

#### Harvard John A. Paulson School of Engineering and Applied Sciences

## **The Power of Resets!** Learning better, one reset at a time.

Kianté Brantley, Assistant Professor, Harvard University

The Power of Resets in Online Reinforcement Learning by Mhammedi et al. 2024



"A large language model is a deep learning model trained on massive text datasets to understand and generate human language."

# Large Language Model

How do *different* objectives shape LLM behavior?





Long Ouyang et al. Training language models to follow instructions with human feedback OpenAI 2022

## Learning to reason with LLMs

"We are introducing OpenAl o1, a new large language model trained with reinforcement learning to perform complex reasoning." - Openai



### MinxEnt RL

Maximize reward



### MinxEnt RL

## $J(\pi) = \mathbb{E}_{\pi} \left[ r(x, y) \right]$

- Sparse or Delayed Rewards
- Exploration vs. Exploitation
- Credit Assignment
- Sample Inefficiency

When you say, "This is a reinforcement learning problem," you should say it with the same excitement as "This is NP-hard." - Tim Vieira

$$\left[ -\frac{1}{\eta} \mathbb{D}_{\mathrm{KL}}(\pi \parallel \pi_{\mathrm{ref}}) \right]$$









Go-Explore: a New Approach for Hard-Exploration Problems by Ecoffet et. al 2019

Learning Montezuma's Revenge from a Single Demonstration by Salimans et. al 2018





# How can we develop efficient algorithms for solving the MinxEnt RL objective, assuming access to environment resets?

### Outline

- Resetting with reference policy
- Resetting with demonstration data
- Resetting with the current policy

#### blicy on data bolicy

### Outline

### Resetting with reference policy (to boost exploration)

Reset allows us to rollout a policy from partial sentences

- **1.** Sample a prompt from  $x \sim D$
- **2.** Sample a response from  $y \sim \pi$



## Inject additional data sources into experience collection

Reset allows us to rollout a policy from partial sentences

- **1. Sample a prompt from**  $x \sim D$
- **2.** Sample a response from  $y \sim \pi$
- 3. Reset and sample a continuation of the response from  $y \sim \pi$





#### Inject additional data sources into experience collection



Reset allows us to rollout a policy from partial sentences

- **1. Sample a prompt from**  $x \sim D$
- **2.** Sample a response from  $y \sim \pi$
- 3. Reset and sample a continuation of the response



**Transition:** P(s'|s,a)Deterministic

#### Inject additional data sources into experience collection





### Learning to Generate Better Than Your Teacher







**Rajkumar Ramamurthy** 





Dipendra Misra

Jonathan D. Chang

### **Rollin and Rollout** approaches



- •

Proximal Policy Optimization, John Schulman et al., 2017

$$\mathbb{E}_{\pi} \left[ \hat{r}(x, y) \right] + \frac{1}{\eta} D_{\text{KL}}(\pi | | \pi_{\text{ref}})$$
Constrains policy to stay near  $\pi_{\text{ref}}$ 

Does not leverage problem-specific structure • Samples prompts  $x \sim D$ • Scores action with  $\hat{r}(x, y)$ 

#### **Rollin and Rollout** approaches



Proximal Policy Optimization, John Schulman et al., 2017

$$\mathbb{E}_{\pi} \left[ \hat{r}(x, y) \right] + \frac{1}{\eta} D_{\text{KL}}(\pi | | \pi_{\text{ref}})$$
Constrains policy to stay near  $\pi_{\text{ref}}$ 

Does not leverage problem-specific structure Samples prompts  $x \sim D$ Scores action with  $\hat{r}(x, y)$ 

• Sample prompts from a mixture  $x \sim \beta D + (1 - \beta) d^{\pi_{ref}}$ • Scores actions with  $\hat{r}(x, y)$ • Intuition: Richer initial states boost exploration



## **Theory of PPO++**

#### Let $\pi^{\star}$ be a high quality policy covered by $\pi_{ref}$

Performance gap  $J(\pi^{\star}) - J(\pi^{t}) \le O\left(\begin{array}{c}H^2 \max \\ s\end{array}\right)$ 

Assume bound density ratio and  $\pi_{\rm ref}$  provides coverage for  $\pi^{\star}$ 

Approximately Optimal Approximate Reinforcement Learning Kakade and Langford 2002

$$\mathbb{E}_{s \sim \beta \rho^{\pi} \text{ref} + (1 - \beta)D}$$

$$\frac{d^{\pi^{\star}}(s)}{d^{\pi^{\text{ref}}}(s)} \in \mathcal{E}$$

$$\left[\max_{a} A^{\pi^{t}}(s,a)\right] \leq \epsilon$$

Assume that one-step local improvement over  $\pi^t$  is small



## **Experimental Setup**

#### Task Statement

Given a reddit post, write a TL;DR (short summary).

#### **Example Post**

#### **SUBREDDIT:** r/dogs

**TITLE:** [HELP] Not sure how to deal with new people/dogs and my big ole pup **POST:** I have a three year old Dober/Pit mix named Romulus ("Rome" for short). I live with 3 other dogs: a 10 year old labrador, a 2 year old French Bulldog and a 8 year old maltese mix. The four of them get along just fine, Rome and the Frenchie are best best best best friends. He isn't the best at meeting new people, but not ALWAYS....Then, the crux of the matter: I want to have a 4th of July party. Several people want to bring their dogs. I doubt I can say "no dogs allowed" and I don't want to let everyone else bring their dog and make mine stay at day care all day.

Learning to summarization from human feedback Stiennon et al. 2017

#### **Example Human Label**

**TL;DR:** HOW do I introduce new people? HOW do l introduce new dogs? WHAT do l do about 4th of July??



## **Experimental Setup**

#### **Task Statement**

Given a reddit post, write a TL;DR (short summary).

#### **Dataset Composition**

- 210K Prompts total

# 117K Prompts with Human Labels Used to do RL fine-tuning 93K Prompts with Human Preference Labels Used to pre-train a reward model

## **Experimental Results: TL;DR**

#### **GPT4 Winrate Prompt Template**

Which of the following summaries does a better job of summarizing the most important points in the given forum Post? FIRST provide a one-sentence comparison of the two summaries, explaining which you prefer and why. SECOND, on a new line, state only "A" or "B" to indicate your choice.

Post: <Post> A: <TLDR A> B: <TLDR B>



**GPT-4 Win Rate** 

### Summary

- PPO++ does not leverage a problem-specific structure.
- The ability to reset is a unique property of MDPs for LLMs.
- PPO is a simple algorithm that leverages resets.

#### Can we improve performance by mixing in a different distribution?



#### Can we improve performance by mixing in a different distribution?





#### **GPT-4 Win Rate**

### Outline

- Reset with reference policy
- Reset with the demonstration data

### Outline

### **Reset with the demonstration data (to boost exploration)**



#### Can we improve performance by mixing in a different distribution?

Most text generation tasks have offline label response

- **1.** Sample a prompt from and response from  $(x, y) \sim D$
- **2.** Reset using the response from D



Most text generation tasks have offline label response

- 1. Sample a prompt from and response from
- **2.** Reset using the response from D
- 3. Sample a continuation of the response from  $y \sim \pi$



$$\mathbf{n}(x,y) \sim D$$

### **Dataset Reset Policy Optimization**



Jonathan D. Chang



Wenhao Zhan



Owen Oertell





Dipendra Misra



Jason Lee



Wen Sun

### **Rollin and Rollout** approaches



Does not leverage problem-specific structure
Samples prompts x ~ D
Scores action with r(x, y)

Sample prompts from a mixture x ~ βD<sub>x</sub> + (1 - β)D<sub>x||y</sub>
Scores actions with r̂(x, y)
Intuition: Richer initial states boost exploration

## **Informal Theory of DR-PO**

#### **Informal statement:**

When using NPG as the policy optimization oracle, DR-PO learns a policy that is at least as good as any policy covered by the offline data D



**Trajectory-wise density** 

#### **Coverage assumptions:**

$$\frac{d^{\pi^*}(x,y)}{d^{\pi_{\text{ref}}}(x,y)} \le C_2 < \infty$$

**State-action sample-wise density** 

## **Experimental Setup**

#### Task Statement

Given a reddit post, write a TL;DR (short summary).

#### **Example Post**

#### **SUBREDDIT:** r/dogs

**TITLE:** [HELP] Not sure how to deal with new people/dogs and my big ole pup **POST:** I have a three year old Dober/Pit mix named Romulus ("Rome" for short). I live with 3 other dogs: a 10 year old labrador, a 2 year old French Bulldog and a 8 year old maltese mix. The four of them get along just fine, Rome and the Frenchie are best best best best friends. He isn't the best at meeting new people, but not ALWAYS....Then, the crux of the matter: I want to have a 4th of July party. Several people want to bring their dogs. I doubt I can say "no dogs allowed" and I don't want to let everyone else bring their dog and make mine stay at day care all day.

Learning to summarization from human feedback Stiennon et al. 2017

#### **Example Human Label**

**TL;DR:** HOW do I introduce new people? HOW do l introduce new dogs? WHAT do l do about 4th of July??



### **Experimental Results: TL;DR**

#### 80% **Takeaways**: 60% 1. Algorithms that use resets perform better than those that do not. 40% 2. DR-PO consistently outperforms all baseline algorithms. 20%

0%





#### **GPT-4** Win Rate

## **Experimental Results: TL;DR**

#### TL;DR Summarization



#### **Takeaways:**

DR-PO consistently achieves better RM scores with lower KL than PPO.



## **Experimental Setup**

#### **Task Statement**

Anthropic's Helpful Harmful task where our model tries to produce an engaging and helpful response to dialogue sequences.

#### **Example Dialogue**

Human: What do I do if I crack a molar?

**Chosen Assistant:** If you cracked a molar, I imagine you're quite concerned, but there's no need to panic, you just need to schedule an appointment with your dentist.

Rejected Assistant: I'm sorry to hear that.

### **Experimental Results**





#### **Takeaways:**

- 1. Online methods outperform offline baselines.
- 2. DR-PO outperforms all baselines at every model scale.

### Summary

- Resetting from offline demonstration data enhances performance.
- DR-PO is provably efficient and improves upon PPO++ in theory.
- DR-PO is as simple as PPO and requires no additional computation.

### Summary

- DR-PO is provably efficient and improves upon PPO++ in theory.
- DR-PO is as simple as PPO and requires no additional computation. •

#### Resetting DR-PO from offline demonstration data enhances performance.



### Outline

- Resetting with reference policy
- Resetting with demonstration data
- Resetting with the current policy

#### blicy on data bolicy

### Outline

### Resetting with the current policy (to reduce computation)



#### How can resets be used to make RL algorithms more efficient in computation and memory?

### MinXent RL

$$\pi_{i+1} = \arg \max_{\substack{\pi \in \Pi \\ y \sim \pi}} \mathbb{E} \left[ r(x, y) \right] - \frac{1}{\eta} D_{\text{KL}}(\pi \parallel \pi_t)$$

$$\forall x, y : \pi_{t+1}(y \mid x) = \frac{\pi_t \exp(\eta r(x, y))}{Z(x)}$$

$$\forall x, y : r(x, y) = \frac{1}{\eta} \left( \ln Z(x) + \ln \left( \frac{\pi_{t+1}(y \mid x)}{\pi_t(y \mid x)} \right) \right)$$

$$L(r,D) = -\mathbb{E}_{(x,y_w,y_l)\sim D}\left[\log\sigma\left(r(x,y_w) - r(x,y_l)\right)\right]$$

Direct Preference Optimization: Your Language Model is Secretly a Reward Model by Rafailov et al. 2023 Maximum Entropy Inverse Reinforcement Learning by Ziebart et al. 2010

**Closed-form solution** [Ziebart et al., 2008]:

Rewrite the reward in terms of the policy [Rafailov et al., 2023]:

**Assume reward follows Bradley Terry** 



### MinXent RL

$$\pi_{i+1} = \arg \max_{\substack{\pi \in \Pi \\ y \sim \pi}} \mathbb{E} \left[ r(x, y) \right] - \frac{1}{\eta} D_{\text{KL}}(\pi \parallel \pi_t)$$

$$\forall x, y : \pi_{t+1}(y \mid x) = \frac{\pi_t \exp(\eta r(x, y))}{Z(x)}$$

$$\forall x, y : r(x, y) = \frac{1}{\eta} \left( \ln Z(x) + \ln \left( \frac{\pi_{t+1}(y \mid x)}{\pi_t(y \mid x)} \right) \right)$$

Direct Preference Optimization: Your Language Model is Secretly a Reward Model by Rafailov et al. 2023 Maximum Entropy Inverse Reinforcement Learning by Ziebart et al. 2010

**Closed-form solution** [Ziebart et al., 2008]:

Rewrite the reward in terms of the policy [Rafailov et al., 2023]:

**Assume nothing about the reward** structure, but instead assume access to environment resets.



### **Rebel: Reinforcement learning via regressing relative rewards**

#### Zhaolin Gao



Gokul Swamy



#### Jonathan D. Chang

Wenhao Zhan





#### Owen Oertell

Thorsten Joachims

Jason Lee



J. Andrew Bagnel



#### Wen Sun





### REBEL algorithm overiew

#### At iteration t with policy $\pi_t$

1. Sample (hybrid) data using resets:

$$D_t: \{x, y, y'\} \qquad x \sim D, y \sim \pi_t(\cdot \mid$$

2. Regressing the relative rewards (least squares regression):

$$\pi_{t+1} = \arg\min_{\pi} \mathbb{E}_{D_t} \left( \frac{1}{\eta} \left( \ln \frac{\pi(y \mid x)}{\pi_t(y \mid x)} - \ln \frac{\pi(y' \mid x)}{\pi_t(y' \mid x)} \right) - \left( r(x, y) - r(x, y') \right) \right)^2$$
Predictor
Relative reward

e.g., offline data or reference policy or **best-of-N of**  $\pi_t$ 



49

## Informal Theory of REBEL

$$\pi_{t+1} = \arg\min_{\pi} \mathbb{E}_{D_t} \left( \frac{1}{\eta} \left( \ln \frac{\pi(y \mid x)}{\pi_t(y \mid x)} - \ln \frac{\pi(y' \mid x)}{\pi_t(y' \mid x)} \right) - \left( r(x, y) - r(x, y') \right) \right)^2$$

then we can do as well as any policy that is covered by the training data distributions

$$\forall t, \max_{x,y} \frac{\pi^*(y \mid x)}{\pi_t(y \mid x) + \mu(y \mid x)} \le C < \infty$$

#### **Informal statement:** If we can solve each regression problem well (in-distribution),

## **Experimental Setup**

#### Task Statement

Given a reddit post, write a TL;DR (short summary).

#### **Example Post**

#### **SUBREDDIT:** r/dogs

**TITLE:** [HELP] Not sure how to deal with new people/dogs and my big ole pup **POST:** I have a three year old Dober/Pit mix named Romulus ("Rome" for short). I live with 3 other dogs: a 10 year old labrador, a 2 year old French Bulldog and a 8 year old maltese mix. The four of them get along just fine, Rome and the Frenchie are best best best best friends. He isn't the best at meeting new people, but not ALWAYS....Then, the crux of the matter: I want to have a 4th of July party. Several people want to bring their dogs. I doubt I can say "no dogs allowed" and I don't want to let everyone else bring their dog and make mine stay at day care all day.

Learning to summarization from human feedback Stiennon et al. 2017

#### **Example Human Label**

**TL;DR:** HOW do I introduce new people? HOW do l introduce new dogs? WHAT do l do about 4th of July??



### **Experimental Results: TL;DR**







**2.8B** 



### **Experimental Results**





#### Takeaways:

- Offline methods are more efficient but yield lower win rates.
- 2. REBEL outperforms PPO in both win rate and resource efficiency.
- 3. REBEL outperforms iterative DPO in win rate while matching its efficiency.

# Scaling to larger model (8B) on more modern benchmarks

#### **Experimental Results** Fine-tuning Llama 3-8B model for general chat

**Length-controlled Winrate** 



#### **Dataset:** ultrafeedback [Cui et al]; **Reward Model:** ArMo [Wang et al]



#### **GPT-4 Omni (05/24) REBEL-Llama-3-8B-IT**

#### Scaling to Reasoning Tasks

#### **Experimental Results** Fine-tuning Qwen-2.5B model

Dataset: Math and GSM8K; Reward Model: Verifiable Reward



GSM8K

## Summary

- tasks.
- Empirically, REBEL outperforms PPO in performance, computational efficiency, and memory usage.
- REBEL achieves strong results on standard LLM benchmarks.

REBEL reduces the RL problem to a sequence of relative reward regression

#### How can we develop efficient algorithms for solving the MinxEnt RL objective, assuming access to environment resets?

- Resetting with reference policy (to boost exploration)

- Resetting with demonstration data (to boost exploration) - Resetting with the current policy (to reduce computation)

### Acknowledgement























