

# Multi-Adapter joint fine-tuning of Diffusion Models with LoRA for Visual Illusions

Nossaiba Kheiri, Hao Liu, Hansen Lin

Under Supervision of Dr. Ali Hirta, Miao Wang And Gary Kazantsev at Bloomberg



## MOTIVATION

We present a novel approach to fine-tuning large diffusion models for generating illusion-based images using a GPT-based reward model. We build on the illusion generation mechanism in the **Visual Anagrams** paper [?] to train the model to generate images that are difficult for GPT-4 to distinguish from non-illusions. Our approach incorporates **LoRA (Low-Rank Adaptation)** for efficient fine-tuning and **multi-adapter training** using the **PEFT** framework, allowing the model to handle complex visual styles and transformations as training tasks.

We explore the extent to which state-of-the-art LLMs understand complex visual phenomena that challenge human perception. We generate optical illusions and assess the performance of LLMs (specifically GPT-4) in illusion detection.

## PROBLEM: FAILURE MODE OF GENAI FOR IMAGE GENERATION

- 1 Originating from the diffusion model [1]
  - Independent Synthesis
  - Noise shift
  - Correlated noise
- 2 Originating from the illusion pairing model

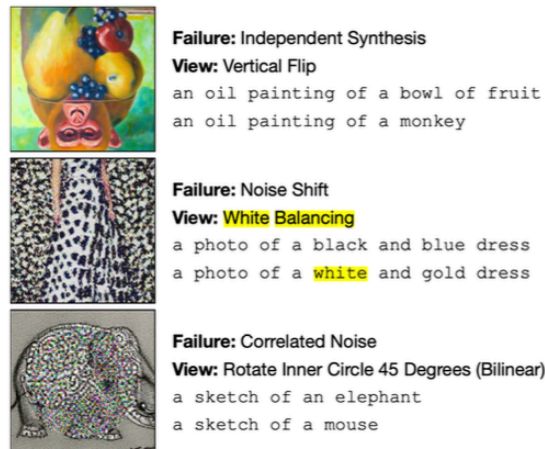


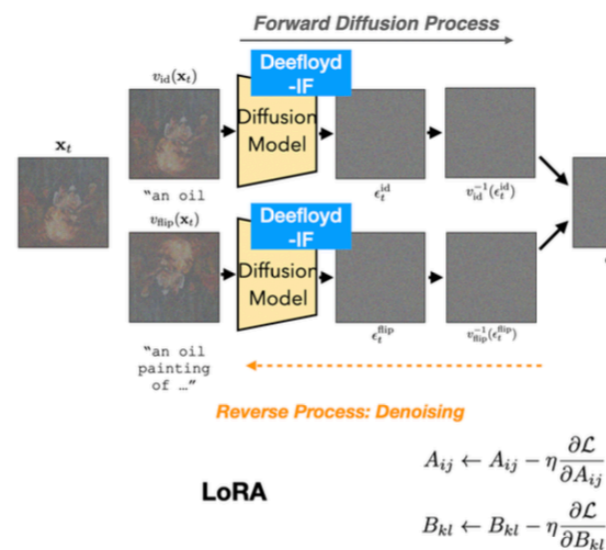
Figure 9. Failures. We highlight three interesting failure cases,

## METHODOLOGY

- 1 Text conditioned Diffusion Process
- 2 Independent Low Rank Adaptation (LoRA) training
- 3 Joint task Fine tuning

Using **classifier free guidance**

$$\epsilon_{CFG} = \epsilon(x, t, \emptyset) + \gamma (\epsilon(x, t, y) - \epsilon(x, t, \emptyset))$$



**Joint Training**

$$W_{final} = \alpha W_{adapter1} + \beta W_{adapter2}$$

**Trainable parameters**

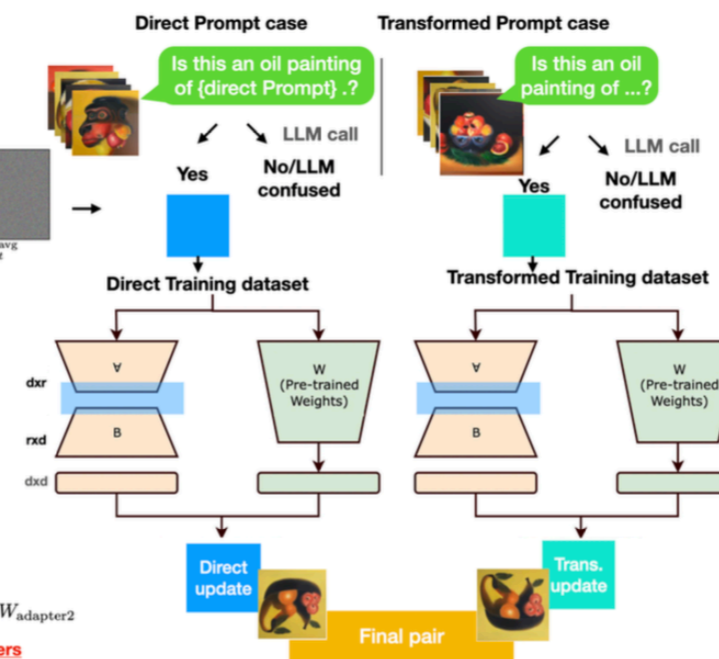
After the diffusion model generates an image  $I$ , GPT evaluates the probability  $P_{GPT}$  that the image contains an illusion. The reward function  $R(I)$  is defined as:

$$R(I) = 1 - P_{GPT}(\text{illusion}|I)$$

Where:  $P_{GPT}(\text{illusion}|I)$  is the probability that the LLM detects the illusion in the image  $I$ .

## LoRA AND GPT AS THE REWARD MODEL (RL AGENT)

Low-Rank Adaptation (LoRA) **reduces the number of trainable parameters** during fine-tuning by introducing low-rank matrices updates to the weight matrices of a pre-trained model to capture important **task-specific information**.



## POSSIBLE FUTURE WORK: ONLINE REINFORCEMENT LEARNING

The proposed approach can be extended to an **online reinforcement learning (RL)** framework. In this setup, the model generates images in real-time and receives feedback from GPT, updating its parameters dynamically. In this scenario, the model continuously improves by learning from its performance during inference.

## CONCLUSION

The model learns to generate illusionary images that are increasingly difficult for GPT to detect. Multi-adapter training allows flexible blending of original and transformed prompts, enhancing the model's ability to produce complex visual illusions. This approach is computationally scalable and adaptable to real-time online reinforcement learning, making it suitable for advanced image generation tasks.

## OBJECTIVES/INNOVATION/USAGE

- 1 Adjust the finetuning by replicating the human behavior on illusions (as in [2]) to the LLM's (GPT model) behavior.
- 2 Illusion generation using diffusion model and parallel denoising
- 3 Fine-tuning illusions to be harder to detect

## REFERENCES

- [1] "DreamBooth Fine-Tuning with LoRA," peft/main/en/task\_guides/dreambooth\_lora.
- [2] Ying Fan et al., "DPOK: Reinforcement Learning for Fine-Tuning Text-to-Image Diffusion Models", <https://doi.org/10.48550/arXiv.2305.16381>.
- [3] Daniel Geng, Inbum Park, and Andrew Owens, "Visual Anagrams: Generating Multi-View Optical Illusions with Diffusion Models", <http://arxiv.org/abs/2311.17919>.



Weights & Biases (Wandb) Training 1



Github