
Aligning Text-to-Image Models using Human Feedback

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Deep generative models have shown impressive results in text-to-image synthesis.
2 However, current text-to-image models often generate images that are poorly
3 aligned with text prompts. We propose a fine-tuning method for aligning such
4 models using human feedback, comprising three stages. First, we collect human
5 feedback that assesses model output alignment over a set of diverse text prompts.
6 We then use the human-labeled image-text dataset to train a reward function
7 that predicts human feedback. Lastly, the text-to-image model is fine-tuned by
8 maximizing *reward-weighted* likelihood to improve image-text alignment. Our
9 method generates objects with specified colors, counts, and backgrounds more
10 accurately than the pre-trained model. We also analyze several design choices and
11 show that careful investigation of such design choices is crucial when balancing
12 the alignment-fidelity tradeoff. Our results demonstrate the potential for learning
13 from human feedback to significantly improve text-to-image models.

14 1 Introduction

15 Deep generative models have recently shown remarkable success in generating high-quality images
16 from text prompts [33, 34, 37, 51, 35, 6]. This success has been driven in part by the scaling of deep
17 generative models to large-scale datasets from the web such as LAION [39, 40]. However, major
18 challenges remain in domains where large-scale text-to-image models fail to generate images that are
19 well-aligned with text prompts [9, 25, 26]. For instance, current text-to-image models often fail to
20 produce reliable visual text [26] and struggle with *compositional* image generation [9].

21 In language modeling, *learning from human feedback* has emerged as a powerful technique for
22 aligning model behavior with human intent [53, 43, 46, 29, 30, 3]. Such methods first learn a *reward*
23 *function* intended to reflect what humans care about in the task—they do so by exploiting *human*
24 *feedback* on model outputs. The language model is then optimized using the learned reward function
25 by a *reinforcement learning (RL)* algorithm, such as proximal policy optimization (PPO [41]). This
26 *RL with human feedback (RLHF)* framework has successfully aligned large-scale language models
27 (e.g., GPT-3 [5]) with complex human quality assessments.

28 Motivated by the success of RLHF in language domains, we propose a fine-tuning method for *aligning*
29 *text-to-image models using human feedback*. Our method consists of the three main stages, illustrated
30 in Figure 1: (1) We first generate diverse images from a set of text prompts designed to test output
31 alignment of a text-to-image model. Specifically, we examine prompts where pre-trained models
32 are more prone to error—generating objects with specific colors, counts, and backgrounds. We then
33 collect binary human feedback assessing model outputs. (2) Using this human-labeled dataset, we
34 train a reward function to predict human feedback given the image and text prompt. We propose an
35 auxiliary task—identifying the original text prompt from a set of *perturbed* text prompts—to more
36 effectively exploit human feedback for reward learning. This technique improves the generalization

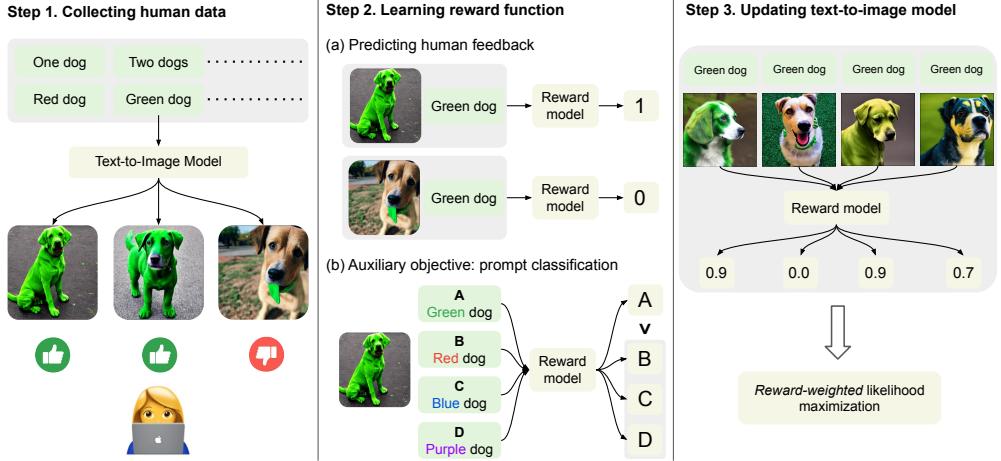


Figure 1: The three stages of our text-to-image fine-tuning method. (1) Multiple images are sampled from the model using the same text prompt, followed by collection of (binary) human feedback. (2) A reward function is learned from these assessments to predict image-text alignment. We develop and exploit *prompt classification*, an auxiliary objective which identifies the original text prompt from among a set of *perturbed* text prompts. (3) We update the text-to-image model via reward-weighted likelihood maximization.

37 of the reward function to unseen images and text prompts. (3) We update the text-to-image model
 38 via *reward-weighted likelihood maximization* to better align it with human feedback. Unlike prior
 39 work that uses RL for optimization [43, 30], we update the model using semi-supervised learning to
 40 measure model-output quality w.r.t. the learned reward function.

41 We illustrate our approach by fine-tuning the *stable diffusion* model [35] using 27K image-text
 42 pairs with human feedback. Our fine-tuned model shows improvement in generating objects with
 43 specified colors, counts, and backgrounds. Moreover, it improves compositional generation (i.e.,
 44 better generates unseen objects¹ given unseen combinations of color, count, and background prompts).
 45 We also observe that the learned reward function is better aligned with human assessments of
 46 alignment than the CLIP score [31] on tested text prompts. We analyze several design choices, such
 47 as using an auxiliary loss for reward learning, and the effect of using “diverse” datasets for fine-tuning.
 48 Our main contributions are as follows:

- 49 • We propose a simple yet efficient fine-tuning method for aligning a text-to-image model using
 50 human feedback.
- 51 • We show that fine-tuning with human feedback significantly improves the image-text alignment
 52 of a text-to-image model. On human evaluation, our model achieves up to 47% improvement in
 53 image-text alignment at the expense of mildly degraded image fidelity.
- 54 • We show that our learned reward function predicts human assessments of the quality more
 55 accurately than the CLIP score [31]. Moreover, we demonstrate that rejection sampling based on
 56 our learned reward function can significantly improve the image-text alignment.
- 57 • Naive fine-tuning with human feedback can significantly reduce the image fidelity, despite better
 58 alignment. We find that the careful investigation of several key design choices is vital in balancing
 59 the alignment-fidelity tradeoff.

60 Though our results do not address all the failure modes of the existing text-to-image models, our
 61 work highlights the potential of learning from human feedback for improving and better aligning
 62 these models.

¹“Unseen objects” are those included in the pre-training dataset but not in the human evaluation dataset.

63 **2 Related Work**

64 We briefly review related work on text-to-image models and human-feedback learning.

65 **Text-to-image models.** Various deep generative models, such as variational auto-encoders [16],
66 generative adversarial networks [11], auto-regressive models [44], and diffusion models [42, 14]
67 have been proposed for image distributions. Combined with the large-scale language encoders [31],
68 [32], these models have shown impressive results in text-to-image generation [34, 37, 51, 35, 6].
69 However, text-to-image models frequently struggle to generate images that are well-aligned with
70 text prompts [9, 25, 26]. Liu et al. [26] show that current models fail to produce reliable visual
71 text and often perform poorly w.r.t. compositional generation [9, 25]. Several techniques, such as
72 character-aware text encoders [49, 26] and structured representations of language inputs [9], have
73 been investigated to address these issues. We study learning from human feedback, which aligns
74 text-to-image models directly using human feedback on model outputs.

75 Fine-tuning with few images [36, 20, 10] for personalization of text-to-image diffusion models is also
76 related to our work. The DreamBooth approach [36] shows that text-to-image models can generate
77 diverse images in a personalized way by fine-tuning with just a few images, while Kumari et al. [20]
78 propose a more memory and computationally efficient method. In this work, we demonstrate that one
79 can fine-tune text-to-image models using simple binary (good/bad) human feedback.

80 **Learning from human feedback.** Human feedback has been used to improve various AI systems
81 (e.g., translation [2, 18], web question-answering [29], story generation [52], training RL agents
82 without hand-designed rewards [28, 7, 45, 15, 21]); and to induce more truthful, less harmful
83 instruction-following and dialogue [30, 3, 53, 46, 43, 24, 38, 4]. Increasingly, learning from human
84 feedback has been used to improve language models and game agents. Several lines of work
85 concurrent with ours have collected human-feedback datasets, trained reward functions [48, 47, 17],
86 and shown that learned rewards are better-aligned with human evaluation than existing score functions
87 (e.g., CLIP [31], BLIP [22]). We demonstrate that certain aspects of text-to-image models can be
88 improved by fine-tuning the model using a reward function trained on human feedback.

89 **3 Main Method**

90 To improve the alignment of generated images with their text prompts, we fine-tune a pre-trained
91 text-to-image model [34, 37, 35] by repeating the following steps (see Figure 1). We first generate a
92 set of diverse images from a collection of text prompts designed to test various capabilities of the
93 text-to-image model. Human raters provide binary feedback on these images (Section 3.1). Next
94 we train a *reward model* to predict human feedback given a text prompt and an image as inputs
95 (Section 3.2). Finally, we fine-tune the text-to-image model using *reward-weighted* log likelihood to
96 improve text-image alignment (Section 3.3).

97 **3.1 Human Data Collection**

98 **Image-text dataset.** To test specific capabilities of a given text-to-image model, we consider three
99 categories of text prompts that generate objects with a specified count, color, or background.² For
100 each category, we generate prompts by combining a word or phrase from that category with some
101 object; e.g., combining green (or in a city) with dog. We also consider combinations of the
102 three categories (e.g., two green dogs in a city). From each prompt, we generate up to
103 60 images using a pre-trained text-to-image model—in this work, we use Stable Diffusion v1.5 [35].

104 **Human feedback.** We collect simple binary feedback from multiple human labelers on the image-
105 text dataset. Labelers are presented with three images generated from the same prompt and are asked
106 to assess whether each image is well-aligned with the prompt (“good”) or not (“bad”).³ We use
107 binary feedback given the simplicity of our prompts—the evaluation criterion are fairly clear. More
108 informative human feedback, such as ranking [43, 30], should prove useful when more complex or
109 subjective text prompts are used (e.g., artistic or open-ended generation).

²For simplicity, we consider a limited class of text categories in this work, deferring the study of broader and more complex categories to future work.

³Labelers are instructed to skip a query if it is hard to answer. Skipped queries are not used in training.

110 **3.2 Reward Learning**

111 To measure image-text alignment, we learn a *reward function* $r_\phi(\mathbf{x}, \mathbf{z})$, parameterized by ϕ , that
 112 maps the CLIP embeddings [31] of an image \mathbf{x} and a text prompt \mathbf{z} to a scalar value.⁴ It is trained to
 113 predict human feedback $y \in \{0, 1\}$ (1 = good, 0 = bad). Formally, given the human feedback dataset
 114 $\mathcal{D}^{\text{human}} = \{(\mathbf{x}, \mathbf{z}, y)\}$, r_ϕ is trained by minimizing the mean-squared-error (MSE):

$$\mathcal{L}^{\text{MSE}}(\phi) = \mathbb{E}_{(\mathbf{x}, \mathbf{z}, y) \sim \mathcal{D}^{\text{human}}} [(y - r_\phi(\mathbf{x}, \mathbf{z}))^2]. \quad (1)$$

115 **Prompt classification.** Data augmentation can significantly improve the data-efficiency and per-
 116 formance of learning [19, 8]. To effectively exploit the feedback dataset, we design a simple data
 117 augmentation scheme and auxiliary loss for reward learning. For each image-text pair that has been
 118 labeled *good*, we generate $N - 1$ text prompts with different semantics than the original text prompt.
 119 For example, we might generate {Blue dog, . . . , Green dog} given the original prompt Red
 120 dog.⁵ This process generates a dataset $\mathcal{D}^{\text{txt}} = \{(\mathbf{x}, \{\mathbf{z}_j\}_{j=1}^N, i')\}$ with N text prompts $\{\mathbf{z}_j\}_{j=1}^N$,
 121 including the original, for each image \mathbf{x} . We denote by $i' \in [N]$ the index of the original text prompt.
 122 We use the augmented prompts in an auxiliary task, namely, classifying the original prompt for
 123 reward learning. Our prompt classifier uses the reward function r_ϕ as follows:

$$P_\phi(i | \mathbf{x}, \{\mathbf{z}_j\}_{j=1}^N) = \frac{\exp(r_\phi(\mathbf{x}, \mathbf{z}_i)/T)}{\sum_j \exp(r_\phi(\mathbf{x}, \mathbf{z}_j)/T)}, \quad \forall i \in [N],$$

124 where $T > 0$ is a temperature parameter. Our auxiliary loss is then defined as

$$\mathcal{L}^{\text{pc}}(\phi) = \mathbb{E}_{(\mathbf{x}, \{\mathbf{z}_j\}_{j=1}^N, i') \sim \mathcal{D}^{\text{txt}}} [\mathcal{L}^{\text{CE}}(P_\phi(i | \mathbf{x}, \{\mathbf{z}_j\}_{j=1}^N), i')], \quad (2)$$

125 where \mathcal{L}^{CE} is the standard cross-entropy loss. This encourages r_ϕ to produce low values for prompts
 126 with different semantics than the original. Our experiments show this auxiliary loss improves
 127 generalization to unseen images and text prompts. Finally, we define our reward learning loss by
 128 combining the losses in (1) and (2) as

$$\mathcal{L}^{\text{reward}}(\phi) = \mathcal{L}^{\text{MSE}}(\phi) + \lambda \mathcal{L}^{\text{pc}}(\phi), \quad (3)$$

129 where λ is a penalty parameter. The pseudo-code for reward learning is provided in Appendix E.

130 **3.3 Updating the Text-to-Image Model**

131 We use our learned r_ϕ to update the text-to-image model p with parameters θ by minimizing the loss

$$\mathcal{L}(\theta) = \mathbb{E}_{(\mathbf{x}, \mathbf{z}) \sim \mathcal{D}^{\text{model}}} [-r_\phi(\mathbf{x}, \mathbf{z}) \log p_\theta(\mathbf{x} | \mathbf{z})] + \beta \mathbb{E}_{(\mathbf{x}, \mathbf{z}) \sim \mathcal{D}^{\text{pre}}} [-\log p_\theta(\mathbf{x} | \mathbf{z})]. \quad (4)$$

132 In (4), $\mathcal{D}^{\text{model}}$ is the model-generated dataset (i.e., images generated by the text-to-image model on
 133 the tested text prompts), \mathcal{D}^{pre} is the *pre-training dataset*, and β is a penalty parameter. The first
 134 term in (4) minimizes the *reward-weighted* negative log-likelihood (NLL) on $\mathcal{D}^{\text{model}}$. To increase
 135 diversity, we collect an unlabeled dataset $\mathcal{D}^{\text{unlabel}}$ by generating more images from the text-to-image
 136 model, and use both the human-labeled dataset $\mathcal{D}^{\text{human}}$ and the unlabeled dataset $\mathcal{D}^{\text{unlabel}}$ for training,
 137 i.e., $\mathcal{D}^{\text{model}} = \mathcal{D}^{\text{human}} \cup \mathcal{D}^{\text{unlabel}}$. By evaluating the quality of the outputs using a reward function
 138 aligned with the text prompts, this term improves the image-text alignment of the model.

139 Minimizing the loss in (4) in autoregressive models like Parti [51] is relatively straightforward.
 140 However, this minimization is challenging for diffusion models, such as Imagen [37], Stable Diffusion
 141 [35], and Dalle-2 [34], as computing the exact log-likelihood is difficult in these models. In this
 142 work, we minimize the reward-weighted MSE loss [4] to fine-tune diffusion-based text-to-image
 143 models (Stable Diffusion in our experiments).

144 Typically, the diversity of the model-generated dataset $\mathcal{D}^{\text{model}}$ is limited, which can result in overfitting.
 145 To mitigate this, similar to Ouyang et al. [30], we also minimize the *pre-training loss*, the second

⁴To improve generalization ability, we use CLIP embeddings pre-trained on various image-text samples.

⁵We use a rule-based strategy to generate different text prompts (see Appendix E for more details).

Category	Total # of images	Human feedback (%)			Examples
		Good	Bad	Skip	
Count	6480	34.4	61.0	4.6	One dog; Two dogs; Three dogs; Four dogs; Five dogs;
Color	3480	70.4	20.8	8.8	A green colored dog; A red colored dog;
Background	2400	66.9	33.1	0.0	A dog in the forest; A dog on the moon;
Combination	15168	35.8	59.9	4.3	Two blue dogs in the forest; Five white dogs in the city;
Total	27528	46.5	48.5	5.0	

Table 1: Details of image-text datasets and human feedback.

146 term in (4). This reduces NLL on the pre-training dataset \mathcal{D}^{pre} . In our experiments, we observed that
147 regularization in the loss function $\mathcal{L}(\theta)$ in (4) enables the generation of more natural images.

148 Different objective functions and algorithms (e.g., PPO; [41]) could be considered for updating the
149 text-to-image model, much like RLHF fine-tuning [30]. We believe RLHF fine-tuning may lead to
150 better models because it uses online sample generation during updates and KL-regularization over
151 the prior model. However, RL usually requires extensive hyperparameter tuning and engineering,
152 thus, we defer proper investigation of RLHF fine-tuning in text-to-image models to future work.

153 4 Experiments

154 In this section, we describe a set of experiments designed to test the efficacy of our approach to
155 fine-tune text-to-image models with human feedback and report our findings and results.

156 4.1 Experimental Setup

157 **Models.** For our baseline generative model, we use Stable Diffusion v1.5 [35], which has been
158 pre-trained on large image-text datasets [39, 40].⁶ For fine-tuning, we freeze the CLIP language
159 encoder [31] and fine-tune only the diffusion module. For the reward model, we use ViT-L/14 CLIP
160 model [31] to extract image and text embeddings and train an MLP using these embeddings as input.
161 Experimental details (e.g., model architectures, hyperparameters) are given in Appendix B.

162 **Datasets.** From a set of 2700 English prompts (see Table I for some instances), we generate 27K
163 images using the stable diffusion model (see Appendix C for further details). Table I shows the
164 feedback distribution provided by multiple human labelers, which has been class-balanced. We note
165 that the stable diffusion model struggles to generate the number of objects specified by the prompt,
166 but reliably generates specified colors and backgrounds.

167 We use 23K samples for training and leave the remaining samples for validation. We also use 16K
168 unlabeled samples for the reward-weighted loss and a 625K subset⁷ of LAION-5B [40] filtered by an
169 aesthetic score predictor⁸ for the pre-training loss.

170 4.2 Text-Image Alignment Results

171 **Human evaluation.** We measure human ratings of image alignment with 120 text prompts (60
172 seen text prompts and 60 unseen text prompts), testing the ability of the models to render different
173 colors, number of objects, and backgrounds (see Appendix C for the full set of prompts). Given two
174 (anonymized) sets of images, one from our fine-tuned model and one from the stable diffusion model,

⁶Our fine-tuning method can be used readily with other text-to-image models, such as Imagen [37], Parti [51], and Dalle-2 [34].

⁷https://huggingface.co/datasets/ChristophSchuhmann/improved_aesthetics_6.5plus

⁸<https://github.com/christophschuhmann/improved-aesthetic-predictor>

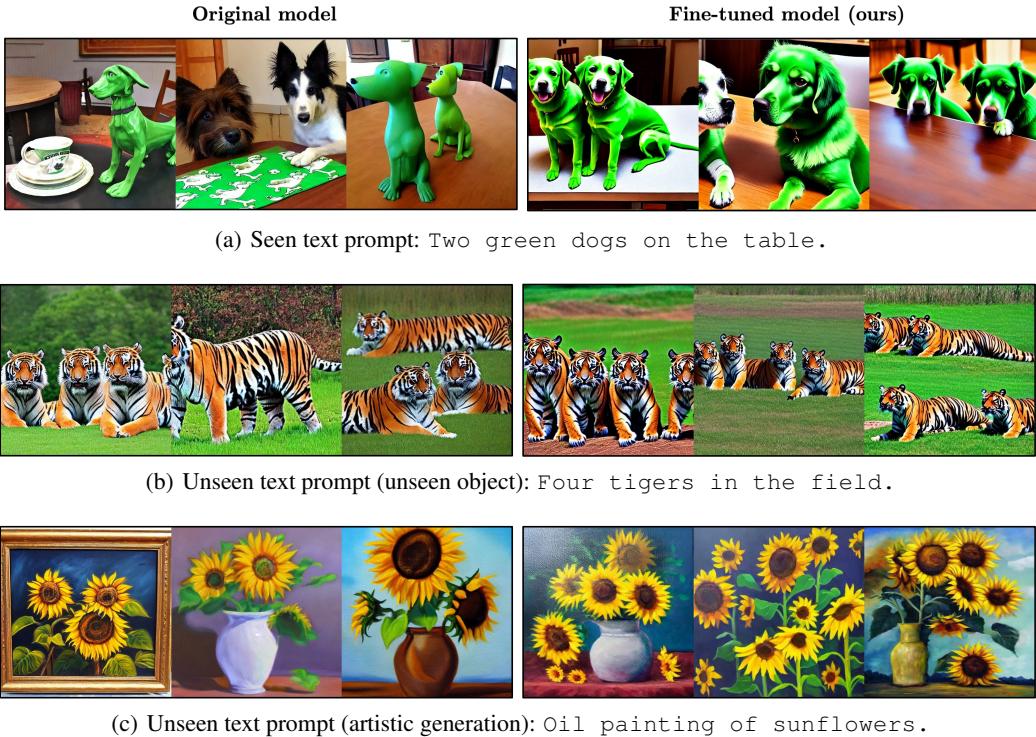


Figure 2: Samples from the original Stable Diffusion model (left) and our fine-tuned model (right). (a) Our model generates a high-quality seen object (dog) with specified color, count, and background. (b) Our model generates an unseen object (tiger) with specified color, count, and background. (c) Our model generates reasonable images even given unseen text categories (artistic generation).

175 we ask human raters to assess which is better w.r.t. image-text alignment and fidelity (i.e., image
 176 quality). We ask raters to declare a *tie* if they have similar quality. Each query is evaluated by 9
 177 independent human raters. We show the percentage of queries based on the number of positive votes.

178 As shown in Figure 3(a), our method significantly improves image-text alignment relative to the
 179 original model. Specifically, 50% of samples from our model receive at least two-thirds of the votes
 180 (7 or more positive votes) for image-text alignment. However, fine-tuning somewhat degrades image
 181 fidelity (15% compared to 10%). We expect that this is because (i) we asked the labelers to provide
 182 feedback mainly on alignment, (ii) the diversity of our human data is limited, and (iii) we used a
 183 small subset of the pre-training dataset for fine-tuning.⁹ This issue can presumably be mitigated with
 184 greater number of raters and a larger pre-training dataset.

185 **Qualitative comparison.** Figure 2 shows image samples from the original model and our fine-tuned
 186 counterpart (see Appendix A for more image examples). While the original often generates images
 187 with missing details (e.g., color, background or count) (Figure 2(a)), our model generates objects that
 188 adhere to the prompt-specified colors, counts, and backgrounds. Of special note, our model generates
 189 high-quality images on unseen text prompts that specify unseen objects (Figure 2(b)). Our model also
 190 generates reasonable images given unseen text categories, such as artistic generation (Figure 2(c)).

191 Despite its strong performance, we do observe some difficulties with our fine-tuned model. For
 192 certain text prompts, our fine-tuned model generates oversaturated and non-photorealistic images.
 193 Our model occasionally duplicates entities within the generated images, or produces images for the
 194 same prompt with low diversity. We believe these issues can be addressed with larger (and more
 195 diverse) human datasets and better optimization (e.g., with RL).

⁹Similar issue, which is akin to the *alignment tax*, has been observed in language domains [1, 30].

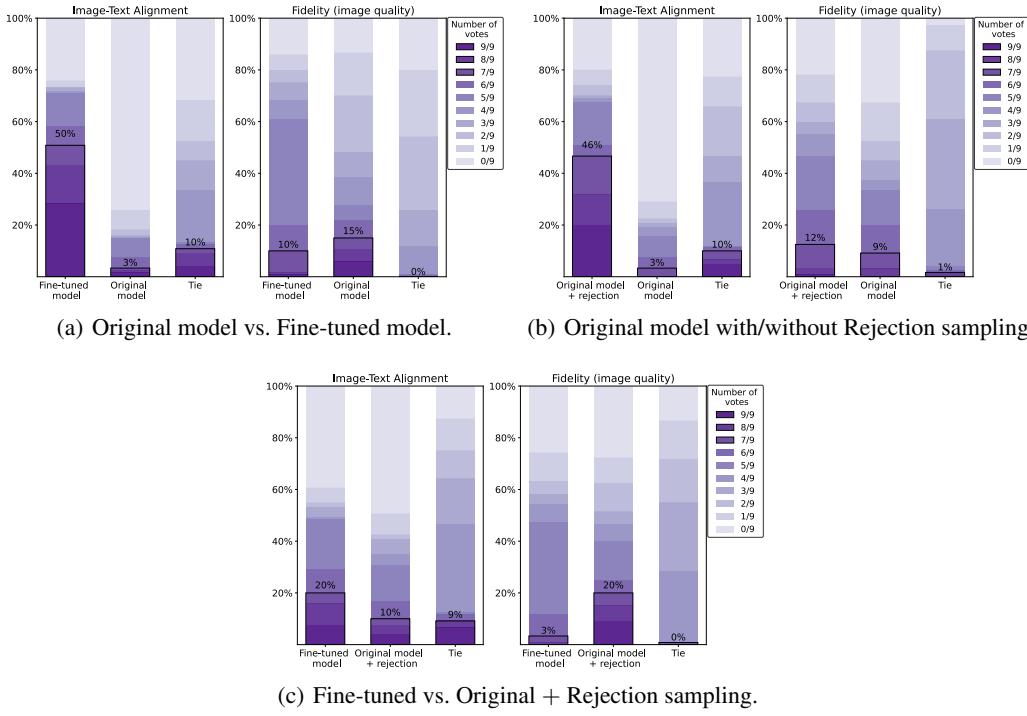


Figure 3: Human evaluation on 120 tested text prompts (60 seen and 60 unseen). We generate two sets of images with the same text prompt. Human raters indicate which set is better (or if it is a tie, i.e., the two sets are similar), in terms of both image-text alignment and image fidelity. Each query is evaluated by nine raters and we report the percentage of queries that receive a given number of positive votes (0–9). We also highlight the percentage of queries receiving at least two-thirds (7 or more) of the positive votes in the black box. (a) Comparison of the fine-tuned and original models. (b) Comparison of the original model with and without the rejection sampling. (c) Comparison of the fine-tuned model and the original model with rejection sampling.

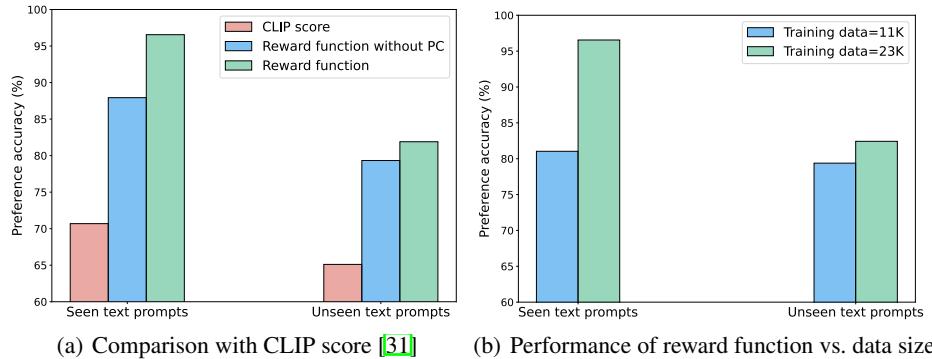


Figure 4: (a) Accuracy of CLIP score [31] and our reward functions when predicting the assessments of human labelers. We consider a variant of our reward function which is *not* trained with the prompt classification (PC) loss in [2]. (b) Performance of reward functions when varying the size of the training dataset.

196 4.3 Results on Reward Learning

197 **Predicting human preferences.** We investigate the quality of our learned reward function by
 198 evaluating its prediction of the assessments of human labelers. Given two images x_1 and x_2 from the
 199 same text prompt z , we check whether our reward r_ϕ generates a higher score for the human-preferred

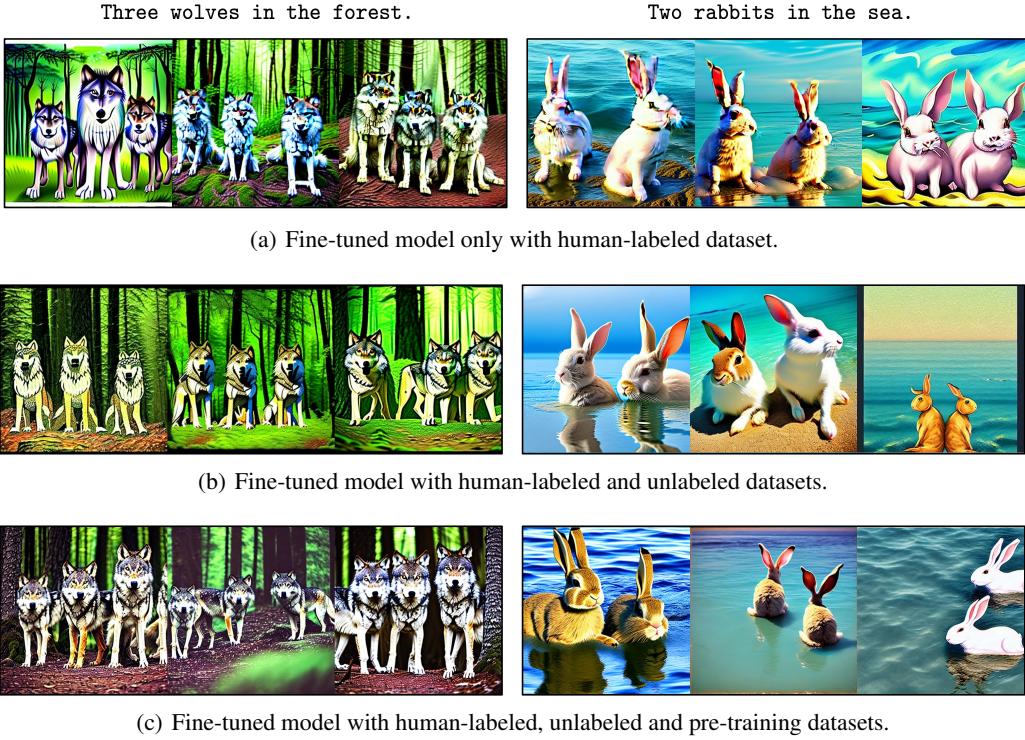


Figure 5: Samples from fine-tuned models trained with different datasets on unseen text prompts. (a) Model fine-tuned with only human data generates low-quality images due to overfitting. (b) Unlabeled samples improve the quality of generated images. (c) Fine-tuned model can generate high-fidelity images by using the pre-training dataset.

200 image, i.e., $r_\phi(\mathbf{x}_1, \mathbf{z}) > r_\phi(\mathbf{x}_2, \mathbf{z})$ when rater prefers \mathbf{x}_1 to \mathbf{x}_2 . As a baseline, we compare it with the
201 CLIP score [12], which measures image-text similarity in the CLIP embedding space [31].

202 Figure 4(a) compares the accuracy of r_ϕ and the CLIP score on unseen images from both seen and
203 unseen text prompts. Our reward (green) more accurately predicts human evaluation than the CLIP
204 score (red), hence is better aligned with typical human assessments. To show the benefit of our
205 auxiliary loss (prompt classification) in 2, we also assess a variant of our reward function which
206 ignores this loss (blue). The results show that the auxiliary classification task improves reward
207 performance on both seen and unseen text prompts. The gain from the auxiliary loss clearly shows the
208 importance of text diversity and our auxiliary loss in improving data efficiency. Although our reward
209 function is more accurate than the CLIP score, its performance on unseen text prompts ($\sim 80\%$)
210 suggests that it may be necessary to use more diverse and larger human datasets.

211 **Rejection sampling.** Similar to Parti [51] and DALL-E [33], we evaluate a rejection sampling
212 technique, which selects the best output w.r.t. the learned reward function¹⁰. Specifically, we generate
213 16 images per text prompt from the original stable diffusion model and select the four with the
214 greatest reward scores. We compare these to four randomly sampled images in Figure 3(b). Rejection
215 sampling significantly improves image-text alignment (46% with two-thirds preference vote by raters)
216 without sacrificing image fidelity. This result illustrates the significance of learning a reward function
217 in improving text-to-image models even *without any fine-tuning*.

218 We also compare our fine-tuned model to the original with rejection sampling in Figure 3(c)¹¹.
219 Our fine-tuned model achieves a 10% gain in image-text alignment (20%-10% two-thirds vote) but
220 sacrifices 17% in image fidelity (3%-20% two-thirds vote). However, as discussed in Section 4.2,

¹⁰Parti and DALL-E use similarity scores of image and text embeddings from CoCa [50] and CLIP [31], respectively.

¹¹We note that rejection sampling can also be applied on top of our fine-tuned model.

	Original model	Fine-tuned model w.o unlabeled & pre-train	Fine-tuned model w.o pre-train	Fine-tuned model
FID on MS-CoCo (\downarrow)	13.97	26.59	21.02	16.76
Average rewards on tested prompts (\uparrow)	0.43	0.69	0.79	0.79

Table 2: Comparison with the original Stable Diffusion. To evaluate image fidelity, we measure FID scores on MS-CoCo. To evaluate image-text alignment, we measure reward scores and CLIP scores on 120 tested text prompts. (\uparrow (resp., \downarrow) indicates that higher (resp., lower) metrics are better.)

221 we expect degradation in fidelity to be mitigated with larger human datasets and better hyper-
 222 parameters. Note also that rejection sampling has several drawbacks, including increased inference-
 223 time computation and the inability to improve the model (since it is a post-processing technique).

224 4.4 Ablation Studies

225 **Effects of the size of the human dataset.** To investigate how the quality of human data affects
 226 reward learning, we conduct an ablation study, reducing the number of images per text prompt by
 227 half before training the reward function. Figure 4(b) shows that model accuracy decreases on both
 228 seen and unseen prompts as data size decreases, clearly demonstrating the importance of diversity
 229 and the size of the rater data.

230 **Effects of using diverse datasets.** To verify the importance of data diversity, we incrementally
 231 include unlabeled and pre-training datasets during fine-tuning. We measure the reward score (image-
 232 text alignment) on 120 tested text prompts and FID score [13]—the similarity between generated
 233 images and real images—on MS-CoCo validation data [23]. Table 2 shows that FID score is
 234 significantly reduced when the model is fine-tuned using only human data, despite better image-text
 235 alignment. However, by adding the unlabeled and pre-training datasets, FID score is improved
 236 without impacting the image-text alignment. We provide image samples from unseen text prompts in
 237 Figure 5. We observe that fine-tuned models generate more natural images when exploiting more
 238 diverse datasets.

239 5 Discussion

240 In this work, we showed that fine-tuning with human feedback can effectively improve the image-text
 241 alignment in three domains: generating objects with a specified count, color, or backgrounds. We
 242 analyzed several design choices (such as using an auxiliary loss and collecting diverse training data)
 243 and found that it is challenging to balance the alignment-fidelity tradeoff without careful investigation
 244 of such design choices. While our results do not address all the failure modes of existing text-to-image
 245 models, our work serves as a starting point for the deeper study of learning from human feedback for
 246 improving text-to-image models.

247 **Limitations, future directions.** Several limitations of our work suggest interesting future directions:

248 *More nuanced human feedback.* Some of the poor image generation we observe (e.g., highly saturated
 249 image colors) are likely due to similar images being highly ranked in our training set. Instructing
 250 raters to watch for a more diverse set of failure modes (e.g., oversaturated colors, unrealistic animal
 251 anatomy, physics violations) should improve performance along these dimensions.

252 *Different objectives and algorithms.* To update a text-to-image model, we used a reward-weighted
 253 likelihood maximization. However, similar to prior work in language domains [30], it would be an
 254 interesting direction to use RL [41]. We expect RLHF fine-tuning to lead to better models because (a)
 255 it uses online sample generation during updates and (b) KL-regularization over the prior model can
 256 mitigate overfitting to the reward function.

257 **References**

- 258 [1] Askell, Amanda, Bai, Yuntao, Chen, Anna, Drain, Dawn, Ganguli, Deep, Henighan, Tom, Jones,
259 Andy, Joseph, Nicholas, Mann, Ben, DasSarma, Nova, et al. A general language assistant as a
260 laboratory for alignment. *arXiv preprint arXiv:2112.00861*, 2021.
- 261 [2] Bahdanau, Dzmitry, Brakel, Philemon, Xu, Kelvin, Goyal, Anirudh, Lowe, Ryan, Pineau, Joelle,
262 Courville, Aaron, and Bengio, Yoshua. An actor-critic algorithm for sequence prediction. *arXiv
263 preprint arXiv:1607.07086*, 2016.
- 264 [3] Bai, Yuntao, Jones, Andy, Ndousse, Kamal, Askell, Amanda, Chen, Anna, DasSarma, Nova,
265 Drain, Dawn, Fort, Stanislav, Ganguli, Deep, Henighan, Tom, et al. Training a helpful
266 and harmless assistant with reinforcement learning from human feedback. *arXiv preprint
267 arXiv:2204.05862*, 2022.
- 268 [4] Bai, Yuntao, Kadavath, Saurav, Kundu, Sandipan, Askell, Amanda, Kernion, Jackson, Jones,
269 Andy, Chen, Anna, Goldie, Anna, Mirhoseini, Azalia, McKinnon, Cameron, et al. Constitutional
270 ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- 271 [5] Brown, Tom B, Mann, Benjamin, Ryder, Nick, Subbiah, Melanie, Kaplan, Jared, Dhariwal,
272 Prafulla, Neelakantan, Arvind, Shyam, Pranav, Sastry, Girish, Askell, Amanda, et al. Language
273 models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- 274 [6] Chang, Huiwen, Zhang, Han, Barber, Jarred, Maschinot, AJ, Lezama, Jose, Jiang, Lu, Yang,
275 Ming-Hsuan, Murphy, Kevin, Freeman, William T, Rubinstein, Michael, et al. Muse: Text-
276 to-image generation via masked generative transformers. *arXiv preprint arXiv:2301.00704*,
277 2023.
- 278 [7] Christiano, Paul F, Leike, Jan, Brown, Tom, Martic, Miljan, Legg, Shane, and Amodei, Dario.
279 Deep reinforcement learning from human preferences. In *Advances in Neural Information
280 Processing Systems*, 2017.
- 281 [8] Cubuk, Ekin D, Zoph, Barret, Mane, Dandelion, Vasudevan, Vijay, and Le, Quoc V. Autoaug-
282 ment: Learning augmentation strategies from data. In *Proceedings of the IEEE/CVF Conference
283 on Computer Vision and Pattern Recognition*, 2019.
- 284 [9] Feng, Weixi, He, Xuehai, Fu, Tsu-Jui, Jampani, Varun, Akula, Arjun, Narayana, Pradyumna,
285 Basu, Sugato, Wang, Xin Eric, and Wang, William Yang. Training-free structured diffusion
286 guidance for compositional text-to-image synthesis. *arXiv preprint arXiv:2212.05032*, 2022.
- 287 [10] Gal, Rinon, Alaluf, Yuval, Atzmon, Yuval, Patashnik, Or, Bermano, Amit H, Chechik, Gal, and
288 Cohen-Or, Daniel. An image is worth one word: Personalizing text-to-image generation using
289 textual inversion. *arXiv preprint arXiv:2208.01618*, 2022.
- 290 [11] Goodfellow, Ian, Pouget-Abadie, Jean, Mirza, Mehdi, Xu, Bing, Warde-Farley, David, Ozair,
291 Sherjil, Courville, Aaron, and Bengio, Yoshua. Generative adversarial networks. *Communica-
292 tions of the ACM*, 63(11):139–144, 2020.
- 293 [12] Hessel, Jack, Holtzman, Ari, Forbes, Maxwell, Bras, Ronan Le, and Choi, Yejin. Clipscore: A
294 reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021.
- 295 [13] Heusel, Martin, Ramsauer, Hubert, Unterthiner, Thomas, Nessler, Bernhard, and Hochreiter,
296 Sepp. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In
297 *Advances in neural information processing systems*, 2017.
- 298 [14] Ho, Jonathan, Jain, Ajay, and Abbeel, Pieter. Denoising diffusion probabilistic models. In
299 *Advances in Neural Information Processing Systems*, 2020.
- 300 [15] Ibarz, Borja, Leike, Jan, Pohlen, Tobias, Irving, Geoffrey, Legg, Shane, and Amodei, Dario.
301 Reward learning from human preferences and demonstrations in atari. In *Advances in Neural
302 Information Processing Systems*, 2018.
- 303 [16] Kingma, Diederik P and Welling, Max. Auto-encoding variational bayes. *arXiv preprint
304 arXiv:1312.6114*, 2013.

- 305 [17] Kirstain, Yuval, Polyak, Adam, Singer, Uriel, Matiana, Shahbuland, Penna, Joe, and Levy, Omer.
306 Pick-a-pic: An open dataset of user preferences for text-to-image generation. *arXiv preprint*
307 *arXiv:2305.01569*, 2023.
- 308 [18] Kreutzer, Julia, Khadivi, Shahram, Matusov, Evgeny, and Riezler, Stefan. Can neural machine
309 translation be improved with user feedback? *arXiv preprint arXiv:1804.05958*, 2018.
- 310 [19] Krizhevsky, Alex, Sutskever, Ilya, and Hinton, Geoffrey E. Imagenet classification with deep
311 convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017.
- 312 [20] Kumari, Nupur, Zhang, Bingliang, Zhang, Richard, Shechtman, Eli, and Zhu, Jun-Yan. Multi-
313 concept customization of text-to-image diffusion. *arXiv preprint arXiv:2212.04488*, 2022.
- 314 [21] Lee, Kimin, Smith, Laura, and Abbeel, Pieter. Pebble: Feedback-efficient interactive rein-
315forcement learning via relabeling experience and unsupervised pre-training. In *International*
316 *Conference on Machine Learning*, 2021.
- 317 [22] Li, Junnan, Li, Dongxu, Xiong, Caiming, and Hoi, Steven. Blip: Bootstrapping language-
318 image pre-training for unified vision-language understanding and generation. In *International*
319 *Conference on Machine Learning*, 2022.
- 320 [23] Lin, Tsung-Yi, Maire, Michael, Belongie, Serge, Hays, James, Perona, Pietro, Ramanan, Deva,
321 Dollár, Piotr, and Zitnick, C Lawrence. Microsoft coco: Common objects in context. In *European conference on computer vision*, 2014.
- 323 [24] Liu, Hao, Sferrazza, Carmelo, and Abbeel, Pieter. Chain of hindsight aligns language models
324 with feedback. *arXiv preprint arXiv: Arxiv-2302.02676*, 2023.
- 325 [25] Liu, Nan, Li, Shuang, Du, Yilun, Torralba, Antonio, and Tenenbaum, Joshua B. Compositional
326 visual generation with composable diffusion models. *arXiv preprint arXiv:2206.01714*, 2022.
- 327 [26] Liu, Rosanne, Garrette, Dan, Saharia, Chitwan, Chan, William, Roberts, Adam, Narang, Sharan,
328 Blok, Irina, Mical, RJ, Norouzi, Mohammad, and Constant, Noah. Character-aware models
329 improve visual text rendering. *arXiv preprint arXiv:2212.10562*, 2022.
- 330 [27] Loshchilov, Ilya and Hutter, Frank. Decoupled weight decay regularization. *arXiv preprint*
331 *arXiv:1711.05101*, 2017.
- 332 [28] MacGlashan, James, Ho, Mark K, Loftin, Robert, Peng, Bei, Roberts, David, Taylor, Matthew E,
333 and Littman, Michael L. Interactive learning from policy-dependent human feedback. In *International Conference on Machine Learning*, 2017.
- 335 [29] Nakano, Reiichiro, Hilton, Jacob, Balaji, Suchir, Wu, Jeff, Ouyang, Long, Kim, Christina,
336 Hesse, Christopher, Jain, Shantanu, Kosaraju, Vineet, Saunders, William, et al. Webgpt:
337 Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*,
338 2021.
- 339 [30] Ouyang, Long, Wu, Jeff, Jiang, Xu, Almeida, Diogo, Wainwright, Carroll L, Mishkin, Pamela,
340 Zhang, Chong, Agarwal, Sandhini, Slama, Katarina, Ray, Alex, et al. Training language models
341 to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*, 2022.
- 342 [31] Radford, Alec, Kim, Jong Wook, Hallacy, Chris, Ramesh, Aditya, Goh, Gabriel, Agarwal,
343 Sandhini, Sastry, Girish, Askell, Amanda, Mishkin, Pamela, Clark, Jack, et al. Learning
344 transferable visual models from natural language supervision. In *International Conference on*
345 *Machine Learning*, 2021.
- 346 [32] Raffel, Colin, Shazeer, Noam, Roberts, Adam, Lee, Katherine, Narang, Sharan, Matena,
347 Michael, Zhou, Yanqi, Li, Wei, Liu, Peter J, et al. Exploring the limits of transfer learning with
348 a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21(140):1–67, 2020.
- 349 [33] Ramesh, Aditya, Pavlov, Mikhail, Goh, Gabriel, Gray, Scott, Voss, Chelsea, Radford, Alec,
350 Chen, Mark, and Sutskever, Ilya. Zero-shot text-to-image generation. In *International Conference on*
351 *Machine Learning*, 2021.

- 352 [34] Ramesh, Aditya, Dhariwal, Prafulla, Nichol, Alex, Chu, Casey, and Chen, Mark. Hierarchical
 353 text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- 354 [35] Rombach, Robin, Blattmann, Andreas, Lorenz, Dominik, Esser, Patrick, and Ommer, Björn.
 355 High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF
 356 Conference on Computer Vision and Pattern Recognition*, 2022.
- 357 [36] Ruiz, Nataniel, Li, Yuanzhen, Jampani, Varun, Pritch, Yael, Rubinstein, Michael, and Aberman,
 358 Kfir. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation.
 359 *arXiv preprint arXiv:2208.12242*, 2022.
- 360 [37] Saharia, Chitwan, Chan, William, Saxena, Saurabh, Li, Lala, Whang, Jay, Denton, Emily,
 361 Ghasemipour, Seyed Kamyar Seyed, Ayan, Burcu Karagol, Mahdavi, S Sara, Lopes, Rapha Gon-
 362 tijo, et al. Photorealistic text-to-image diffusion models with deep language understanding. In
 363 *Advances in Neural Information Processing Systems*, 2022.
- 364 [38] Scheurer, Jérémie, Campos, Jon Ander, Chan, Jun Shern, Chen, Angelica, Cho, Kyunghyun,
 365 and Perez, Ethan. Training language models with language feedback. *arXiv preprint arXiv:
 366 Arxiv-2204.14146*, 2022.
- 367 [39] Schuhmann, Christoph, Vencu, Richard, Beaumont, Romain, Kaczmarczyk, Robert, Mullis,
 368 Clayton, Katta, Aarush, Coombes, Theo, Jitsev, Jenia, and Komatsuzaki, Aran. Laion-400m:
 369 Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*,
 370 2021.
- 371 [40] Schuhmann, Christoph, Beaumont, Romain, Vencu, Richard, Gordon, Cade, Wightman, Ross,
 372 Cherti, Mehdi, Coombes, Theo, Katta, Aarush, Mullis, Clayton, Wortsman, Mitchell, et al.
 373 Laion-5b: An open large-scale dataset for training next generation image-text models. *arXiv
 374 preprint arXiv:2210.08402*, 2022.
- 375 [41] Schulman, John, Wolski, Filip, Dhariwal, Prafulla, Radford, Alec, and Klimov, Oleg. Proximal
 376 policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 377 [42] Sohl-Dickstein, Jascha, Weiss, Eric, Maheswaranathan, Niru, and Ganguli, Surya. Deep
 378 unsupervised learning using nonequilibrium thermodynamics. In *International Conference on
 379 Machine Learning*, 2015.
- 380 [43] Stiennon, Nisan, Ouyang, Long, Wu, Jeff, Ziegler, Daniel M, Lowe, Ryan, Voss, Chelsea,
 381 Radford, Alec, Amodei, Dario, and Christiano, Paul. Learning to summarize from human
 382 feedback. *arXiv preprint arXiv:2009.01325*, 2020.
- 383 [44] Van Den Oord, Aäron, Kalchbrenner, Nal, and Kavukcuoglu, Koray. Pixel recurrent neural
 384 networks. In *International conference on machine learning*, pp. 1747–1756, 2016.
- 385 [45] Warnell, Garrett, Waytowich, Nicholas, Lawhern, Vernon, and Stone, Peter. Deep tamer: Inter-
 386 active agent shaping in high-dimensional state spaces. In *Conference on Artificial Intelligence*,
 387 2018.
- 388 [46] Wu, Jeff, Ouyang, Long, Ziegler, Daniel M, Stiennon, Nisan, Lowe, Ryan, Leike, Jan, and
 389 Christiano, Paul. Recursively summarizing books with human feedback. *arXiv preprint
 390 arXiv:2109.10862*, 2021.
- 391 [47] Wu, Xiaoshi, Sun, Keqiang, Zhu, Feng, Zhao, Rui, and Li, Hongsheng. Better aligning
 392 text-to-image models with human preference. *arXiv preprint arXiv:2303.14420*, 2023.
- 393 [48] Xu, Jiazheng, Liu, Xiao, Wu, Yuchen, Tong, Yuxuan, Li, Qinkai, Ding, Ming, Tang, Jie, and
 394 Dong, Yuxiao. Imagereward: Learning and evaluating human preferences for text-to-image
 395 generation. *arXiv preprint arXiv:2304.05977*, 2023.
- 396 [49] Xue, Linting, Barua, Aditya, Constant, Noah, Al-Rfou, Rami, Narang, Sharan, Kale, Mihir,
 397 Roberts, Adam, and Raffel, Colin. Byt5: Towards a token-free future with pre-trained byte-
 398 to-byte models. *Transactions of the Association for Computational Linguistics*, 10:291–306,
 399 2022.

- 400 [50] Yu, Jiahui, Wang, Zirui, Vasudevan, Vijay, Yeung, Legg, Seyedhosseini, Mojtaba, and Wu,
401 Yonghui. Coca: Contrastive captioners are image-text foundation models. *arXiv preprint*
402 *arXiv:2205.01917*, 2022.
- 403 [51] Yu, Jiahui, Xu, Yuanzhong, Koh, Jing Yu, Luong, Thang, Baid, Gunjan, Wang, Zirui, Vasudevan,
404 Vijay, Ku, Alexander, Yang, Yinfei, Ayan, Burcu Karagol, et al. Scaling autoregressive models
405 for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2022.
- 406 [52] Zhou, Wangchunshu and Xu, Ke. Learning to compare for better training and evaluation of open
407 domain natural language generation models. In *Conference on Artificial Intelligence*, 2020.
- 408 [53] Ziegler, Daniel M, Stiennon, Nisan, Wu, Jeffrey, Brown, Tom B, Radford, Alec, Amodei, Dario,
409 Christiano, Paul, and Irving, Geoffrey. Fine-tuning language models from human preferences.
410 *arXiv preprint arXiv:1909.08593*, 2019.