

Distributional Policy Optimization: An Alternative Approach for Continuous Control



Chen Tessler*, Guy Tennenholtz* and Shie Mannor

Technion Institute of Technology, Israel

Continuous Control

Setting:

ullet MDP $(\mathcal{S}, \mathcal{A}, P, r, \gamma)$, where \mathcal{A} is continuous.

Objective:

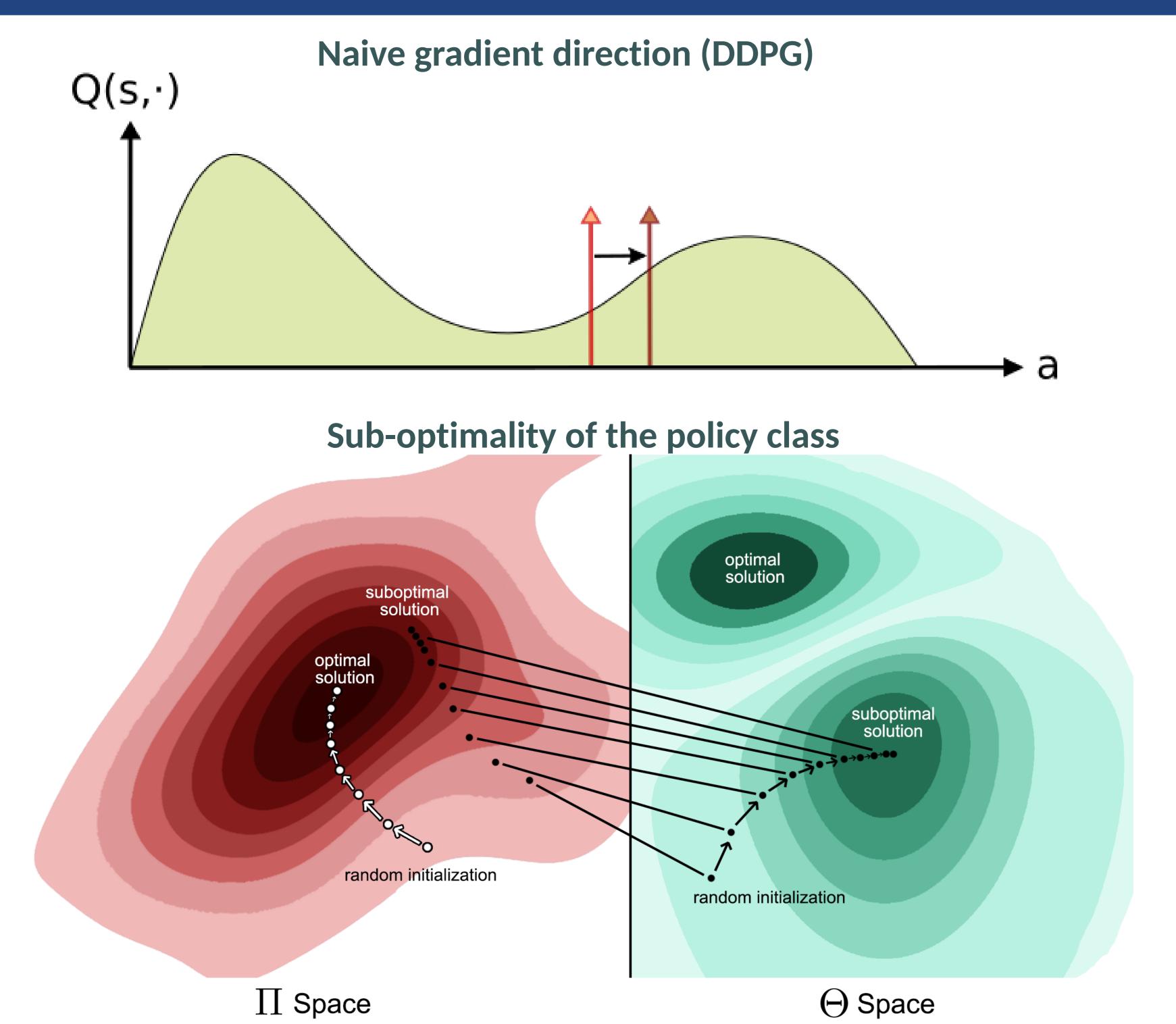
Find π^* which satisfies

$$\pi^*(s) \in rg \max_{\pi \in \Pi} Q^\pi(s,\pi(s)) = rg \max_{\pi \in \Pi} \mathbb{E}^\pi \left[\gamma^t r_t | s_0 = s
ight]$$

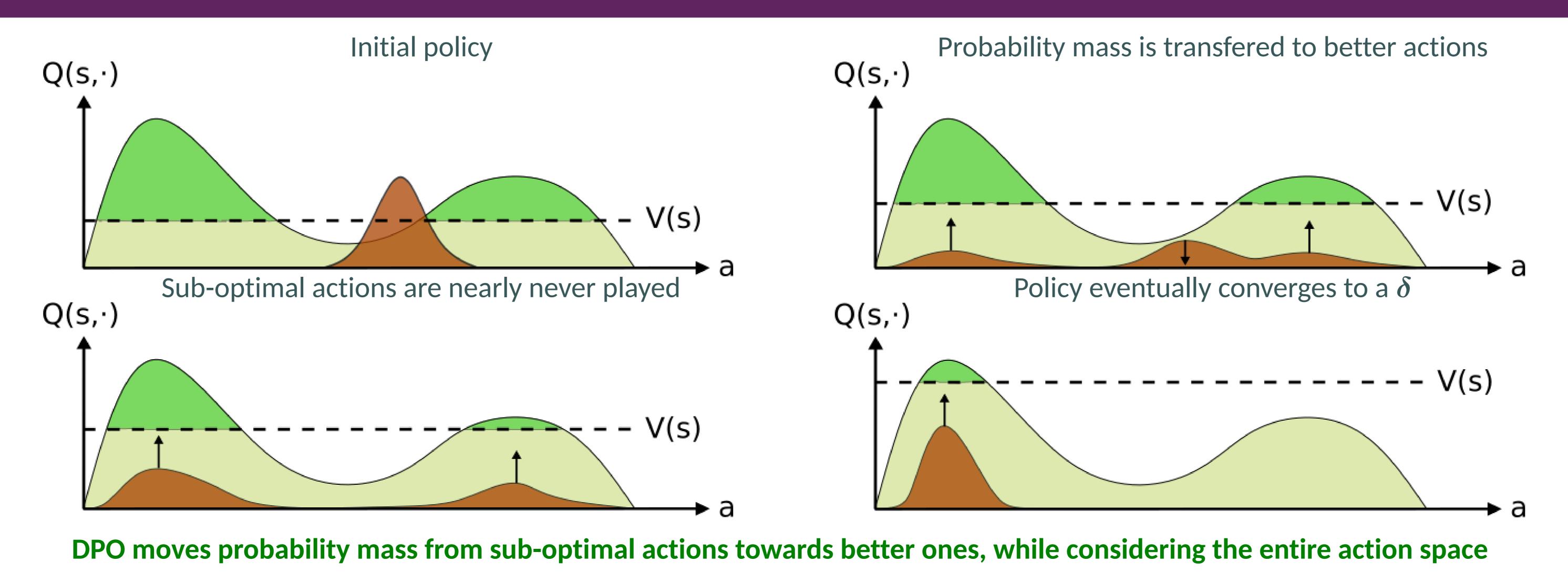


In this work we provide a framework for tackling continuous control problems with non-convex returns

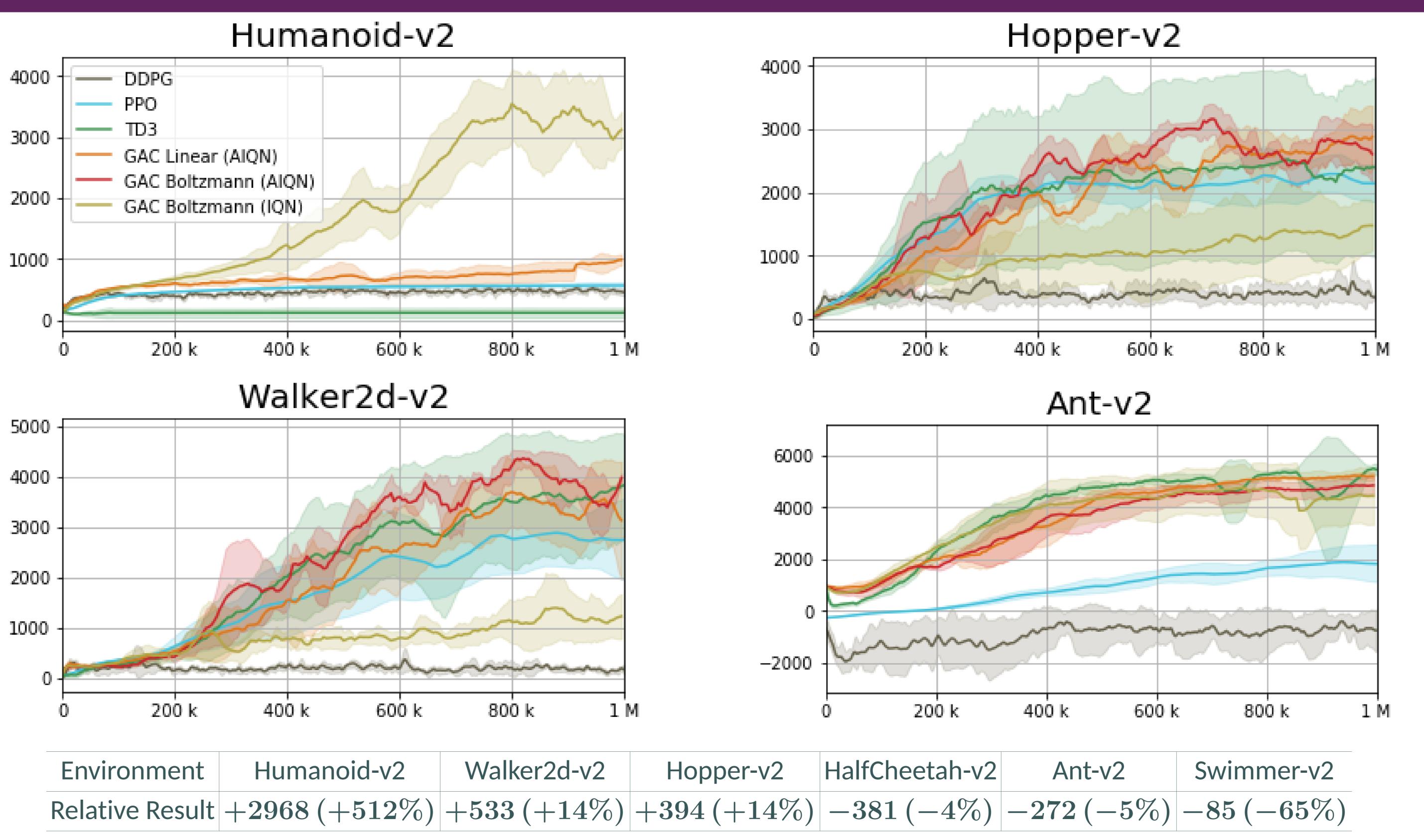
Pitfalls of Current Approaches



Distributional Policy Optimization (Illustration)



Experiments



Relative best GAC results compared to the best policy gradient baseline

Distributional Policy Optimization (Algorithm)

1: Input: learning rates $\alpha_k\gg\beta_k\gg\delta_k$ 2: $\pi_{k+1}=\Gamma\left(\pi_k-\alpha_k
abla_\pi d(\mathcal{D}_{I^{\pi'_k}}^{\pi'_k},\pi)\mid_{\pi=\pi_k}
ight)$ 3: $Q_{k+1}^{\pi'}(s,a)=Q_k^{\pi'}(s,a)+\beta_k\left(r(s,a)+\gamma v_k^{\pi'}(s)-Q_k^{\pi'}(s,a)\right)$ 4: $v_{k+1}^{\pi'}(s)=v_k^{\pi'}+\beta_k\int_{\mathcal{A}}\left(Q_k^{\pi'}(s,a)-v_k^{\pi'}(s)\right)$ 5: $\pi'_{k+1}=\pi'_k+\delta_k(\pi_k-\pi'_k)$

Summary

- 1. We proposed the Distributional Policy Optimization (DPO) framework for tackling non-convex returns.
- 2. DPO requires the ability to represent arbitrary policies and optimize by minimizing the distance to a target distribution.
- 3. We achieve this by modeling the actor (policy) using an Autoregressive Implicit Quantile Network, a generative model.
- 4. Empirical tests attain results competitive to policy gradient methods, while remaining as sample efficient.

Future Work and Possible Extensions

- 1. In this work we modeled the policy using an Autoregressive Implicit Quantile Network, however it can also be modeled using other methods such as GANs, VAEs and Normalizing Flows.
- 2. While the approach is efficient in the number of samples (environment interactions) it is not as computational efficient as alternative methods.
- 3. As the optimization considers a target probability distribution, this can be leveraged for improving exploration and sample efficiency using UCB-like methods.

Additional Details

- 1. Email: chen.tessler@campus.technion.ac.il, guytenn@gmail.com
- 2.Code: github.com/tesslerc/GAC