

Directly Fine-Tuning Diffusion Models on Differentiable Rewards

Kevin Clark* Paul Vicol* Kevin Swersky David Fleet

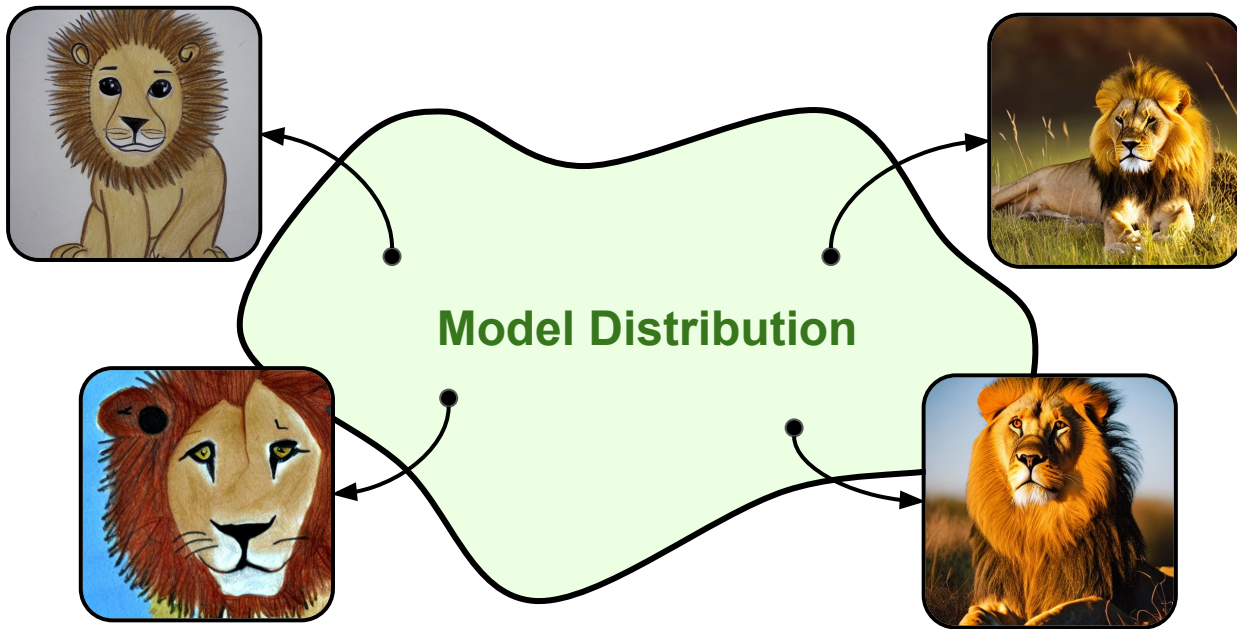
ICLR 2024



Google DeepMind

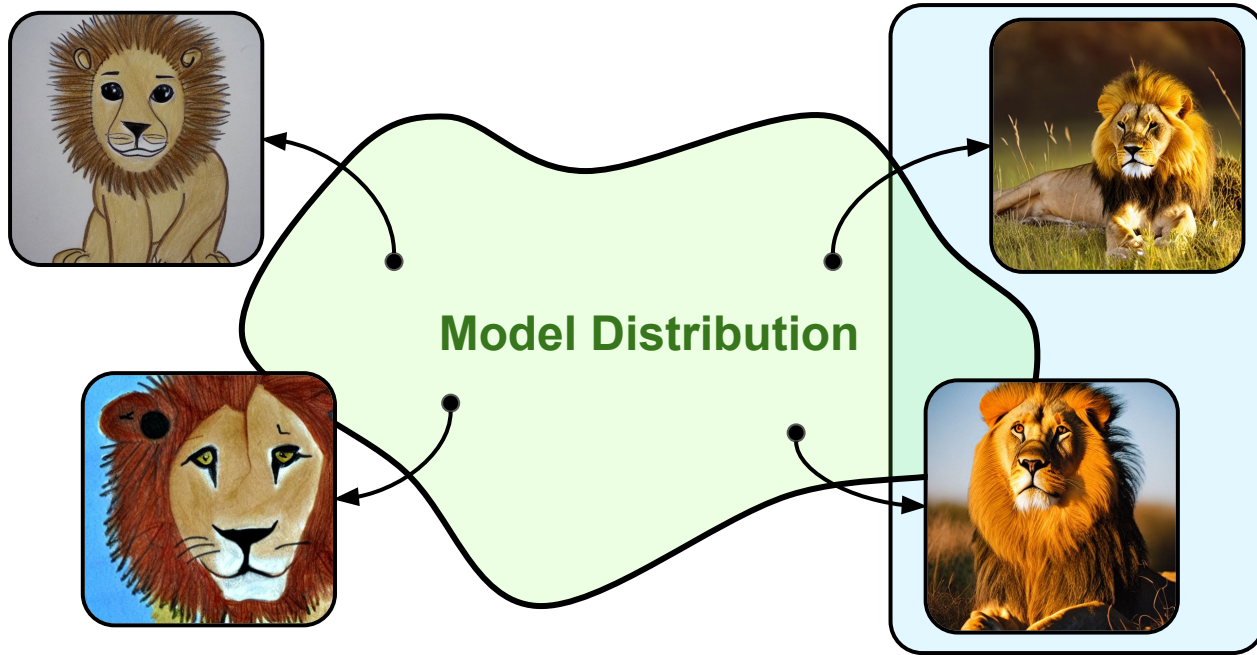
Motivation: Fine-Tuning Diffusion Models

- *Diffusion models are pre-trained to model the data distribution (e.g., diverse web images)*



Motivation: Fine-Tuning Diffusion Models

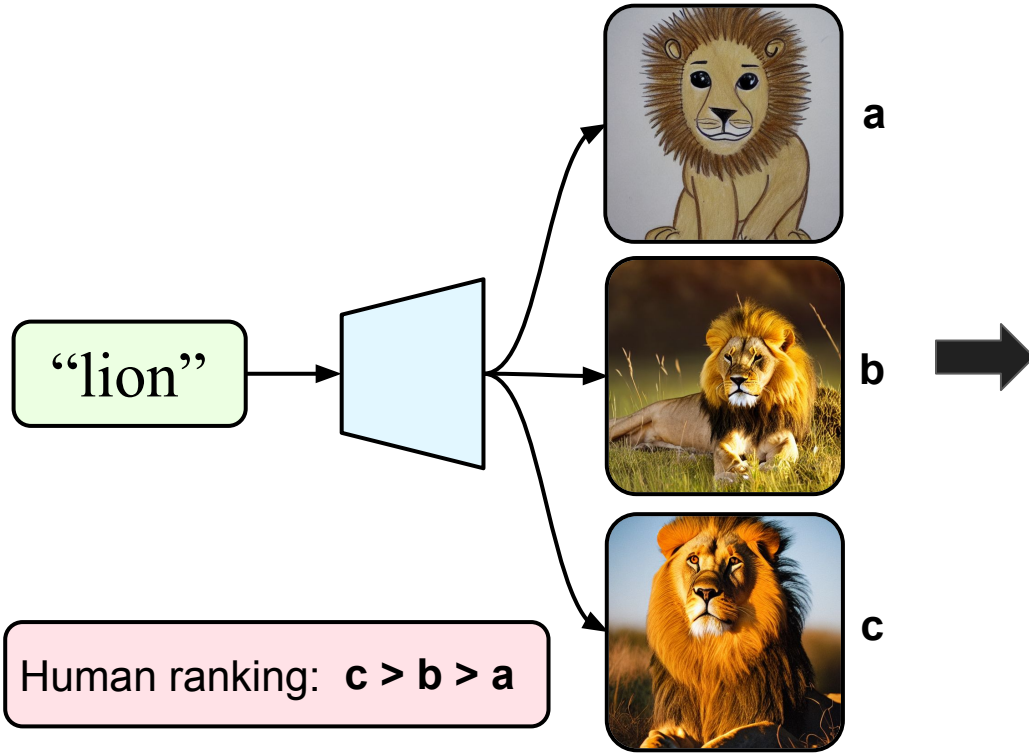
- *Diffusion models are pre-trained to model the data distribution (e.g., diverse web images)*



- But modeling the data distribution exactly often does not align with desired behavior
- E.g., *generating images with certain aesthetic qualities*

Human Preference Datasets

Human Preference Data



Reward Models

$$r_{\phi} \left(\text{Image a}, \text{"lion"} \right) = 0.1$$

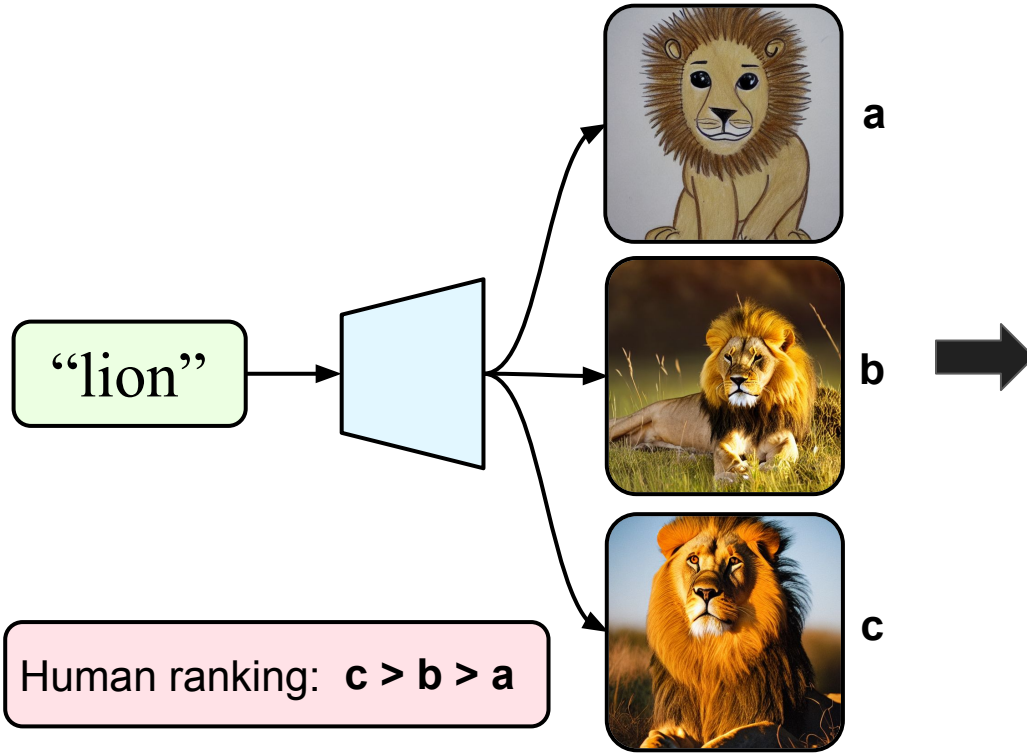
$$r_{\phi} \left(\text{Image c}, \text{"lion"} \right) = 0.6$$

Human Preference Dataset (HPDv2) (Wu et al., 2023)

Pick-a-Pic Dataset (Kirstain et al., 2023)

Human Preference Datasets

Human Preference Data



Reward Models

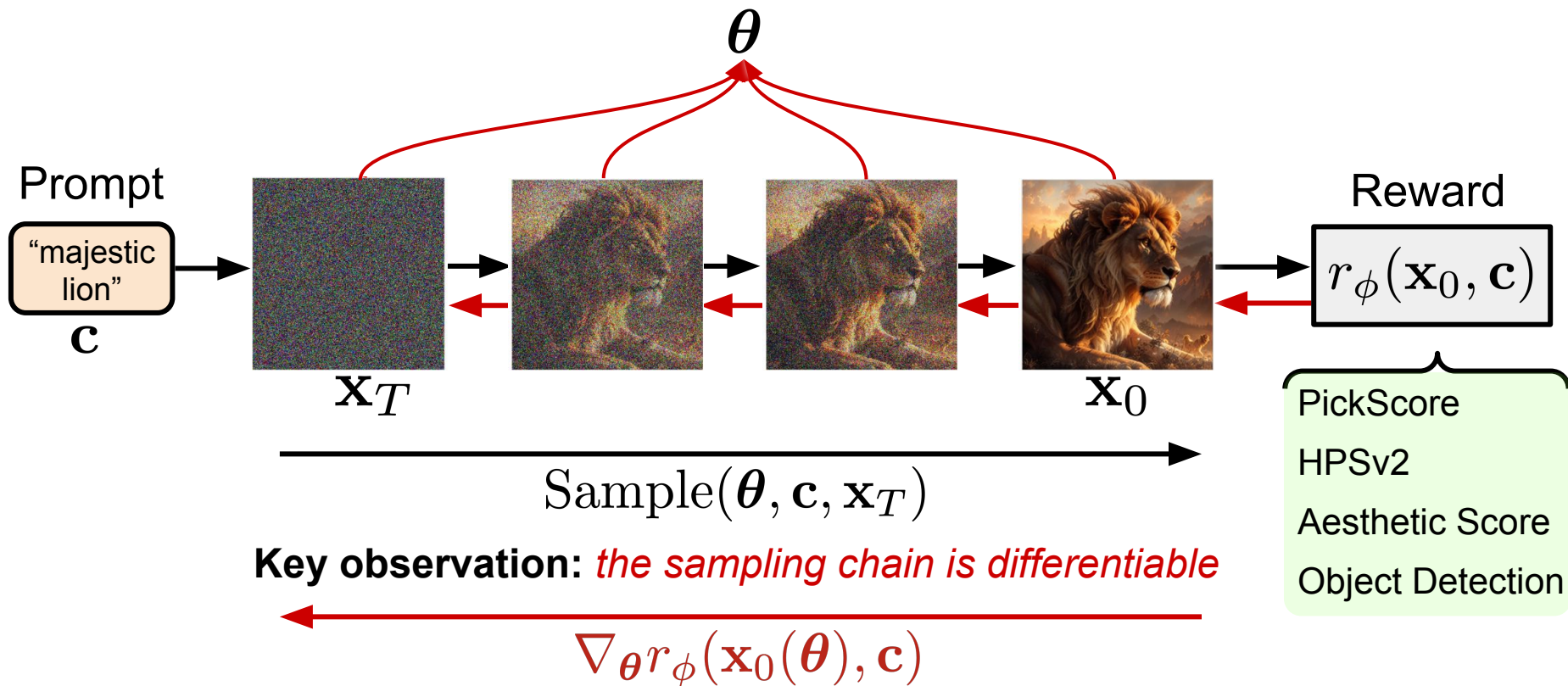
$$r_{\phi} \left(\text{Image a}, \text{"lion"} \right) = 0.1$$

$$r_{\phi} \left(\text{Image b}, \text{"lion"} \right) = 0.6$$

Prior work on diffusion fine-tuning used *RL-based techniques* (Black et al., 2023, Fan et al., 2023) → Promising results, but *sample-inefficient*

Direct Reward Fine-Tuning (DRaFT)

Objective: $J(\theta) = \mathbb{E}_{\mathbf{c} \sim p_{\mathbf{c}}, \mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [r_{\phi}(\text{Sample}(\theta, \mathbf{c}, \mathbf{x}_T), \mathbf{c})]$

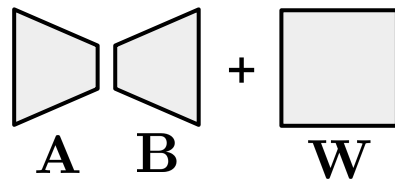


Memory Challenges

- **Main challenge:** *storing intermediate activations* through the unrolled sampling chain for use in backprop is expensive
- DRaFT uses two simple techniques to *keep memory usage tractable*:

① DRaFT fine-tunes *LoRA parameters* (Hu et al., 2022)

- *Reduces memory usage, yields smaller checkpoints*
- **Extra benefit:** Can interpolate between the original and fine-tuned model by re-scaling the LoRA parameters

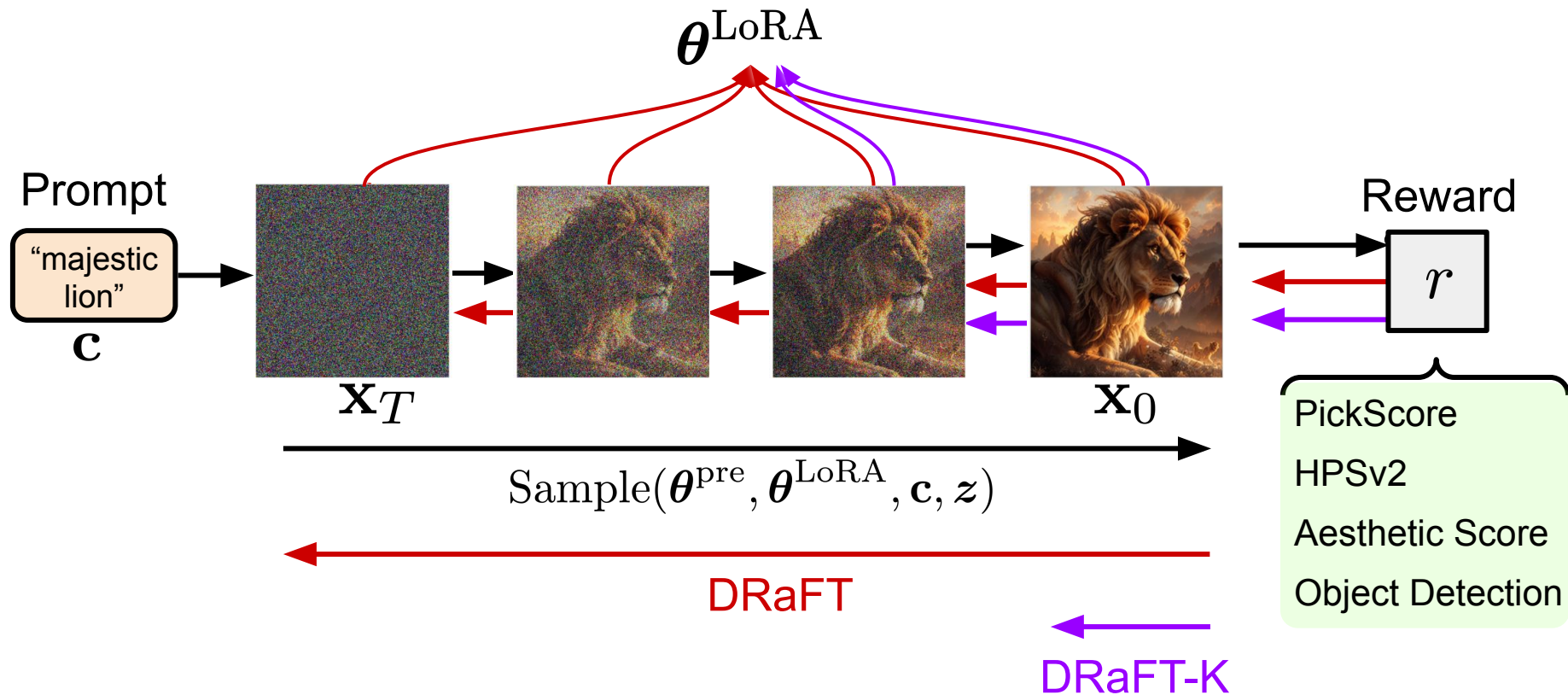


② *Gradient Checkpointing*

- Stores a subset of the intermediate activations in memory, and *recomputes non-stored ones on-the-fly during backprop*
- Only need to add one `@jax.checkpoint` to be able to compute the gradient

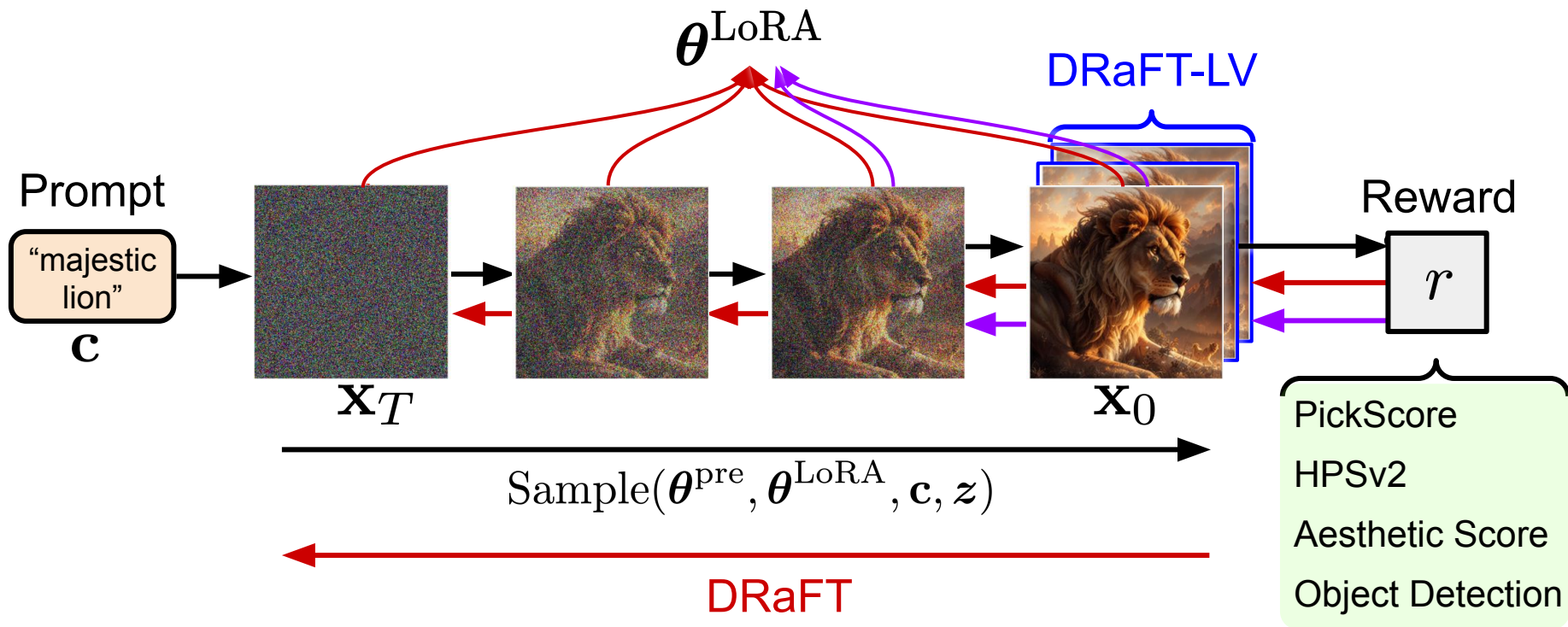
DRaFT + Efficiency Improvements

- **DRaFT-K:** Truncates backpropagation through *only the last K sampling steps*



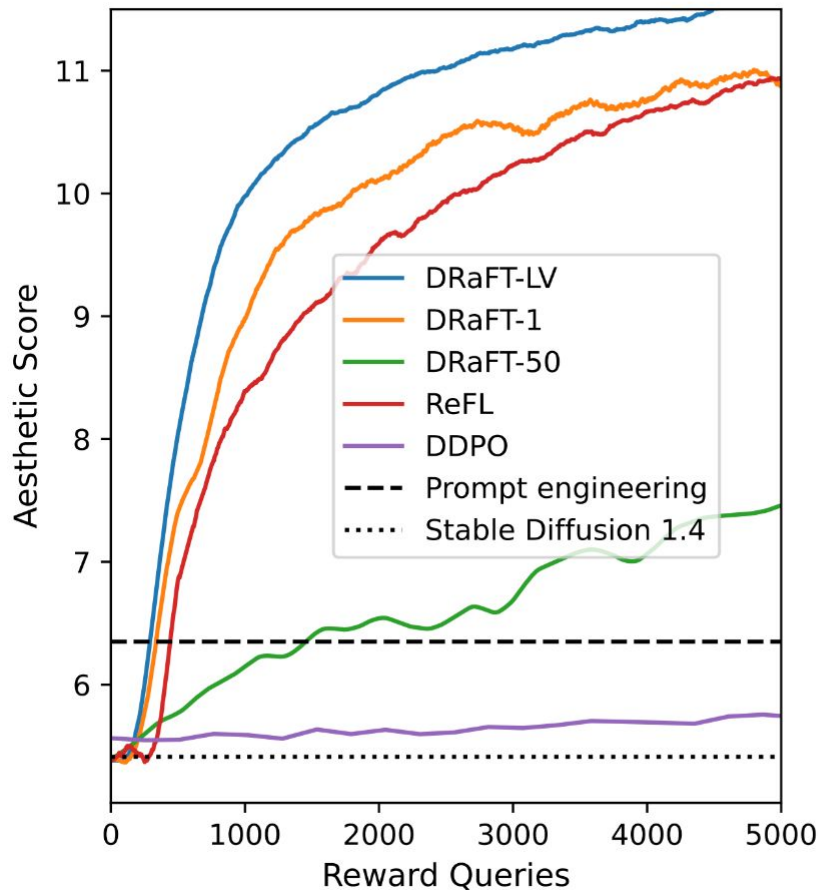
DRaFT + Efficiency Improvements

- **DRaFT-LV**: lower-variance gradient estimator. *Noise the generated image n times, and use the average reward gradient over these examples.*
- Using $n=2$ is around *$2\times$ more efficient than DRaFT-1* while adding around 10% overhead



Quantitative Reward Optimization Comparison

- **Goal:** *optimize aesthetic quality scores of the LAION aesthetic predictor.*
- We compare against DDPO (Black et al., 2023), ReFL (Xu et al., 2023), and a prompt engineering baseline.
- *DRaFT is much more sample-efficient than RL,* as it leverages gradient information
- Because of its low-variance gradient estimate, *DRaFT-LV further improves training efficiency.*



DRaFT: Fine-Tuning for Human Preferences

- DRaFT using human preference reward models yields *more detailed and stylized images than baseline Stable Diffusion*

A stunning beautiful oil painting of a lion, cinematic lighting, golden hour light.



Highly detailed photograph of a meal with many dishes.



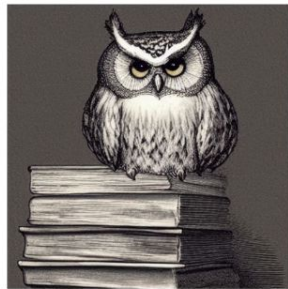
A racoon washing dishes.



Ultra realistic photo of a single light bulb, dramatic lighting.



A fluffy owl sits atop a stack of antique books in a detailed and moody illustration.



Impressionist painting of a cat, textured, hypermodern.



Stable Diffusion 1.4



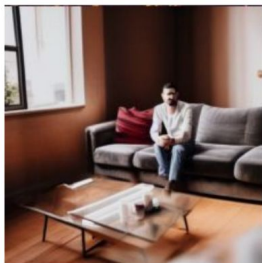
After Reward
Fine-Tuning

Scaling LoRA Parameters

Can *interpolate* between the original pre-trained model and the LoRA-adapted model by *scaling the LoRA weights*

Pre-trained model $\theta^{\text{pre}} + \alpha\theta^{\text{LoRA}}$ *Fine-tuned model*

$\alpha = 0.0$



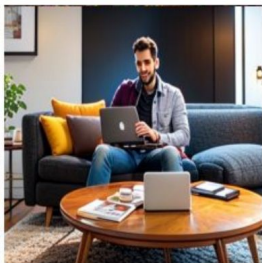
$\alpha = 0.6$



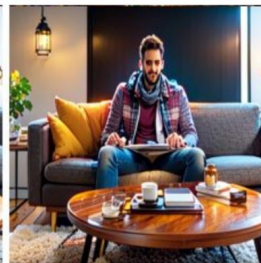
$\alpha = 0.8$



$\alpha = 0.9$



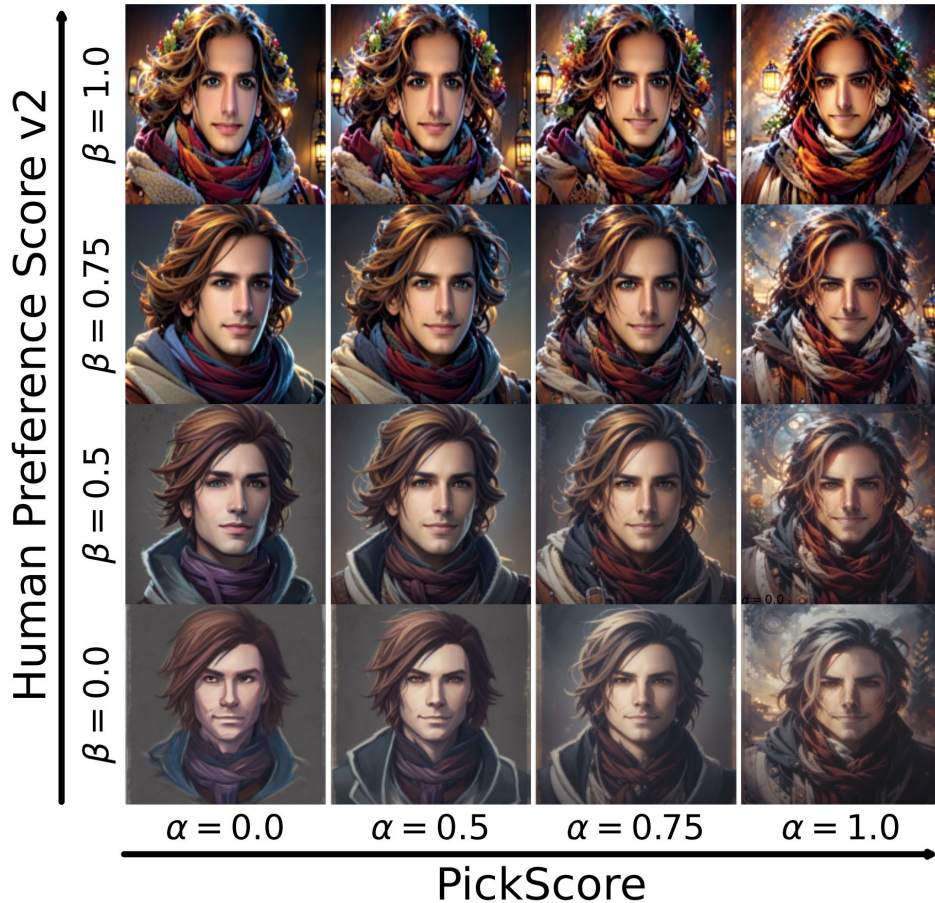
$\alpha = 1.0$



Mixing LoRA Parameters

- LoRA parameters *fine-tuned independently for different rewards* can be combined post-hoc without additional training
- Taking *linear combinations of LoRA parameters*:

$$\theta^{\text{pre}} + \alpha \theta_{\text{PickScore}}^{\text{LoRA}} + \beta \theta_{\text{HPSv2}}^{\text{LoRA}}$$



Object Detection for Addition and Removal

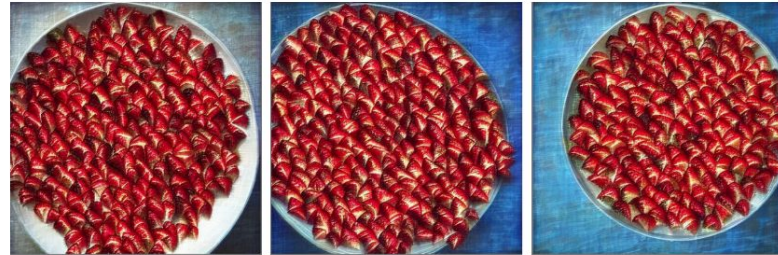
Object Removal: Farmer's Market - People

Object Addition: Fruit Bowl + Strawberries

Before
Fine-Tuning



After
Fine-Tuning



- Well-crafted reward functions can allow us to specify what kinds of images we wish to generate *even when no real images exist that satisfy the criteria*

Diffusion Adversarial Examples

Fine-Tuning Progress



- Fine-tuning a diffusion model such that *images generated based on prompts {"bear", "dog", "mouse"} are classified as target class "cat"* by a ResNet-50 classifier.
- This classifier is *texture-biased*, as the fine-tuned images have cat-like textures while keeping the animal shapes mostly unchanged.

Conclusion

- DRaFT is an efficient framework for *fine-tuning diffusion models on differentiable rewards by leveraging reward gradients*.
 - DRaFT is substantially *more efficient than RL-based fine-tuning approaches*
- We applied DRaFT to a diverse array of reward functions
 - Human preference rewards, object detection, classification, and more
- Just as RLHF has become crucial for deploying LLMs, *reward fine-tuning may become a key step for improving image generation models*.

Thank you!