# Supplementary — Enhancing Diffusion Models with Text-Encoder Reinforcement Learning

Chaofeng Chen<sup>1</sup>\*<sup>©</sup>, Annan Wang<sup>1</sup>\*<sup>©</sup>, Haoning Wu<sup>1</sup><sup>©</sup>, Liang Liao<sup>1</sup><sup>©</sup>, Wenxiu Sun<sup>3</sup><sup>©</sup>, Qiong Yan<sup>3</sup><sup>©</sup>, and Weisi Lin<sup>2</sup><sup>©</sup>

<sup>1</sup> S-Lab, Nanyang Technological University <sup>2</sup> CCDS, Nanyang Technological University <sup>3</sup> Sensetime Research chaofeng.chen@ntu.edu.sg c190190@e.ntu.edu.sg wslin@ntu.edu.sg https://github.com/chaofengc/TexForce

# 1 Results and Test Codes

For the convenience of the reviewers, we provide the test codes and results of our approach on the ImageReward test dataset in the following files: Supplementary Material

juppiementary nateriar
results_imagereward
results_imagereward_sd14 # Results with SDv1.4 backbone
results_imagereward_sd15 # Results with SDv1.5 backbone
results_imagereward_sd21 # Results with SDv2.1 backbone
imagereward_benchmark.json # Test prompts from ImageReward
gpt4v_evaluation_results # GPT4V api response
lora_weights # LoRA weights for TexForce and ReFL
gpt4v_compare.py # Code for GPT4V test
demo.html # Html demo to show all results, open
it with chrome
L   test.ipynb # IPython notebook to test our model, including
instructions
supplementary.pdf # This file

# 2 Method Details

### 2.1 Reinforcement Learning Formulations

We first give the formulations of each term for the PPO algorithm in diffusion models. According to the DDPM paper [4], the t step denoising can be formulated as:

$$\mathbf{x}_{t-1} = \mu_{\theta,\phi}(\mathbf{x}_t, t, \tau_{\phi}(s)) + \sigma_t \epsilon, \tag{1}$$

$$\mu_{\theta,\phi} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta,\phi}(\mathbf{x}_t, t, \tau_{\phi}(s)) \right)$$
(2)

<sup>\*</sup> These authors contributed equally to this work.

where  $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i, \sigma_t^2 = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , and  $\beta_t$  is a predefined variance schedule. Since the U-Net parameter  $\phi$  is fixed during the training, we omit it in the following formulations. The objective function of the PPO algorithm is:

$$J(\phi) = \mathbb{E}\left[\min(r_t(\phi)A, \operatorname{clip}(r_t(\phi), 1 - \lambda, 1 + \lambda)A)\right],\tag{3}$$

where

$$r_t(\phi) = \frac{p_\phi(\mathbf{x}_{t-1}|\mathbf{x}_t)}{p_{\phi_{\text{old}}}(\mathbf{x}_{t-1}|\mathbf{x}_t)}.$$
(4)

Since  $p_{\phi}$  is an isotropic Gaussian distribution, we have:

$$r_{t}(\phi) = \exp(\log p_{\phi}(\mathbf{x}_{t-1}|\mathbf{x}_{t}) - \log p_{\phi_{\text{old}}}(\mathbf{x}_{t-1}|\mathbf{x}_{t}))$$
  
= 
$$\exp\left(-\frac{1}{2\sigma_{t}^{2}}\|\mathbf{x}_{t-1} - \mu_{\phi,t}\|^{2} + \frac{1}{2\sigma_{t}^{2}}\|\mathbf{x}_{t-1} - \mu_{\phi_{\text{old}},t}\|^{2}\right).$$
 (5)

The reward value A is obtained with reward function  $A = R(\mathbf{x}_0, s)$ . With equations above, we can get  $J(\phi)$ .

### 2.2 Face and Hand Rewards

Figure 1 shows the pipeline of face reward and hand reward. For the face quality reward, we finetuned the TOPIQ model<sup>3</sup> [1] with the GFIQA-20k dataset [6] which is specifically designed for face quality. Since the faces of GFIQA are all aligned, we also need to align the generated face before calculating the reward scores. For hand reward function, there is no existing quality model for hands. We found that simple hand detection confidence score can already give reliable reward for the generation quality of hands. Therefore, we directly use the hand detection confidence as rewards. Thanks to the flexibility of RL, the reward function is not required to be differentiable. We use the pretrained YOLOv5<sup>4</sup> for hand detection.

### **3** Experiment Settings

#### 3.1 Training Details

As shown in Tab. 1, we use the same hyper-parameters for all the experiments. The only difference is the batch size and the number of samples per epoch. For the single prompt dataset, we use a batch size of 8 and 256 samples per epoch. For the multi-prompts dataset, we use a batch size of 64 and 2048 samples per epoch.

<sup>&</sup>lt;sup>3</sup> https://github.com/chaofengc/IQA-PyTorch

<sup>4</sup> https://github.com/MahmudulAlam/Unified - Gesture - and - Fingertip -Detection



Fig. 1: Calculation of face quality and hand detection reward.

 Table 1: Hyper-parameters and training settings for single prompt and multi-prompts datasets.

	Hyper-parameters	Single prompt	Multiple prompts		
	Sampler	DDIM [5]			
Diffusion	Guidance Scale	7.5			
	Sampling Steps	50			
	Type	А	AdamW		
	Learning rate	3e-4			
Optimizer	Weight decay	1e-4			
	$(eta_1,eta_2)$	(0.9, 0.999)			
	Gradient clip	1.0			
BL Config	Ratio clip $(\gamma)$	1e-4			
TEL Coming	Advantage clip	10			
	Rank	16			
LoRA	Alpha $(\alpha)$	1			
	Module	q, k, v, out			
Training	Trainable $\#$ Params.	1.18M			
	Numerical precision	torch.float16			
	Batch size	8	64		
	Samples per epoch	256	2048		
	Total epochs	100	100		
	GPUs	1 V100	4 A100		
	Time	$\sim 2 \text{ days}$	$\sim 3 \text{ days}$		

### 3.2 Evaluation Details

We use the seed\_everything() function provided by PyTorch Lightning<sup>5</sup> to set the random seed for all the experiments. We set the start random seed to 234 for ALL prompts.

<sup>&</sup>lt;sup>5</sup> https://lightning.ai/

For the single prompt dataset, we generate 50 examples for each prompt. For the multi-prompts dataset, ImageReward and HPSv2, we generate only one example for each prompt to save time.

# 4 GPT4V Evaluation

We use the gpt-4-1106-vision-preview API<sup>6</sup> provided by OpenAI to evaluate the quality of the generated images. The API takes a text prompt and a list of images as input, and returns a rank of the image names based on their aesthetic quality and their coherence with the prompt. Below is an example we used for the evaluation:



<sup>&</sup>lt;sup>6</sup> https://platform.openai.com/docs/guides/vision



Fig. 2: Complete results of GPT4V evaluation with SDv1.4 backbone on the ImageReward dataset.

As shown in the example above, GPT4V returns two lists of the image names, ranking from good to bad. We assign score 3 for the best image, 0 for the worst image and the final score is normalized to [0, 1]. For reliability, we run the evaluation for 3 times and report the average score. The results for each round and the final test score is reported in Fig. 2. Please refer to the demo.html file for the detailed results of each individual test image.

# 5 Additional Experiments

### 5.1 More Results of Incompression Task

In the main paper, we briefly discussed different behaviors of finetuing U-Net and text encoder in the incompression task. Here, we provide more quantitative and qualitative results about the comparison between finetuning U-Net and text encoder, with the unseen animal prompts below:

**Animal prompts.** Following previous works, we use the 45 simple animals for training. The prompt is defined as A photo of a < animal>. We use the following animals for testing: cheetach, elephant, girraffe, hippo, jellyfish, panda, penguin, swan.

We generate 10 samples for each prompt and report their incompression scores in Tab. 2. We can have the following observations:

- Although the training rewards of text encoder and U-Net are similar, the text encoder is more robust in unseen animals than U-Net, and obtained higher incompression scores.
- The quantitative results also confirm that simply combining U-Net and text encoder is quite effective in improving the reward scores.

Figure 3 shows visual examples of combining U-Net and text encoder at different checkpoints. We can observe that the text encoder can introduce extra reasonable visual concepts to the image to increase complexity, however, U-Net mainly changes the appearance and is easy to disrupt the original structure such as the hippo head.

**Table 2:** Quantitative results of combining U-Net and text encoder in the incompression task. The comparison checkpoints are the epochs for U-Net and text encoder to achieve the same evaluation reward.

C	omparison ch	eckpoints	0	10, 70	20, 10	00 30,	120
	Finetune U	-Net		107.16	143.0	2 17	4.15
I	Finetune text	encoder	84.01	109.48	156.1	0 20	2.79
	Fusion	l		126.59	207.8	0 25	0.18
0,0(origin	nal) 10,70	20,100	30,120	0,0(original)	10,70	20,100	30,120



Fig. 3: Visual examples of combining U-Net and text encoder at different checkpoints in the incompression task.

### 5.2 Results of Aesthetic Reward with Animal Prompts

Aesthetic rewards. Same as previous works, we also conduct experiments using the LAION Aesthetics Predictor [2] as the aesthetic reward function. It should be noted that the aesthetic reward does not consider the coherence with the prompt, and it is easy for the model to hack the reward as shown in [3].

We compare results with DDPO and AlignProb in Fig. 4. Both DDPO and our approach are trained with 10K samples, while AlignProb used early stop to avoid model collapse. Figure 4 presents results for three distinct animals: *jellyfish*, *penguin*, and *swan*. We can observe that both DDPO and AlignProb are over-optimized to the rewards, resulting in over stylized images to maximize



Fig. 4: Comparison with others on unseen animal prompts: *jellyfish*, *penguin* and *swan*. Table 3: Quantitative results of ReFL-LoRA on the ImageReward dataset.

Backbone	Original	$\operatorname{ReFL}$	ReFL-LoRA	TexForce	ReFL+TexForce	ReFL-LoRA + TexForce
SDv1.4	0.2154	0.4485	0.4425	0.4556	0.6553	0.7093
SDv1.5	0.2140	0.5484	0.5558	0.4086	0.6703	0.7438
SDv2.1	0.3891	0.5223	0.5181	0.5084	0.6158	0.6263

aesthetic scores. For instance, DDPO tends to produce images characterized by a blurred yellow background, and AlignProb tends to generate images with oversaturated colors and unrealistic textures. In contrast, our method can generate more lifelike images that closely align with the provided prompts. For example, the composition of the jellyfish image is aesthetically pleasing, and the unrealistic features of the penguin and swan images have been rectified.

### 5.3 Training ReFL with LoRA

To make our model easier to use, we modified the original ReFL to train LoRA weights instead of the entire U-Net, and the results are shown in Tab. 3. We can notice that the performance of ReFL-LoRA is similar to ReFL but much better when combining with TexForce.

# 6 More Qualitative Results

In this section, we select more examples on different backbones in Figs. 5 to 10. Please refer to the demo.html for all results.

 SDv1.5
 ReFL
 TexForce
 ReFL+TexForce

 Portrait of an old sea captain, male, detailed face, fantasy, highly detailed, cinematic,
 Image: Contempt (Contempt (Contemp (Contempt (Contemp (Contematem)(Contematem)(Contemp (Contematem)(Contemp (Contematem)(Cont

art painting by greg rutkowski



A majestic landscape featuring a river, mountains and a forest, A small group of birds is flying in the sky, Harsh winter, very windy, There is a man walking in a deep snow, Cinematic, very, beautiful, painting, in the, style, of Lord of the rings



astronaut drifting afloat in space, in the darkness away from anyone else, alone, black background dotted with stars, realistic



plants, flowers, trees being mixed in a bowl



**Fig. 5:** Results of different methods with Stable Diffusion V1.4 as backbone on ImageReward test dataset.

			1exForce 9
SDv1.5	$\operatorname{ReFL}$	TexForce	$\operatorname{ReFL}+\operatorname{\mathbf{TexForce}}$

small red wooden cottage by the lake, lanterns on the porch, smoke coming out of the chimney, dusk, birch trees, tranquility, two swans swimming on the lake, a wooden rowing boat, cumulus clouds, by charlie bowater, by greg rutkowski



portrait of a cute cyberpunk cat, realistic, professional



field of light blue lotus flowers, minimalistic art, elegant



a coffee mug made of cardboard



Fig. 6: Results of different methods with Stable Diffusion V1.5 as backbone on ImageReward test dataset.

# 7 Limitations

Similar to other RL-based methods, our method also faces the challenges of sample efficiency and the complexity of reward function engineering.

### 10 C. Chen, A. Wang et al. SDv1.5 ReFL **TexForce** ReFL+**TexForce**

landscape photography by marc adamus, mountains with some forests, small lake in the center, fog in the background, sunrays, golden hour, high



an alien planet viewed from space, extremely, beautiful, dynamic, creative, cinematic



photograph of a futuristic cyberpunk flying car at night in the rain



A majestic landscape featuring a river, mountains and a forest, A small group of birds is flying in the sky, Harsh winter, very windy, There is a man walking in a deep snow, Cinematic, very, beautiful, painting, in the, style, of Lord of the rings



Fig. 7: Results of different methods with Stable Diffusion V1.5 as backbone on ImageReward test dataset.

# 8 Broader Impacts

Since TexForce can finetune text-to-image models to satisfy specific rewards, it presents potential societal concerns regarding misinformation, intellectual property rights, and illegal usage of the model.

extremely detailed stunning beautiful futuristic smooth curvilinear museum interior, colorful, hyper, real



beautiful prince, male, golden hair, high, fantasy, art by anato finnstark, joseph leyendecker, peter mohrbacher, ruan jia, marc simonetti, ayami kojima, cedric peyravernay, finnian macmanus, alphonse mucha, victo ngai



small red wooden cottage by the lake, lanterns on the porch, smoke coming out of the chimney, dusk, birch trees, tranquility, two swans swimming on the lake, a wooden rowing boat, cumulus clouds, by charlie bowater, by greg rutkowski



natural suburb with multiple low rise apartment buildings in a park like setting with green hilly lawns and lush trees, pine wooden walls, rustic, large glass windows, cobblestone, grass, white, natural pathways, natural materials, minimalist, swedish design, bright, feng shui, modern, technology, frank lloyd wright



Fig. 8: Results of different methods with Stable Diffusion V1.5 as backbone on ImageReward test dataset.

# References

1. Chen, C., Mo, J., Hou, J., Wu, H., Liao, L., Sun, W., Yan, Q., Lin, W.: Topiq: A topdown approach from semantics to distortions for image quality assessment (2023) 2





 $extremely\ detailed\ stunning\ beautiful\ futuristic\ smooth\ curvilinear\ museum\ interior,$ 



mountains range with waterfall, purple haze, art by greg rutkowski and magali villeneuve



fancy treehouse mansion on top of a mountain overlooking a view of the valley magical realism detailed painting



Fig. 9: Results of different methods with Stable Diffusion V2.1 as backbone on ImageReward test dataset.

- 2. Christoph, S., Romain, B.: Laion-aesthetics. (2022) 6
- 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) 6
- 4. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. In: Advances in Neural Information Processing Systems (2020) 1
- 5. Song, J., Meng, C., Ermon, S.: Denoising diffusion implicit models. arXiv preprint arXiv:2010.02502 (2020) 3

			TexForce	13
SDv1.5	$\operatorname{ReFL}$	TexForce	$\operatorname{ReFL}+\operatorname{\mathbf{TexFor}}$	ce

cinematic movie scene, beautiful Product shot film still of a Syd Mead futuristic modern sleek automobile speeding down a wet street at night in cyperpunk city, motion, hard surface modeling, soft, style of Stanley Kubrick cinematography



field of light blue lotus flowers, minimalistic art, elegant



the grand hall of the sacred library oil painting by james gurney



hyperdetailed samsung store, oak parquet, black walls, digital walls, plants, light



Fig. 10: Results of different methods with Stable Diffusion V2.1 as backbone on ImageReward test dataset.

 Su, S., Lin, H., Hosu, V., Wiedemann, O., Sun, J., Zhu, Y., Liu, H., Zhang, Y., Saupe, D.: Going the extra mile in face image quality assessment: A novel database and model. IEEE Transactions on Multimedia (2023) 2