

# Nonparametric Policy Improvement in Continuous Action Spaces via Expert Demonstrations

**Enrique Mallada**



**JOHNS HOPKINS**  
UNIVERSITY

**Inform**s

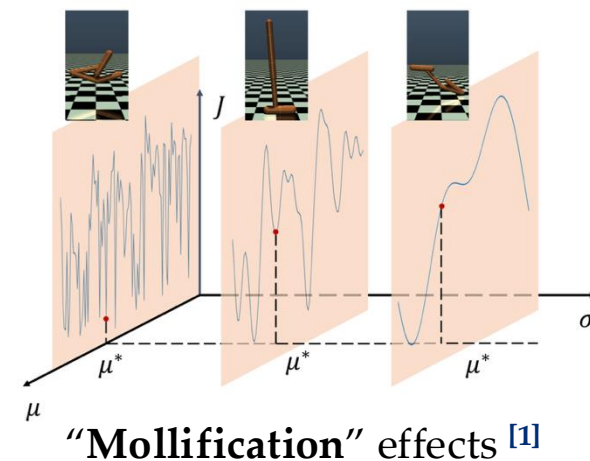
October 26<sup>th</sup>, 2025

# Challenges of “modern” Policy Optimization (P.O.)



## P.O. in continuous spaces:

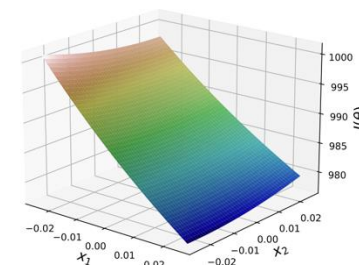
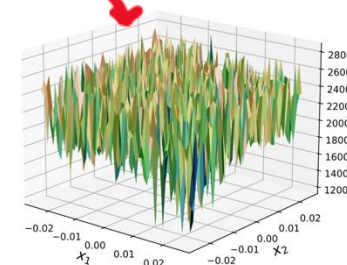
- Largely based on **policy gradients**.
- Choice of **parametrization**:  $a \sim \mathcal{N}(\mu_\theta(s), \sigma_\theta^2(s))$ 
  - Limits expressivity.
  - Local improvement.
  - May yield non-smooth landscapes.



## P.O. in finite spaces was great!

- Policy Iteration = Policy eval. + Policy improvement.  
$$\pi \xrightarrow{\text{eval.}} Q^\pi(s, a) \quad \pi'(s) \in \operatorname{argmax}_{a \in \mathcal{A}} Q^\pi(s, a)$$

Fractal!



Non-smooth landscapes [2]

- Monotonic improvement **everywhere**  $V^{\pi'}(s) \geq V^\pi(s) \quad \forall s \in \mathcal{S}.$

[1] Tao Wang, Sylvia Hebert, Sicun Gao, Mollification effects of policy gradient, ICML 24

[2] —, Fractal landscapes in policy optimization, NeurIPS 23



Could we get:

1. **Benefits of policy iteration**
2. **Avoid drawbacks of gradient methods?**

# Acknowledgements



**Agustin Castellano**



**Sohrab Rezaei**



**Jared Markowitz**



# Problem Setup



**Goal:** find optimal policy

$$\max_{\theta} J(\theta) := E_{s_0 \sim \rho, a_0 \sim \pi_{\theta}(s_0)} \left[ Q^{\pi_{\theta}}(s_0, a_0) \right]$$



# Problem Setup

**Goal:** find optimal **nonparametric** policy

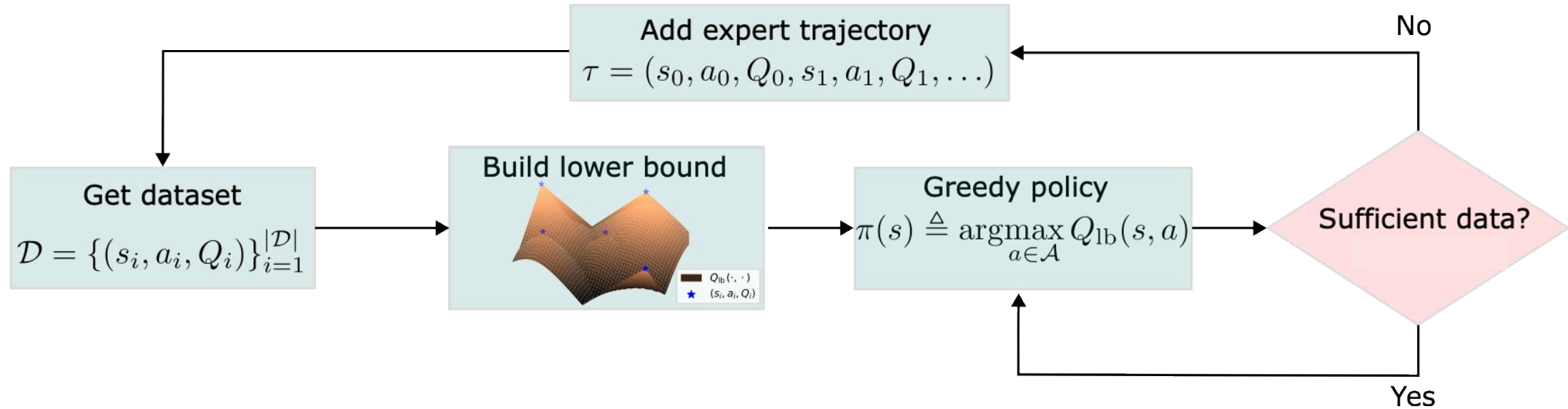
$$\max_{\mathcal{D}} J(\pi_{\mathcal{D}}) := E_{s_0 \sim \rho, a_0 \sim \pi_{\mathcal{D}}(s_0)} [Q^{\pi_{\mathcal{D}}}(s_0, a_0)]$$

**Data set:**  $\mathcal{D} = \{(s_i, a_i, Q_i)\}_{i=1}^{|\mathcal{D}|}$        $Q_i := \sum_t \gamma^t r(s_t, a_t)$

**Assumptions:**

1. How can we leverage these transitions to learn a policy?
2. What guarantees can we get when we add more transitions?
3. Where should we add transitions to improve performance?
- Expert data: we have  $\mathcal{D} = \{(s_i, a_i, Q_i)\}_{i=1}^{|\mathcal{D}|}$ , where  $s_i \sim \rho, a_i \sim \pi^*(s_i), Q_i = Q^*(s_i, a_i)$

# Overview of our method





**1. How can we use these transitions to learn a nonparametric policy?**



# Building bounds & Nonparametric Policy



**Expert data:** we have  $\mathcal{D} = \{(s_i, a_i, Q_i)\}_{i=1}^{|\mathcal{D}|}$ , where:  $a_i = \pi^*(s_i)$ ;  $Q_i = Q^*(s_i, a_i)$

- Use data to define **lower bounds** on optimal values:

$$V_{\text{lb}}(s) \triangleq \max_{1 \leq i \leq |\mathcal{D}|} \{Q_i - L \cdot d_{\mathcal{S}}(s, s_i)\} \quad Q_{\text{lb}}(s, a) \triangleq \max_{1 \leq i \leq |\mathcal{D}|} \{Q_i - L \cdot (d_{\mathcal{S}}(s, s_i) + d_{\mathcal{A}}(a, a_i))\}$$

- **Nonparametric Policy:**

$$\pi(s) \triangleq \operatorname{argmax}_{a \in \mathcal{A}} Q_{\text{lb}}(s, a) = a_{i'}$$

- **Remark:** Note argmax always gives actions in dataset  $(s_{i'}, a_{i'}, Q_{i'})$
- **Question:** What can we say about  $V^\pi(s)$  ?

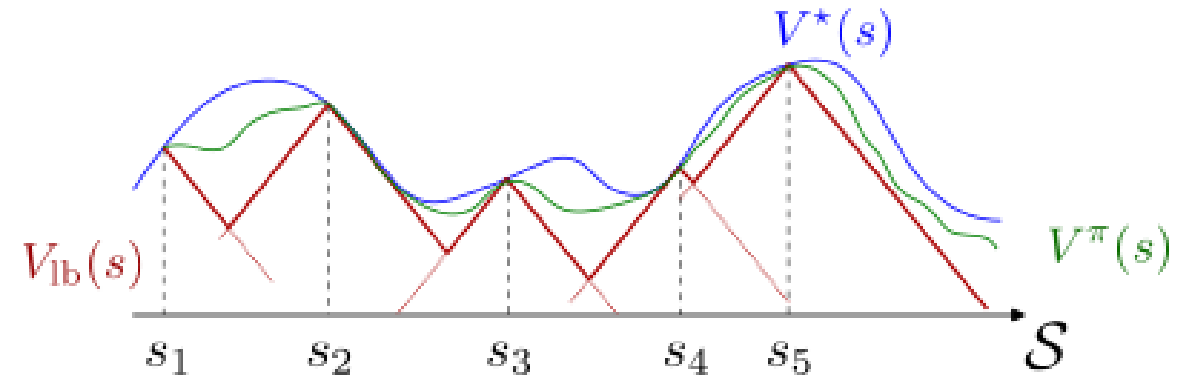


# Nonparametric policy *improves* over lower bound

## Policy Evaluation<sup>\*</sup>:

- Nonparametric  $\pi$  satisfies  $\forall s \in \mathcal{S}$ :

$$V_{lb}(s) \leq V^\pi(s) \leq V^*(s)$$



## Policy Improvement<sup>\*</sup>:

- Given data sets  $\mathcal{D}, \mathcal{D}'$  with  $\mathcal{D} \subset \mathcal{D}'$

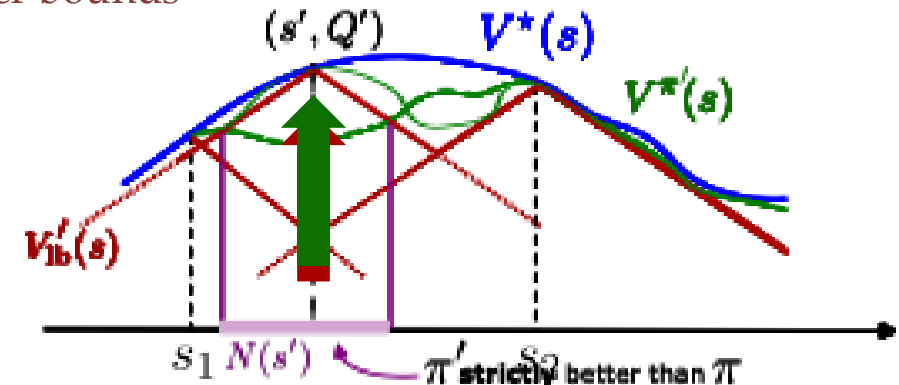
$$V_{lb}(s) \leq V'_{lb}(s) \quad \forall s \in \mathcal{S}$$

More data = better lower bounds

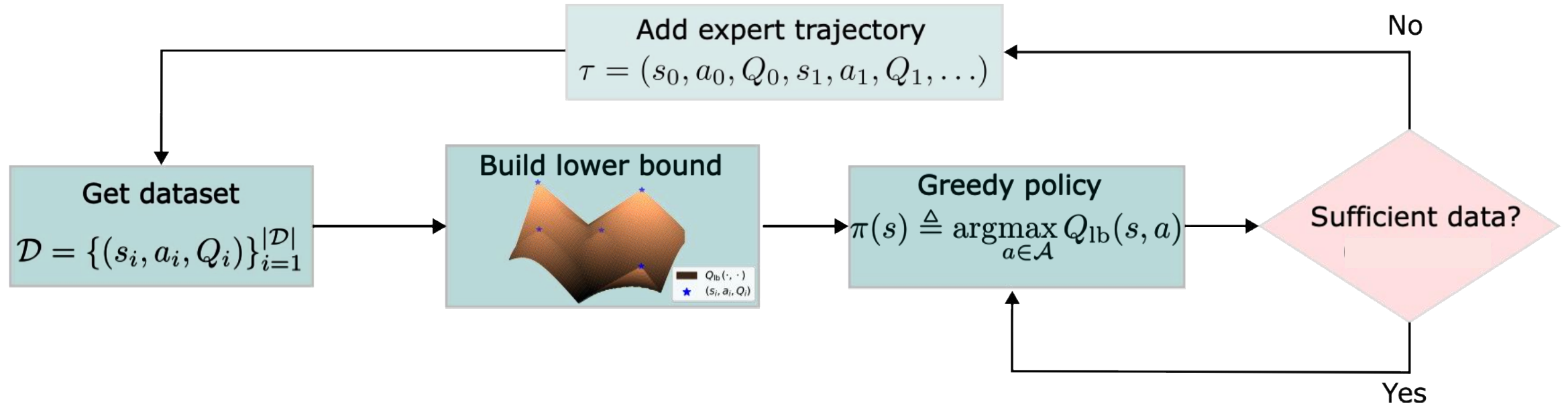
$$V^\pi(s') \leq V^{\pi'}(s') \quad \forall s' \in \mathcal{D}' \setminus \mathcal{D}$$

Improvement on added points

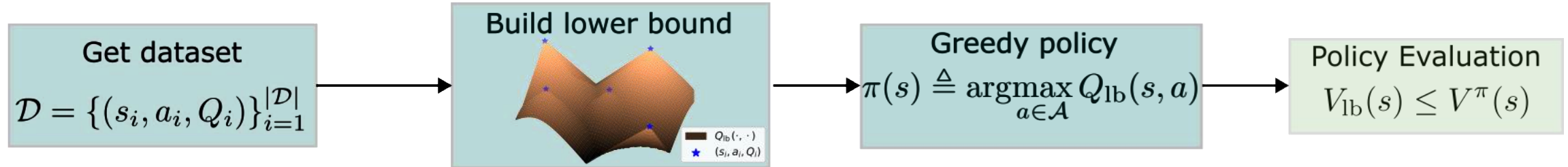
- Strict** on neighbors of new data:  $\forall s \in N(s')$



(<sup>\*</sup>:Data must come from trajectories)



## 1. How to learn a policy?



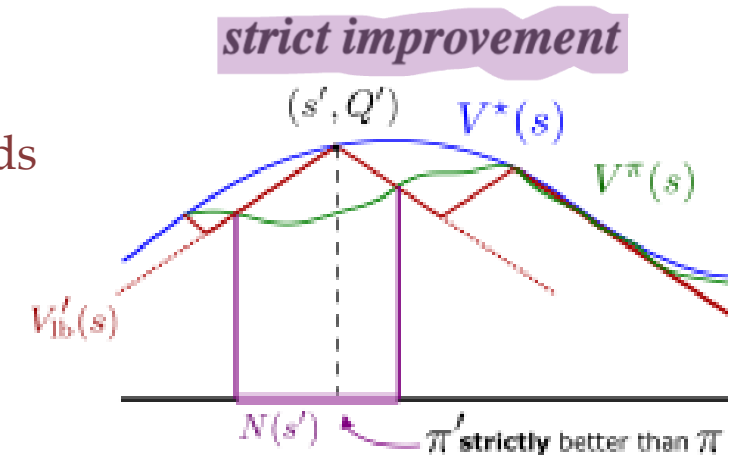
## 2. What guarantees with more transitions?

More data = better lower bounds

$$V_{\text{lb}}(s) \leq V'_{\text{lb}}(s) \quad \forall s \in \mathcal{S}$$

Improvement on added points

$$V^{\pi}(s') \leq V^{\pi'}(s') \quad \forall s' \in \mathcal{D}' \setminus \mathcal{D}$$



## 3. Where to add transitions?

- Only *where* **sufficient improvement** is guaranteed:  $\Delta(s) = V^*(s) - V^{\pi}(s)$

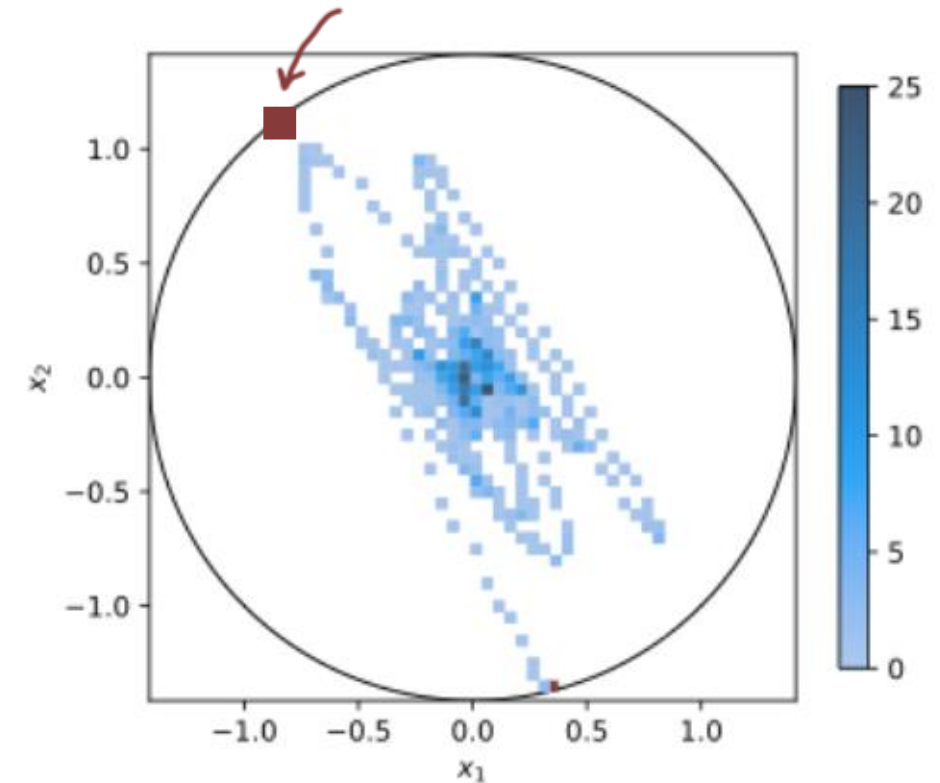
# Algorithm

**Input:** Lipschitz constant  $L$

**For** each episode **do**:

1. Sample  $s \sim \rho(\mathcal{S}_0)$
2. If  $\Delta(s) > \varepsilon$  :
  - Run optimal trajectory with  $\pi^\star$ 

$$\tau = (s_0, a_0, Q_0, s_1, a_1, Q_1, \dots)$$
  - **Repeat:** add tuples to dataset (from  $i = 0$ )
 
$$\mathcal{D} \leftarrow \mathcal{D} \cup \{(s_i, a_i, Q_i)\}$$
  - Until:**  $\Delta(s_i) \ll \varepsilon$ .
3. Else:
  - **Continue**



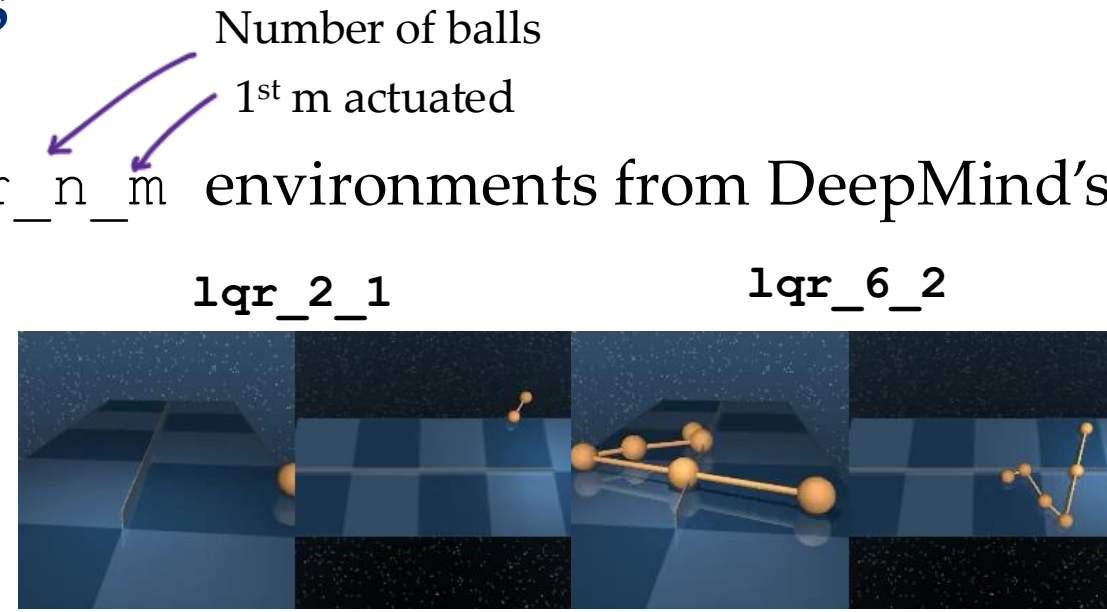
## Convergence guarantees?

$$\sup_{s \in \mathcal{S}_0} |V^\star(s) - V^\pi(s)| \leq \varepsilon$$

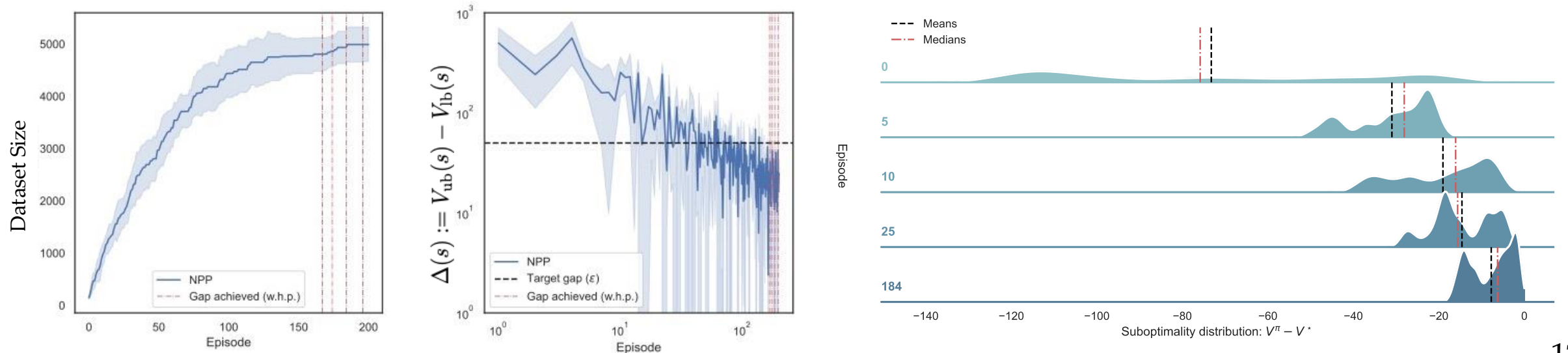
w.p.  $\geq 1 - \delta$  in  $\mathcal{O} \left( N_{\text{cover}} \left( \frac{\varepsilon}{4L} \right) \log N_{\text{cover}} \left( \frac{\varepsilon}{4L} \right) \log \frac{1}{\delta} \right)$  episodes.

# Experiments

- We use the  $lqr\_n\_m$  environments from DeepMind's Control Suite



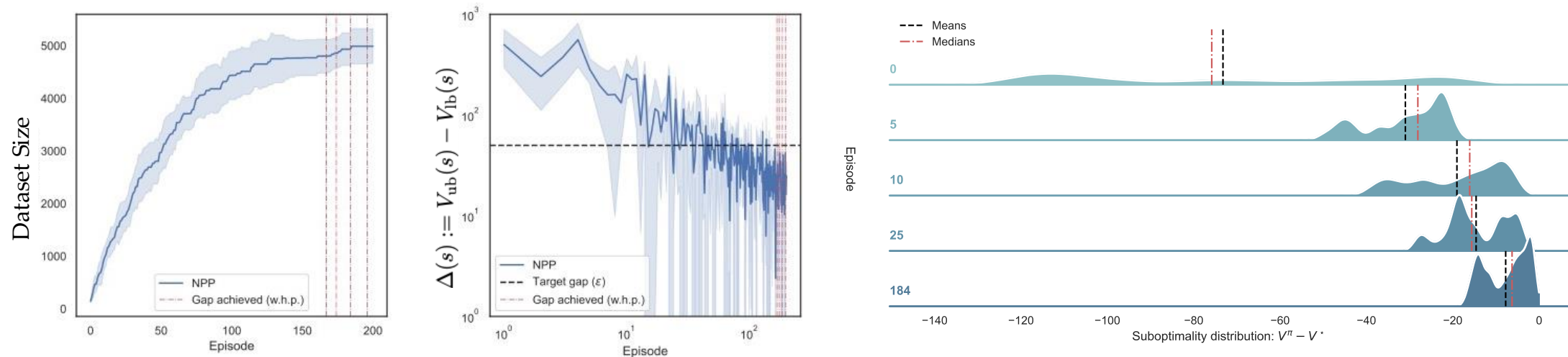
- Results on  $lqr\_2\_1$ :



# Experiments



- We use the `lqr_n_m` environments from DeepMind's Control Suite
- Results on `lqr_2_1`:

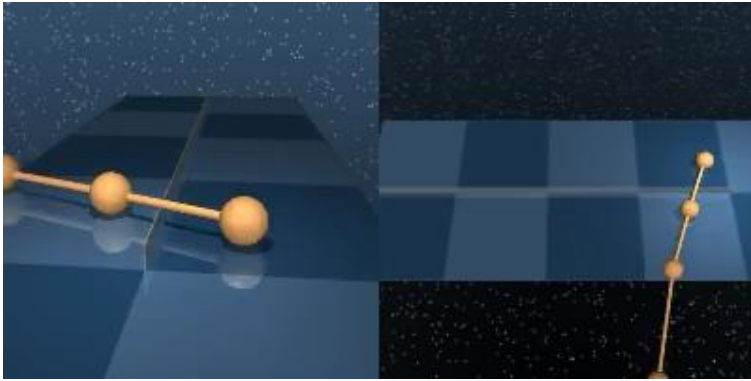


- **Remarks:**
  - **Incremental learning:** No catastrophic forgetting.
  - **Improvement across entire state space** (not in expectation).
  - **Only valuable data is added** (harder to find at times passes)

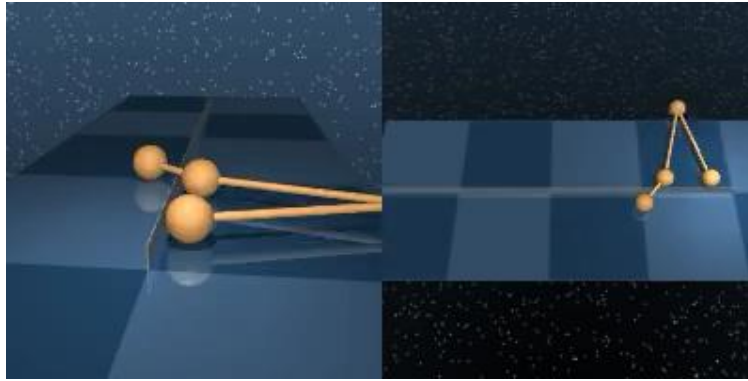
# Incremental Learning



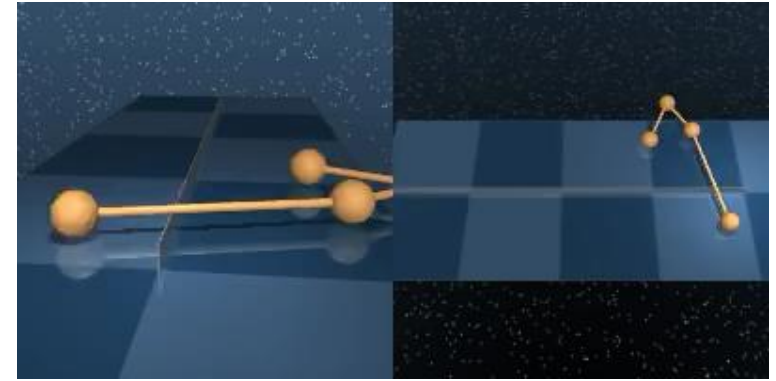
after 10 episode...



after 100 episode...



after 1000 episodes...



after 30K+



optimal control





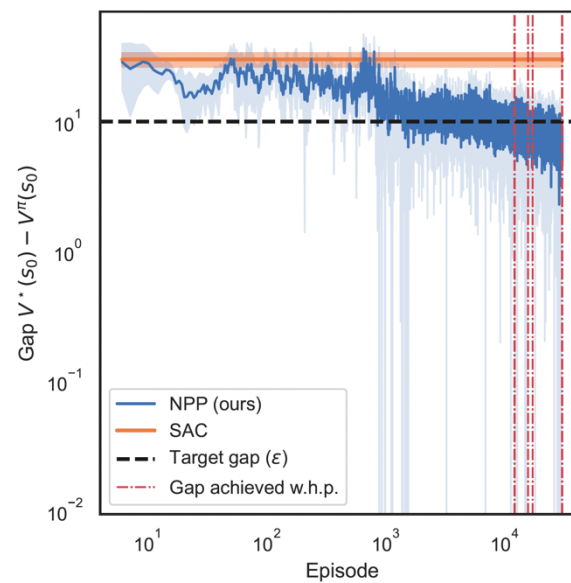
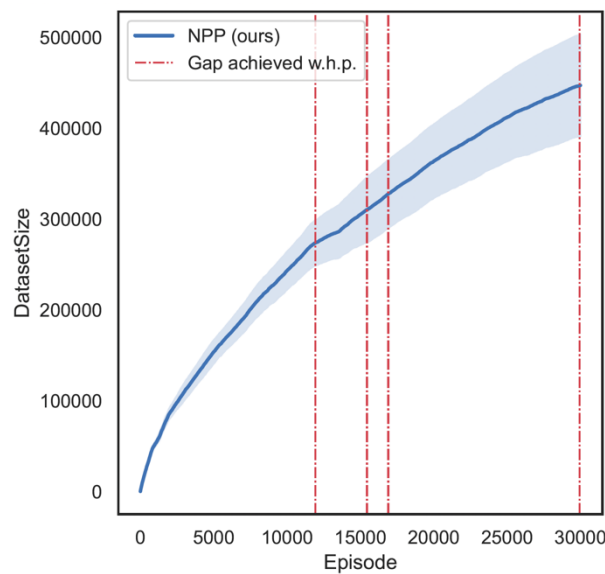
# Incremental Learning



after 30K+



optimal control





# Conclusions

- Proposed **nonparametric policy** based on **expert demonstrations**
- Policy is **greedy** w.r.t. **lower bound** on  $Q^*$ , satisfies:
  - i) **policy evaluation inequality** (everywhere)
  - ii) (strict) **policy improvement** (on new data)
- **Data** collection only **where it's needed**.
- Experiments show **incremental learning, no catastrophic forgetting**

## Future work

- Sub-expert demonstrations.
- Bootstrapping with lower bounds.
- Stochastic MDPs.



# Thank you!

## Questions?



Check out our paper



Check out our Github repo