



Training Diffusion Models with Reinforcement Learning

Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, Sergey Levine
University of California, Berkeley Massachusetts Institute of Technology

Motivation

- Diffusion models are trained with an approximation to the log-likelihood objective. However, most use cases of diffusion models are concerned instead with downstream objectives such as human-perceived image quality.
- This paper investigates reinforcement learning methods for directly optimizing diffusion models for such objectives.
- It describes how posing denoising as a multi-step decision making problem enables a class of policy gradient algorithms, which we refer to as denoising diffusion policy optimization (DDPO).

Compressibility: *llama*



Incompressibility: *bird*



Aesthetic Quality: *rabbit*



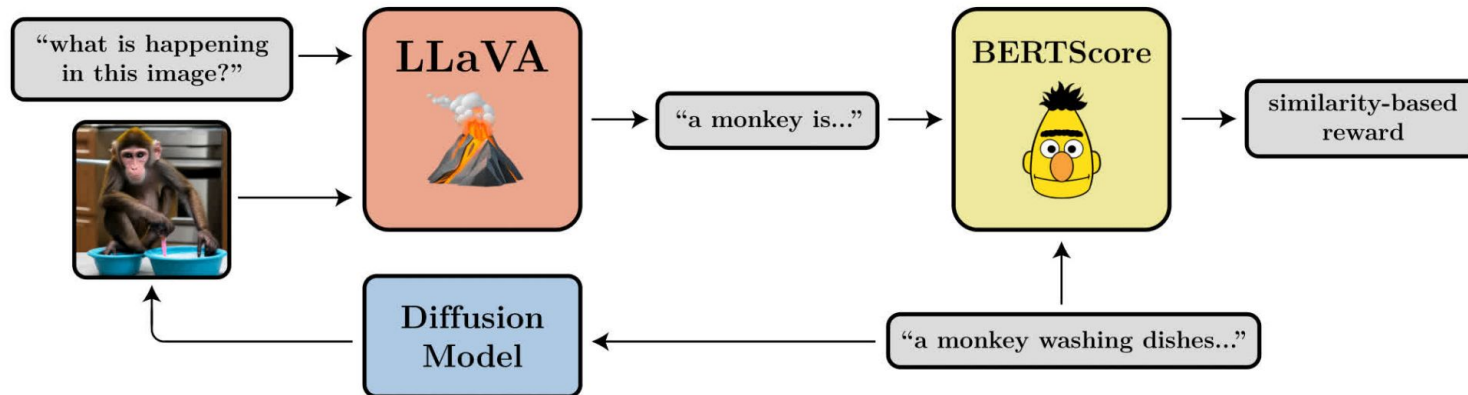
Prompt-Image Alignment: *a raccoon washing dishes*



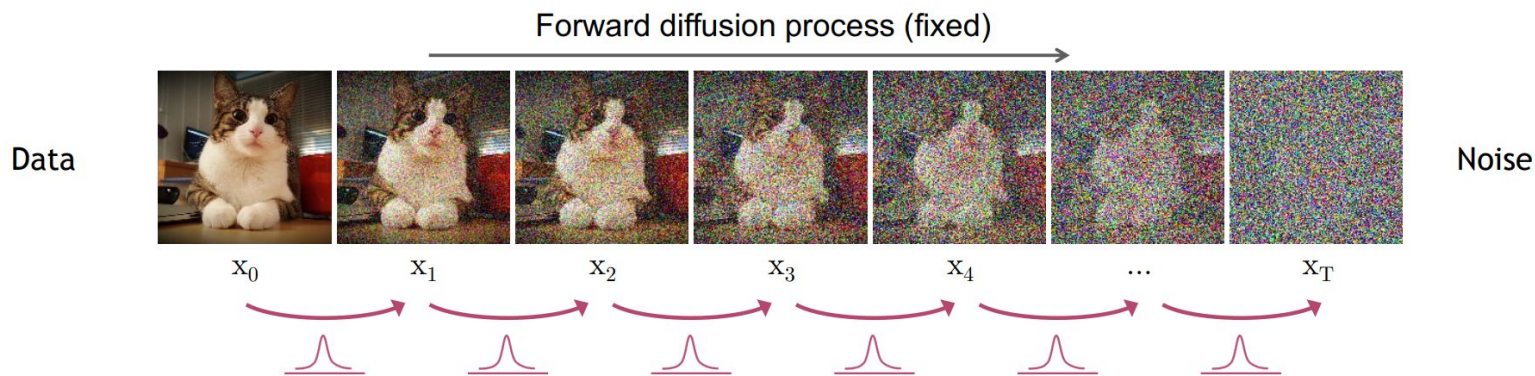
RL training

Approach Summary

- DDPO is used to finetune Stable Diffusion on objectives that are difficult to express via prompting, such as image compressibility, and those derived from human feedback, such as aesthetic quality.
- DDPO can also be used to improve prompt-image alignment without any human annotations using feedback from a vision-language model.

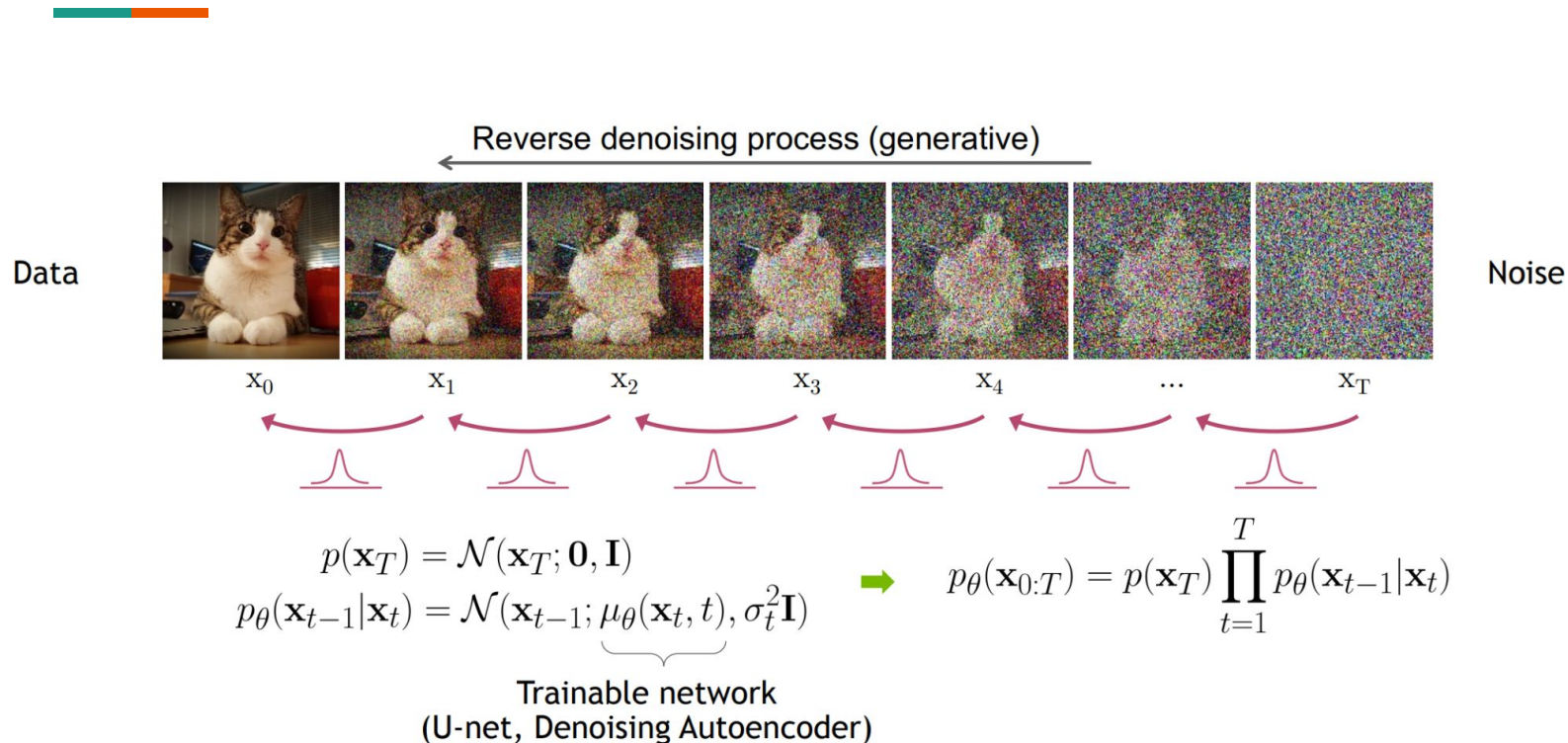


Background - Diffusion Models



$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \quad \rightarrow \quad q(\mathbf{x}_{1:T} | \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1}) \quad (\text{joint})$$

Background - Diffusion Models

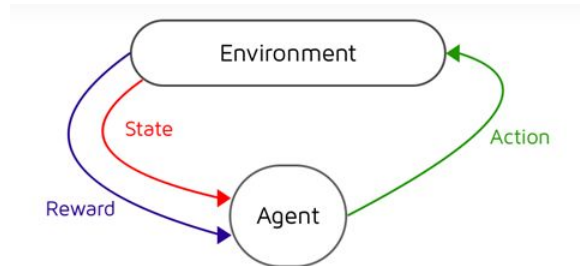


Background - Reinforcement Learning

- **Markov Decision Process:** $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \gamma, \mathcal{R})$,

where \mathcal{S} is the state space, \mathcal{A} is the action space, \mathcal{P} is the transition kernel, γ is the discount factor, and \mathcal{R} is the reward function.

- **Policy** (π): A map from state space to action space. May be stochastic.
- **Reward function** - $\mathcal{R}(s)$: Maps each state-action pair to a scalar called reward.
- **Value function** - $V(s,a)$: Total expected return starting from state s and taking a and following the policy π , discounted by γ .



Problem Statement



- The diffusion model induces a sample distribution $p_{\theta}(\mathbf{x}_0 | \mathbf{c})$.
- The denoising diffusion RL objective is to maximize a reward signal r defined on the samples and contexts.

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

REWARD-WEIGHTED REGRESSION



- Within the RL formalism, the RWR procedure corresponds to the following one-step MDP

$$\mathbf{s} \triangleq \mathbf{c} \quad \mathbf{a} \triangleq \mathbf{x}_0 \quad \pi(\mathbf{a} \mid \mathbf{s}) \triangleq p_\theta(\mathbf{x}_0 \mid \mathbf{c}) \quad \rho_0(\mathbf{s}) \triangleq p(\mathbf{c}) \quad R(\mathbf{s}, \mathbf{a}) \triangleq r(\mathbf{x}_0, \mathbf{c})$$

with a transition kernel that immediately leads to an absorbing termination state. Therefore, maximizing $J_{\text{DDRL}}(\theta)$ is equivalent to maximizing the RL objective $J_{\text{RL}}(\pi)$ in this MDP.

- The reward is then used as a weighting term in the log-likelihood objective function.

$$\mathcal{L}(\theta) = \mathbb{E}_{(\mathbf{x}, \mathbf{z}) \sim \mathcal{D}^{\text{model}}} \left[-r_\phi(\mathbf{x}, \mathbf{z}) \log p_\theta(\mathbf{x} \mid \mathbf{z}) \right]$$

REWARD-WEIGHTED REGRESSION



$$\mathcal{L}(\theta) = \mathbb{E}_{(\mathbf{x}, \mathbf{z}) \sim \mathcal{D}^{\text{model}}} \left[-r_{\phi}(\mathbf{x}, \mathbf{z}) \log p_{\theta}(\mathbf{x}|\mathbf{z}) \right]$$

- Limitations:
 - Diffusion loss does not involve an exact log-likelihood – it is instead derived as a variational bound.
 - It ignores the sequential nature of the denoising process, only using the final sample.
 - $p_{\theta}(\mathbf{x}_0 | \mathbf{c})$ is an arbitrarily complicated distribution.

DENOISING DIFFUSION POLICY OPTIMIZATION

- MDP Formulation:

$$\begin{aligned} \mathbf{s}_t &\triangleq (\mathbf{c}, t, \mathbf{x}_t) & \pi(\mathbf{a}_t | \mathbf{s}_t) &\triangleq p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{c}) & P(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) &\triangleq (\delta_{\mathbf{c}}, \delta_{t-1}, \delta_{\mathbf{x}_{t-1}}) \\ \mathbf{a}_t &\triangleq \mathbf{x}_{t-1} & \rho_0(\mathbf{s}_0) &\triangleq (p(\mathbf{c}), \delta_T, \mathcal{N}(\mathbf{0}, \mathbf{I})) & R(\mathbf{s}_t, \mathbf{a}_t) &\triangleq \begin{cases} r(\mathbf{x}_0, \mathbf{c}) & \text{if } t = 0 \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

δ_y = Dirac delta distribution with nonzero density only at y .

- This provides a general framework (DDPO) with an arbitrary reward function, that can be used to optimize arbitrary downstream objectives.

DENOISING DIFFUSION POLICY OPTIMIZATION



- DDPO_{SF} - Score Function Policy Gradient Estimator

$$\nabla_{\theta} \mathcal{J}_{\text{DDRL}} = \mathbb{E} \left[\sum_{t=0}^T \nabla_{\theta} \log p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{c}) r(\mathbf{x}_0, \mathbf{c}) \right] \quad (\text{DDPO}_{\text{SF}})$$

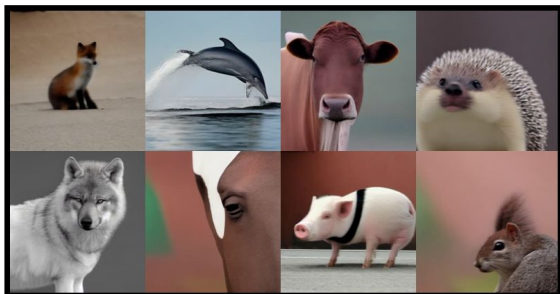
- DDPO alternates collecting denoising trajectories $\{\mathbf{x}_T, \mathbf{x}_{T-1}, \dots, \mathbf{x}_0\}$ via sampling and updating parameters via gradient descent.

Results

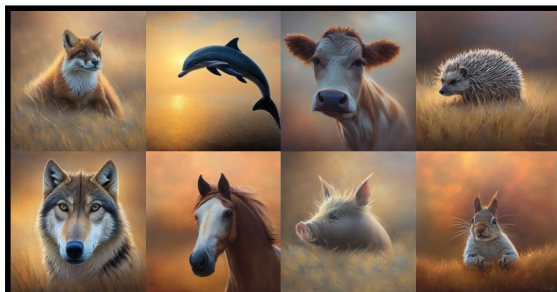
Pretrained



Compressibility



Aesthetic Quality



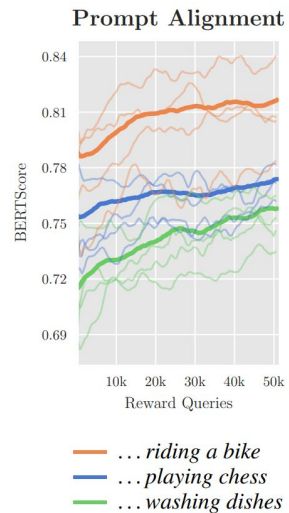
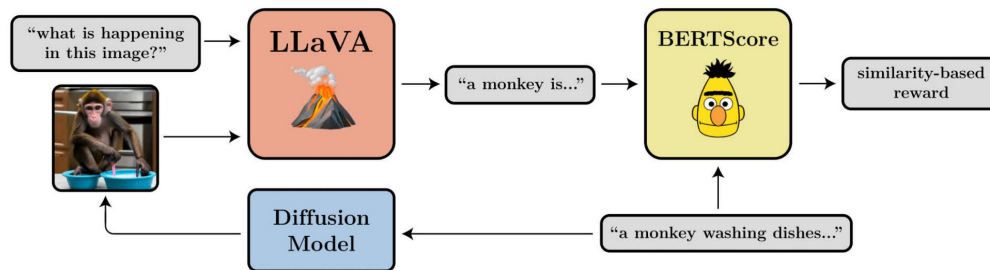
Incompressibility



- **Aesthetic Quality:**
Reward provided by LAION aesthetics predictor.
- **Compressibility/
Incompressibility:**
Reward provided by JPEG compression

Results

Automated Prompt Alignment



Generalization



Finetuning on a limited set of animals generalizes to both new animals and non-animal everyday objects. The prompts for the rightmost two columns are “a capybara washing dishes” and “a duck taking an exam”