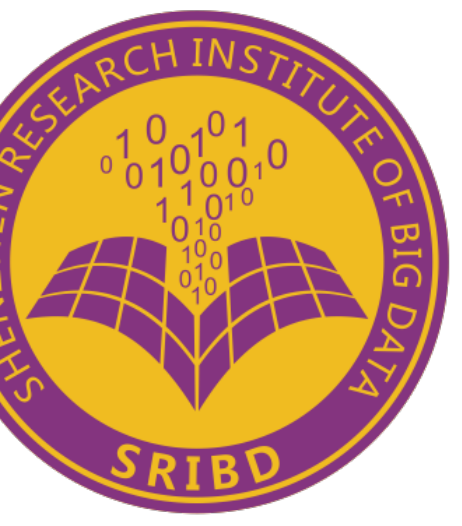




# DIVERGENCE-AUGMENTED POLICY OPTIMIZATION



QING WANG<sup>4</sup>, YINGRU LI<sup>1,2</sup>, JIECHAO XIONG<sup>5</sup>, TONG ZHANG<sup>2,3</sup>

[1] CUHK(SZ) [2] SRIBD [3] HKUST [4] HUYA AI [5] TENCENT AI LAB

## PREVIEWS

We are interested in designing **conservative update** to ensure **stable** policy learning when **off-policy** data are reused.

1. We propose to stabilize policy improvement by constraining the **discounted state-action visitation** induced by consecutive policies to be close to one another.
2. We demonstrate the benefits of our approach through an empirical comparison on Atari games where the reuse of off-policy data is necessary.

## PROBLEM AND CHALLENGES

### Reinforcement learning paradigm:

- Learning agent that interacts repeatedly with the environment.
- The environment is formalized by a *Markov decision process*  $(\mathcal{S}, \mathcal{A}, \mathbb{P}, r, \gamma, \rho)$ .
- Goal: how to map states to distribution over actions (policy learning) in order to maximize the long term cumulative reward.

### Policy optimization

Consider some class of **parametric** stochastic policy  $\{\pi_\theta, \theta \in \Theta\}$ .

**Objective:**  $J(\theta) = (1 - \gamma)\mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r_t | \rho, \pi_\theta]$

- **Policy gradient:** optimize directly the cumulative rewards
    - Reinforce (Williams 88', 92')
    - Actor-critic methods (Sutton et al. 00')
$$\widehat{\nabla}_\theta J(\theta) = \sum_{t=0}^{\tau} \gamma^t r_t \sum_{t=0}^{\tau} \nabla \log \pi_\theta(a_t | s_t).$$

$$\theta \leftarrow \theta + \eta \widehat{\nabla}_\theta J(\theta),$$
  - **Large variance issue! ✗**
- **Conservative approaches:** optimize surrogate objective that can provide local improvement.
    - Conservative policy Iteration (CPI) (Kakade & Langford 02').
    - Trust region policy optimization (TRPO) (Shulman et al. 15')
    - Proximal policy optimization (PPO) (Shulman et al. 17')
    - **Need to ensure that the new policy stays in the vicinity of the current policy!**

## CONSERVATIVE OPTIMIZATION

- Discounted state-visitation:

$$d_\rho(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s | \rho, \pi).$$

- Discounted state-action-visitation:

$$\mu(s, a) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s, a_t = a | \rho, \pi).$$

- **Performance difference lemma:**

$$J(\pi') = J(\pi) + \mathbb{E}_{s \sim d_\rho} \mathbb{E}_{a \sim \pi'(\cdot|s)} [A^\pi(s, a)].$$

- Implementation of the **vicinity requirement:**

- **Surrogate objective:** given a current policy  $\pi_i$ ,
- $$L_{\pi_i}(\pi') = J(\pi_i) + \mathbb{E}_{s \sim d_{\rho}^{\pi_i}} \mathbb{E}_{a \sim \pi'(\cdot|s)} [A^{\pi_i}(s, a)].$$

- CPI returns a stochastic mixture:

$$\pi_{i+1} = \alpha_i \pi_i^+ + (1 - \alpha_i) \pi_i,$$

where  $\pi_i^+ = \arg \max_{\pi'} L_{\pi_i}(\pi')$ .

- TRPO solves a KL-constrained problem:

$$\pi_{i+1} = \arg \max_{\pi'} L_{\pi_i}(\pi')$$

$$\text{s.t } \mathbb{E}_{s \sim d_{\rho}^{\pi_i}} [\text{KL}(\pi'(s) || \pi_i(s))] \leq \delta.$$

- PPO optimizes a clipped objective for ensuring target policy stays near to current policy.

## THEORETICAL INSIGHTS

**Lower bound on policy's performance:** for any two policies  $\pi$  and  $\pi'$

### Long term effects matters

In terms of action probabilities

$$J(\pi') \geq L_\pi(\pi') - \frac{2\gamma\epsilon^\pi}{1-\gamma} \mathbb{E}_{s \sim d_\pi} [D_{\text{TV}}(\pi'(s) || \pi(s))].$$

Pay extra term of **planning horizon!**

- Used in CPI, TRPO and PPO.
- **Looser lower bound!**
- **Instability in long horizon problem!**

In terms of discounted state-action-visitation

$$J(\pi') \geq L_\pi(\pi') - \epsilon^\pi D_{\text{TV}}(\mu_{\pi'} || \mu_\pi).$$

- **Do not suffer long-horizon factor. ✓**
- **Requires state distribution! ✗**

## DIVERGENCE-AUGMENTED POLICY OPTIMIZATION (DAPO)

### Mirror descent formulation of policy optimization

We seek to solve policy optimization problem by looking at **state-action joint space**:

$$\max_{\pi} J(\pi) = \mathbb{E}_{(s,a) \sim \mu_\pi} [r(s, a)] = \langle \mu_\pi, r \rangle$$

Then we try to solve the problem by mirror descent (proximal iterative scheme) to promote interior point to prevent premature convergence and instability

$$\mu_{t+1} = \arg \min_{\mu \in \Delta_\Pi} D_F(\mu_\pi, \mu_t) + \eta \langle -r, \mu_\pi \rangle \quad (1)$$

Potential function as **Conditional Negentropy**  $F(\mu) = \sum_{s,a} \mu(s, a) \log \pi_\mu(a|s)$ , we get

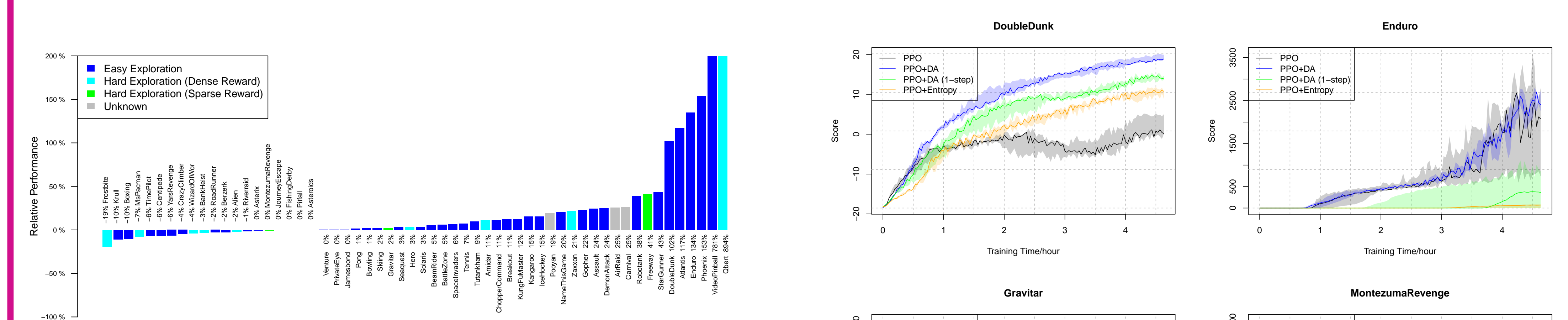
$$D_F(\mu', \mu) = D'_{\text{KL}}(\mu', \mu) = \sum_{s,a} \mu'(s, a) \log \frac{\pi'(a|s)}{\pi(a|s)}$$

### Divergence-augmented pseudo reward

Optimizing equation 1 with conditional negentropy via policy gradient is equivalent to ordinary policy optimization using **pseudo reward** at each state-action pair: ( $\pi_t$  is the a current data-generating policy)

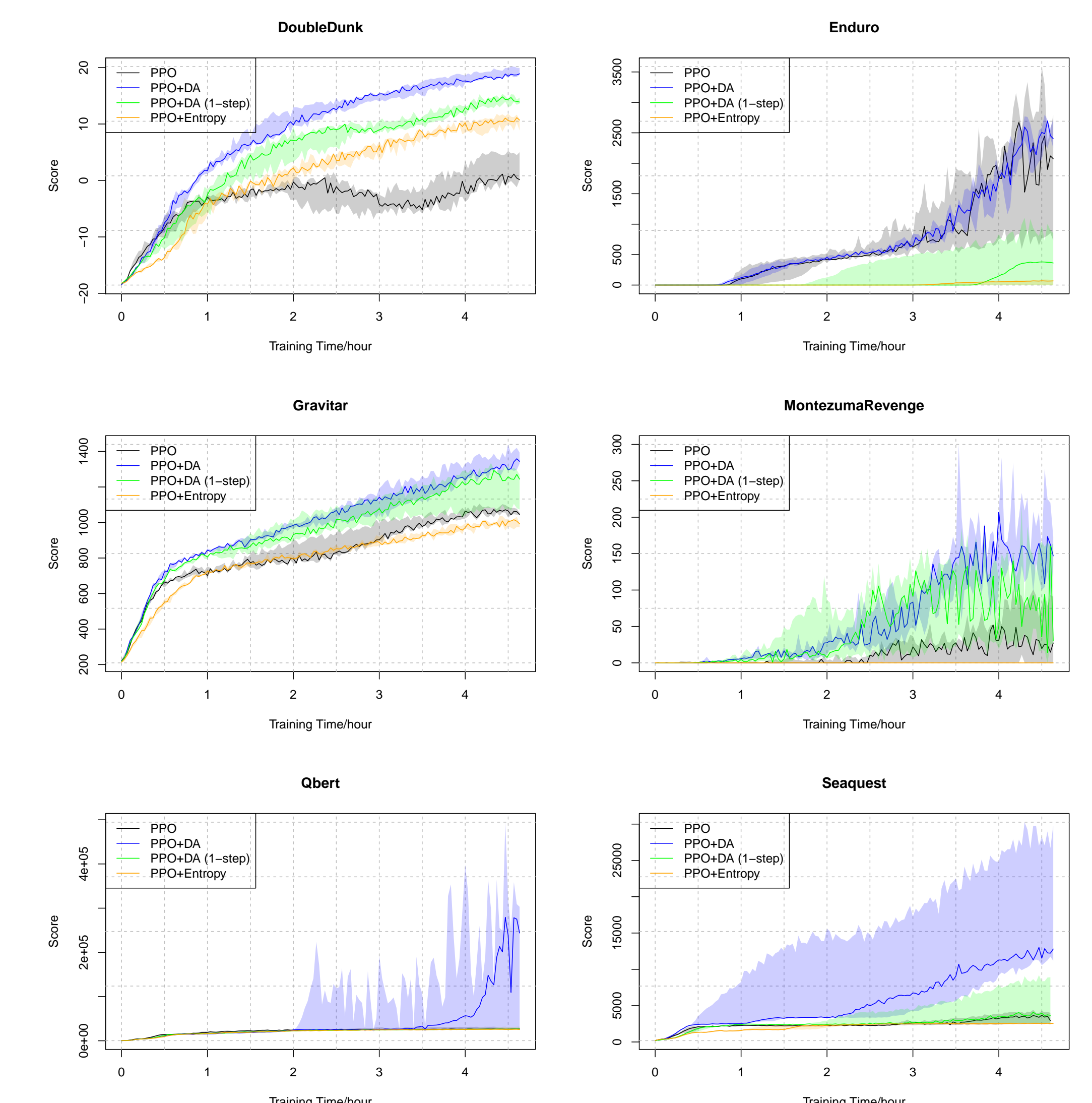
$$r(s, a) - \frac{1}{\eta} \log \frac{\pi(a|s)}{\pi_t(a|s)}$$

## SO WHAT?



**Figure 1:** The relative score  $\frac{\text{proposed} - \text{baseline}}{\max(\text{human}, \text{baseline}) - \text{random}}$  of DAPO (our proposed) v.s. PPO (baseline) in all 58 atari games.

- In **Figure 1**, our algorithms achieve better performance in a large fraction of game environments.
- In **Figure 2**, the performance of **PPO**, **PPO+DA**, **PPO+DA (1-step)**, and **PPO+Entropy** are plotted in different colors.
- **DAPO** corresponds to multi-step divergence augmentations (**PPO+DA**), which performs best.
- Entropy regularization is not a good choice. PPO performs the worst in the comparison.



**Figure 2:** Comparison on selected environments.

## MORE INFORMATION

Paper link: [tinyurl.com/78wmp42w](https://tinyurl.com/78wmp42w)  
Code link: [github.com/lms/dapo](https://github.com/lms/dapo)

WeChat: RichardYRLi  
Personal site: [richardli.xyz](https://richardli.xyz)  
twitter: @RichardYRLi

