon, W., Kim, S., Kim, J.: The power of sound (tpos): Audio reactive video generation with stable diffusion. In: *ICCV* (2023)
22. Kawar, B., Zada, S., Lang, O., Tov, O., Chang, H., Dekel, T., Mosseri, I., Irani, M.: Imagic: Text-based real image editing with diffusion models. In: *CVPR* (2023)
23. Ke, J., Wang, Q., Wang, Y., Milanfar, P., Yang, F.: Musiq: Multi-scale image quality transformer. In: *ICCV* (2021)
24. Ke, J., Ye, K., Yu, J., Wu, Y., Milanfar, P., Yang, F.: Vila: Learning image aesthetics from user comments with vision-language pretraining. In: *CVPR* (2023)
25. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. *ICLR* (2015)
26. Kirstain, Y., Polyak, A., Singer, U., Matiana, S., Penna, J., Levy, O.: Pick-a-pic: An open dataset of user preferences for text-to-image generation. arXiv preprint arXiv:2305.01569 (2023)
27. Lee, K., Liu, H., Ryu, M., Watkins, O., Du, Y., Boutilier, C., Abbeel, P., Ghavamzadeh, M., Gu, S.S.: Aligning text-to-image models using human feedback. arXiv preprint arXiv:2302.12192 (2023)
28. Lee, S.H., Kim, S., Yoo, I., Yang, F., Cho, D., Kim, Y., Chang, H., Kim, J., Kim, S.: Soundini: Sound-guided diffusion for natural video editing. arXiv preprint arXiv:2304.06818 (2023)
29. Li, Y., Liu, H., Wu, Q., Mu, F., Yang, J., Gao, J., Li, C., Lee, Y.J.: Gligen: Open-set grounded text-to-image generation. In: *CVPR* (2023)
30. Lin, X., Yang, Z., Zhang, X., Zhang, Q.: Pareto set learning for expensive multi-objective optimization. *NeurIPS* (2022)
31. Mannor, S., Shimkin, N.: The steering approach for multi-criteria reinforcement learning. *NeurIPS* (2001)
32. Miettinen, K.: *Nonlinear multiobjective optimization*, vol. 12. Springer Science & Business Media (1999)
33. Murray, N., Marchesotti, L., Perronnin, F.: Ava: A large-scale database for aesthetic visual analysis. In: *CVPR* (2012)
34. Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al.: Training language models to follow instructions with human feedback. *NeurIPS* (2022)
35. 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: *ICML* (2021)
36. Rame, A., Couairon, G., Shukor, M., Dancette, C., Gaya, J.B., Soulier, L., Cord, M.: Rewarded soups: towards pareto-optimal alignment by interpolating weights fine-tuned on diverse rewards. *NeurIPS* (2023)
37. Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., Chen, M.: Hierarchical text-conditional image generation with clip latents. arXiv preprint arXiv:2204.06125 (2022)
38. Richardson, E., Goldberg, K., Alaluf, Y., Cohen-Or, D.: Conceptlab: Creative generation using diffusion prior constraints. arXiv preprint arXiv:2308.02669 (2023)
39. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: *CVPR* (2022)



40. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: MICCAI (2015)
41. Saharia, C., Chan, W., Chang, H., Lee, C., Ho, J., Salimans, T., Fleet, D., Norouzi, M.: Palette: Image-to-image diffusion models. In: SIGGRAPH (2022)
42. 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. NeurIPS (2022)
43. Schuhmann, C., Beaumont, R., Vencu, R., Gordon, C., Wightman, R., Cherti, M., Coombes, T., Katta, A., Mullis, C., Wortsman, M., et al.: Laion-5b: An open large-scale dataset for training next generation image-text models. NeurIPS (2022)
44. Serra, A., Carrara, F., Tesconi, M., Falchi, F.: The emotions of the crowd: Learning image sentiment from tweets via cross-modal distillation. ECAI (2023)
45. Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., Ganguli, S.: Deep unsupervised learning using nonequilibrium thermodynamics. In: ICML (2015)
46. Song, J., Meng, C., Ermon, S.: Denoising diffusion implicit models. arXiv preprint arXiv:2010.02502 (2020)
47. Tesauro, G., Das, R., Chan, H., Kephart, J., Levine, D., Rawson, F., Lefurgy, C.: Managing power consumption and performance of computing systems using reinforcement learning. NeurIPS (2007)
48. Tu, Z., Talebi, H., Zhang, H., Yang, F., Milanfar, P., Bovik, A., Li, Y.: Maxvit: Multi-axis vision transformer. In: ECCV (2022)
49. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. NeurIPS (2017)
50. Wu, X., Sun, K., Zhu, F., Zhao, R., Li, H.: Human preference score: Better aligning text-to-image models with human preference. In: ICCV (2023)
51. Xu, J., Liu, X., Wu, Y., Tong, Y., Li, Q., Ding, M., Tang, J., Dong, Y.: Imagereward: Learning and evaluating human preferences for text-to-image generation. arXiv preprint arXiv:2304.05977 (2023)
52. Ying, Z., Niu, H., Gupta, P., Mahajan, D., Ghadiyaram, D., Bovik, A.: From patches to pictures (paq-2-piq): Mapping the perceptual space of picture quality. In: CVPR (2020)
53. 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 (2022)
54. Yu, L., Cheng, Y., Sohn, K., Lezama, J., Zhang, H., Chang, H., Hauptmann, A.G., Yang, M.H., Hao, Y., Essa, I., et al.: Magvit: Masked generative video transformer. In: CVPR (2023)
55. Zhou, Y., Liu, B., Zhu, Y., Yang, X., Chen, C., Xu, J.: Shifted diffusion for text-to-image generation. In: CVPR (2023)

# Parrot: Pareto-optimal Multi-Reward Reinforcement Learning Framework for Text-to-Image Generation

## Supplementary Material

This supplementary material provides:

- Sec. A: implementation details, including the training details, and details of quantitative experiments.
- Sec. B: more ablation studies on original prompt guidance and training scheme of Parrot.
- Sec. C: more visual examples to show the advancements of Parrot.

### A. Implementation Details

**Training Details.** We conduct our experiments with Jax implementation of Stable Diffusion 1.5 [39]. In terms of diffusion-based RL, we sample 256 images per RL-tuning iteration. For policy gradient updates, we accumulate gradients across all denoising timesteps. Our experiments employ a small range of gradient clip  $10^{-4}$ . We keep negative prompt as null text.

**Details of Quantitative Experiments.** From Parti [53] prompts, we generate images of dimensions  $512 \times 512$ . In all experiments in the main paper, we apply reward-specific preference expressed as “ $\langle \text{reward } 1 \rangle, \langle \text{reward } 2 \rangle, \langle \text{reward } 3 \rangle, \langle \text{reward } 4 \rangle$ ”, which is optional to select one or several rewards. Reward models are aesthetics, human preference, text-image alignment and image sentiment. During the inference stage, the guidance scale is also set as 5.0 for the diffusion sampling process. We employ the AdamW [25] as optimizer with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$  and a weight decay of 0.1.

### B. Prompt Guidance and Training Scheme

**Original Prompt Guidance.** In Fig. 10, we provide additional ablation study of original prompt centered guidance by adjusting the guidance scale  $w_1$  and  $w_2$ , which determines visual changes from original prompt and expanded prompt, respectively. We assigned 5 to both  $w_1$  and  $w_2$  in most of our experiments, which shows better text-image alignment performance in Fig. 10.

**Ablation Study on Training Scheme.** Fig. 11 provides ablation study comparing variation of the Parrot: Stable Diffusion, the prompt expansion network (PEN) tuning only, T2I model fine-tuning only, Parrot without joint optimization, and Parrot. In the second column, without fine-tuning the T2I diffusion model does not lead to significant improvements in terms of texture and composition. Furthermore, the third column demonstrates that fine-tuning the

diffusion model enhances texture and perspective, yet this improvement is hindered by the limited information in the text prompt. We also observe that the quality of images from joint optimization surpasses that of combining decoupled generative models.

### C. More Visual Examples

We show additional visual examples of Parrot in Figs. 12 to 17. Note that generated images from Parrot are improved across multiple-criteria. Fig. 12 highlights examples where the Parrot brings improvements in aesthetics. For example, Parrot effectively addresses issues such as poor cropping in the fourth column and improves color in the fifth column. Fig. 13 presents examples of images with improved human preference score generated by Parrot. In Fig. 14, we provide examples of improved text-image alignment achieved by Parrot. Fig. 15 shows examples where Parrot enhances image sentiment, producing emotionally rich images.

Finally, additional comparison results between diffusion-based RL baselines are described in Fig. 16, and Fig. 17. Diffusion-based RL baselines are listed: Stable Diffusion 1.5 [39], DPOK [12] with weighted sum of multiple reward scores, Promptist [11], Parrot without prompt expansion, and Parrot. For Parrot without prompt expansion, we only take original prompt as input.



**Fig. 10:** Original prompt centered guidance. We present visual comparison of 5 different pairs of  $w_1$  and  $w_2$  to demonstrate the effectiveness of guidance scales. For all experiments, we assign  $w_1 = 5$  and  $w_2 = 5$  (3rd row) for the best performance.



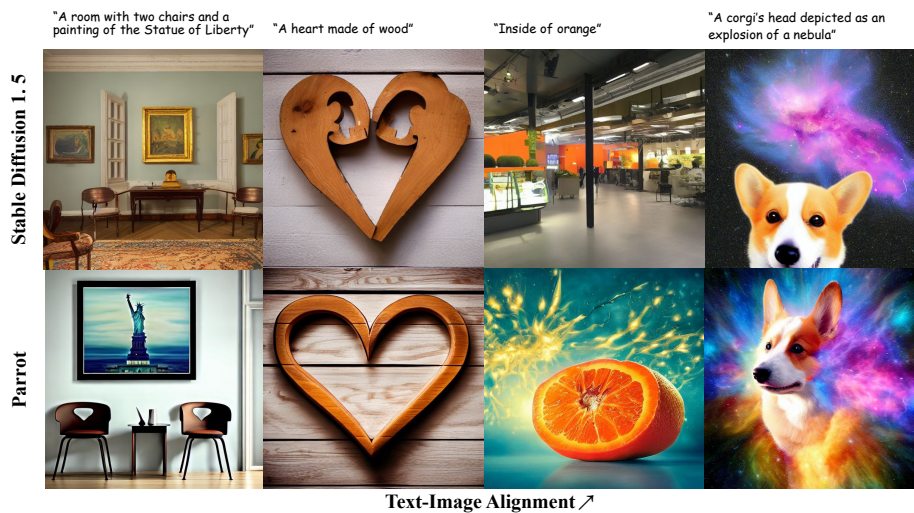
**Fig. 11:** Visual examples of **Parrot** under different settings. From left to right, we provide results of Stable Diffusion 1.5 [39], the only fine-tuned PEN, the only fine-tuned T2I diffusion model, the **Parrot** without joint optimization, and the **Parrot**.



**Fig. 12:** More Examples of aesthetics improvement from the **Parrot**. Given the text prompt, we generate images with Stable Diffusion and **Parrot**. After fine-tuning, the **Parrot** alleviates quality issues such as poor composition (e.g. bad cropping), misalignment with the user input (e.g. missing objects), or generally less aesthetic pleasing.



**Fig. 13:** More Examples of human preference improvement from the **Parrot**. Given the text prompt, we generate images with Stable Diffusion 1.5 [39] and **Parrot**.



**Fig. 14:** More examples of text-image alignment improvement from the **Parrot**. Given the text prompt, we generate images with the Stable Diffusion 1.5 [39] and the **Parrot**.

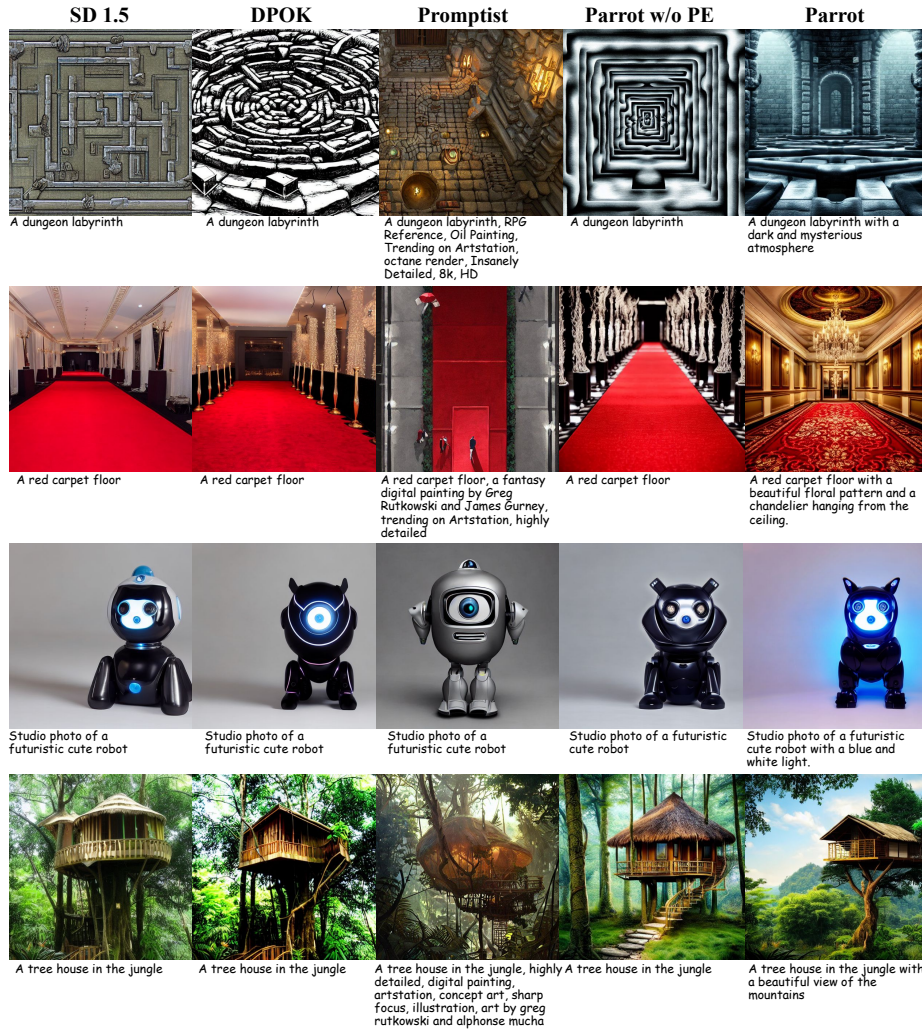


**Fig. 15:** More examples of image sentiment improvement from the **Parrot**. Given the text prompt, we generate images with the Stable Diffusion 1.5 [39] and the **Parrot**.



**Fig. 16:** More results from the **Parrot** and baselines: Stable Diffusion 1.5 [39], DPOK [12] with weighted sum, Promptist [11], Parrot without prompt expansion, and Parrot.





**Fig. 17:** More results from the **Parrot** and baselines: Stable Diffusion 1.5 [39], DPOK [12] with weighted sum, Promptist [11], Parrot without prompt expansion, and Parrot.