



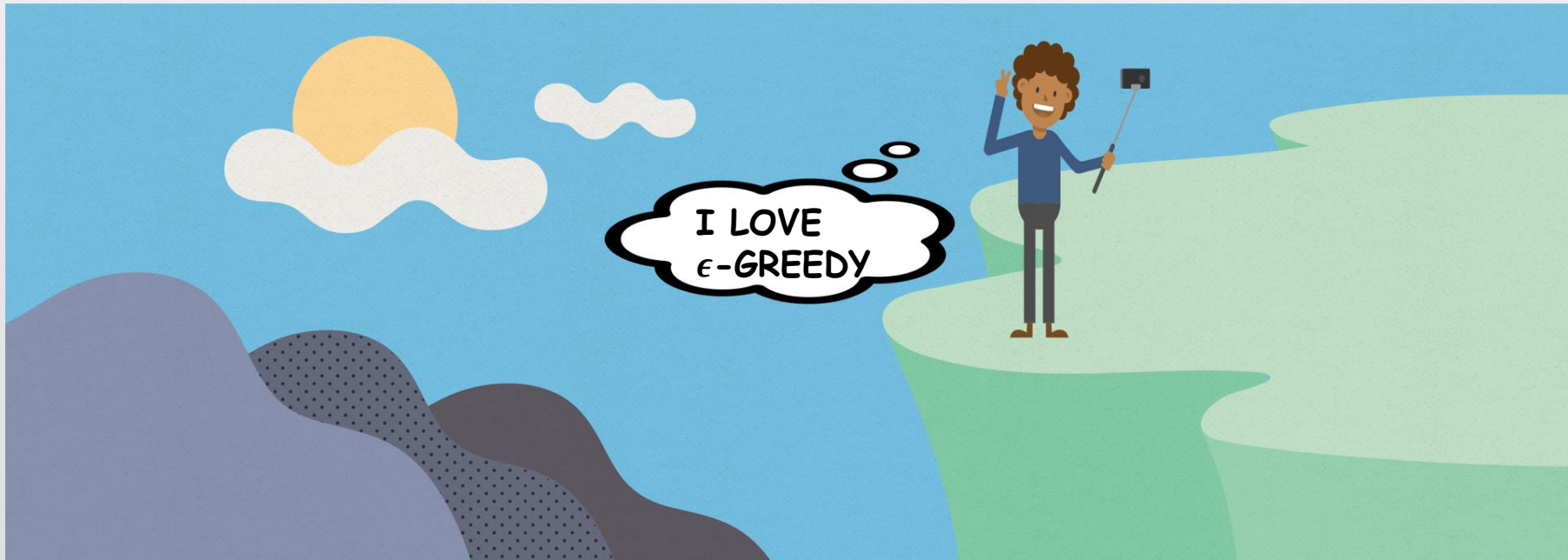
# Exploration Conscious Reinforcement Learning Revisited

Lior Shani\*   Yonathan Efroni\*   Shie Mannor

Technion Institute of Technology

# Why?

- To learn a good policy, an RL agent must explore!
- However, it can cause hazardous behavior during training.



---

# Why?

- To learn a good policy, an RL agent must explore!
- However, it can cause hazardous behavior during training.





# Exploration Conscious Reinforcement Learning

- **Objective: Find the optimal policy knowing that exploration might occur**
- **For example :  $\epsilon$ -greedy exploration ( $\alpha = \epsilon$ )**

$$\pi_{\alpha}^* \in \operatorname{argmax}_{\pi \in \Pi} \mathbb{E}^{(1-\alpha)\pi + \alpha\pi_0} \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)$$



# Exploration Conscious Reinforcement Learning

- **Objective: Find the optimal policy knowing that exploration might occur**
- **For example :  $\epsilon$ -greedy exploration ( $\alpha = \epsilon$ )**

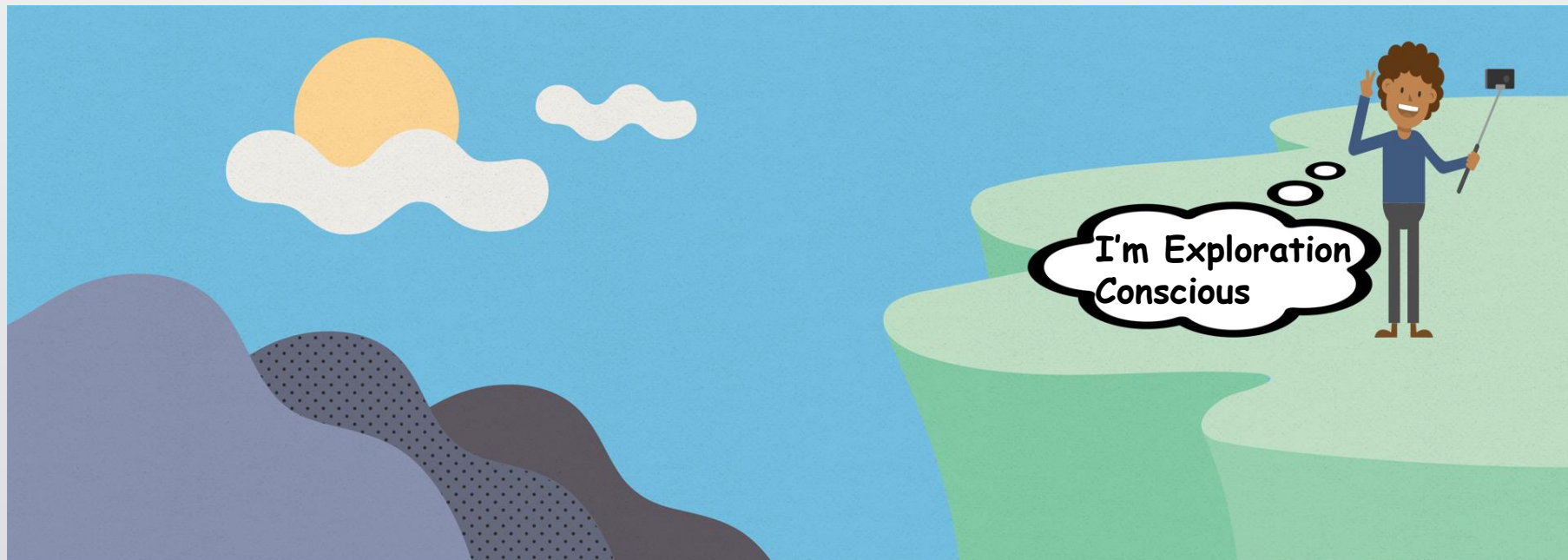
$$\pi_{\alpha}^* \in \operatorname{argmax}_{\pi \in \Pi} \mathbb{E}^{(1-\alpha)\pi + \alpha\pi_0} \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)$$

- **Solving the Exploration-Conscious problem = Solving an MDP**
- **We describe a bias-error sensitivity tradeoff in  $\alpha$**



# Exploration Conscious Reinforcement Learning

- **Objective: Find the optimal policy knowing that exploration might occur**



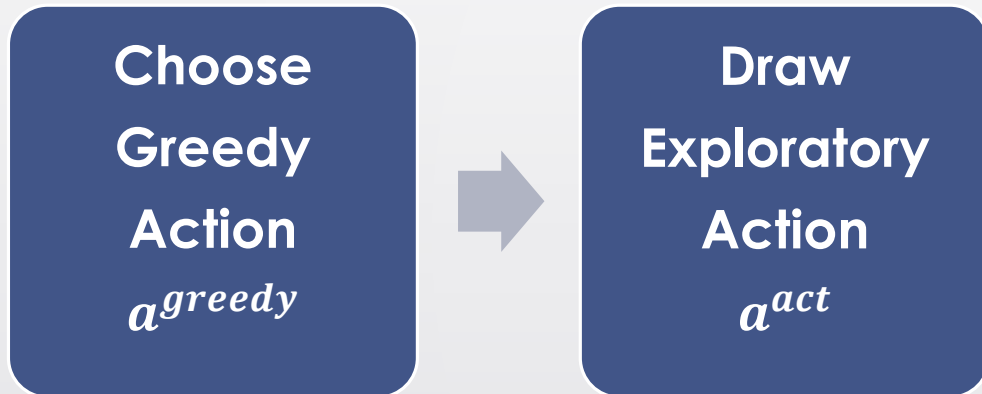


## Fixed Exploration Schemes (e.g. $\epsilon$ -greedy)

Choose  
Greedy  
Action  
 $a^{greedy}$

- $a^{greedy} \in \operatorname{argmax}_a Q^{\pi^a}(s, a)$

# Fixed Exploration Schemes (e.g. $\epsilon$ -greedy)



- $a^{greedy} \in \operatorname{argmax}_a Q^{\pi^\alpha}(s, a)$
- **For  $\alpha$ -greedy:**  $a^{act} \in \begin{cases} a^{greedy} & \text{w.p. } 1 - \alpha \\ \pi_0 & \text{else} \end{cases}$

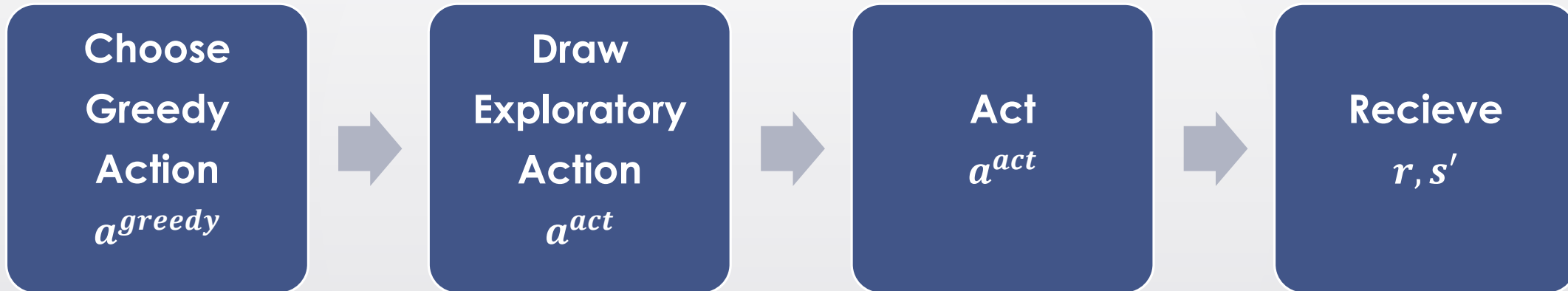


## Fixed Exploration Schemes (e.g. $\epsilon$ -greedy)



- $a^{greedy} \in \operatorname{argmax}_a Q^{\pi^a}(s, a)$
- **For  $\alpha$ -greedy:**  $a^{act} \in \begin{cases} a^{greedy} & \text{w.p. } 1 - \alpha \\ \pi_0 & \text{else} \end{cases}$

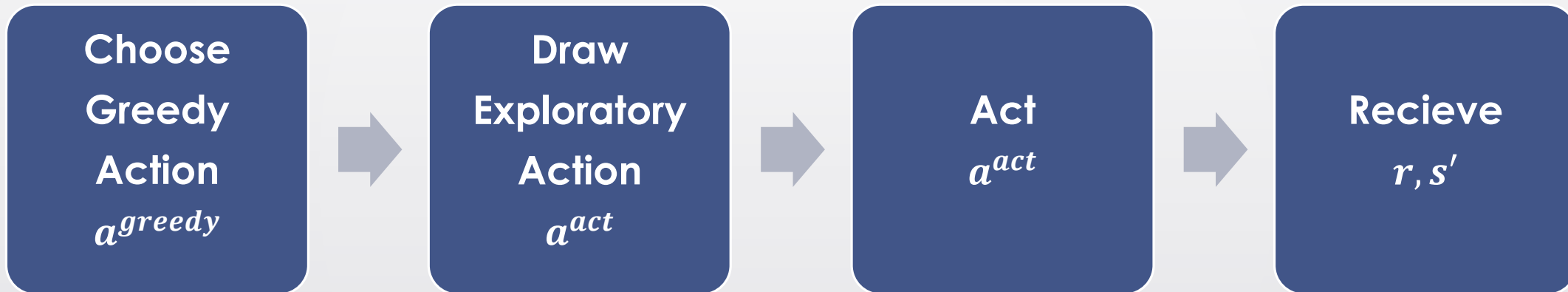
## Fixed Exploration Schemes (e.g. $\epsilon$ -greedy)



- $a^{greedy} \in \operatorname{argmax}_a Q^{\pi^a}(s, a)$
- **For  $\alpha$ -greedy:**  $a^{act} \in \begin{cases} a^{greedy} & \text{w.p. } 1 - \alpha \\ \pi_0 & \text{else} \end{cases}$

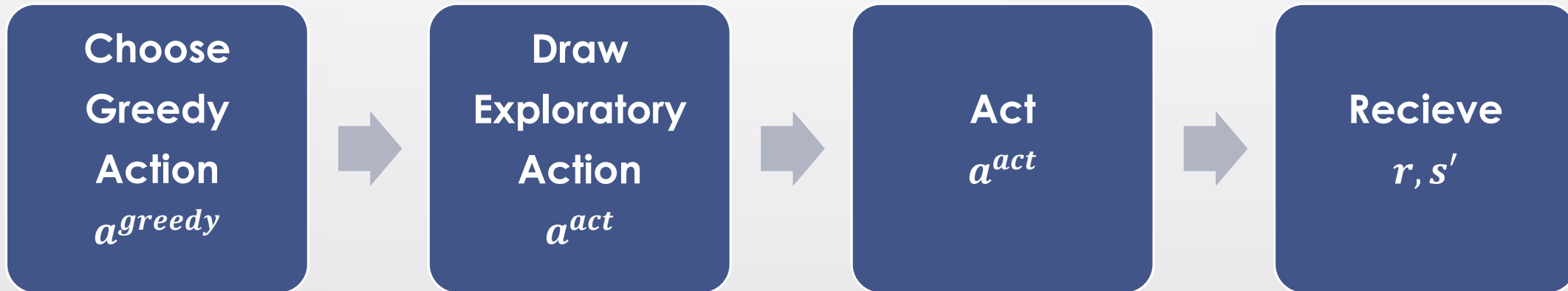


## Fixed Exploration Schemes (e.g. $\epsilon$ -greedy)




- **Normally used information:**  $(s, a^{act}, r, s')$

## Fixed Exploration Schemes (e.g. $\epsilon$ -greedy)




- Normally used information:  ~~$(s, a^{act}, r, s')$~~
- Using information about the exploration process:  $(s, a^{greedy}, a^{act}, r, s')$



## Two Approaches – Expected approach

1. Update  $Q^{\pi^\alpha}(s_t, a_t^{act})$
2. **Expect** that the agent might explore in the next state

$$Q^{\pi^\alpha}(s_t, a_t^{act}) += \eta \left( r_t + \gamma \mathbb{E}^{(1-\alpha)\pi + \alpha\pi_0} Q^{\pi^\alpha}(s_{t+1}, a) - Q^{\pi^\alpha}(s_t, a_t^{act}) \right)$$



## Two Approaches – Expected approach

1. Update  $Q^{\pi^\alpha}(s_t, a_t^{act})$
2. **Expect** that the agent might explore in the next state

$$Q^{\pi^\alpha}(s_t, a_t^{act}) += \eta \left( r_t + \gamma \mathbb{E}^{(1-\alpha)\pi + \alpha\pi_0} Q^{\pi^\alpha}(s_{t+1}, a) - Q^{\pi^\alpha}(s_t, a_t^{act}) \right)$$

- Calculating expectations can be hard.
  - **Requires sampling in the continuous case!**





## Two Approaches – Surrogate approach

- Exploration is incorporated into the environment!
1. Update  $Q^{\pi^\alpha}(s_t, \mathbf{a}_t^{\text{greedy}})$
  2. The rewards and next state  $r_t, s_{t+1}$  are given by the acted action  $\mathbf{a}_t^{\text{act}}$

$$Q^{\pi^\alpha}(s_t, \mathbf{a}_t^{\text{greedy}}) += \eta \left( r_t + \gamma Q^{\pi^\alpha}(s_{t+1}, \mathbf{a}_{t+1}^{\text{greedy}}) - Q^{\pi^\alpha}(s_t, \mathbf{a}_t^{\text{greedy}}) \right)$$



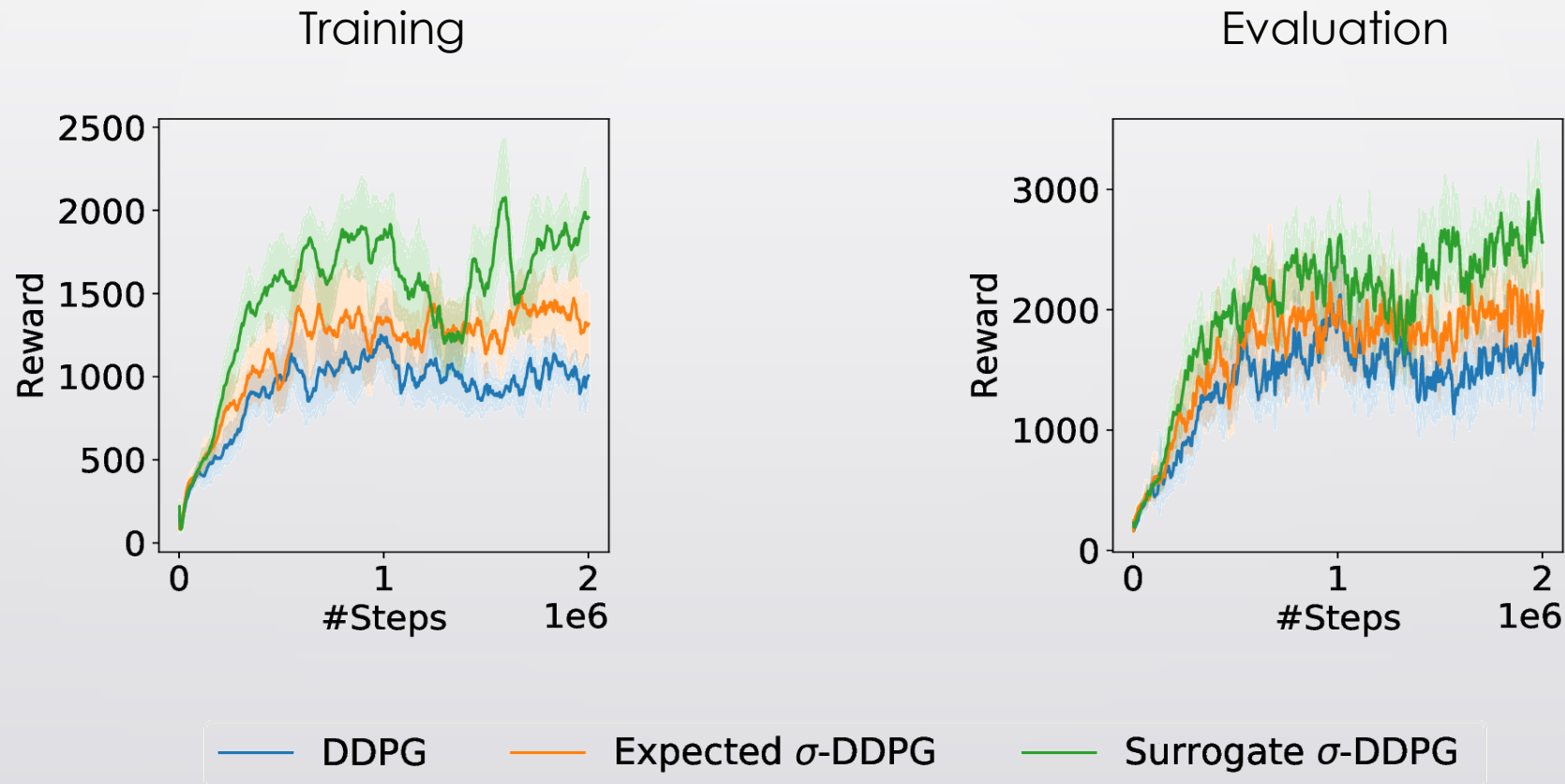
## Two Approaches – Surrogate approach

- Exploration is incorporated into the environment!
1. Update  $Q^{\pi^\alpha}(s_t, \mathbf{a}_t^{\text{greedy}})$
  2. The rewards and next state  $r_t, s_{t+1}$  are given by the acted action  $\mathbf{a}_t^{\text{act}}$

$$Q^{\pi^\alpha}(s_t, \mathbf{a}_t^{\text{greedy}}) += \eta \left( r_t + \gamma Q^{\pi^\alpha}(s_{t+1}, \mathbf{a}_{t+1}^{\text{greedy}}) - Q^{\pi^\alpha}(s_t, \mathbf{a}_t^{\text{greedy}}) \right)$$

- **NO NEED TO SAMPLE!**

# Deep RL Experimental Results





# Summary

- We define Exploration Conscious RL and analyze its properties.
- Exploration Conscious RL can improve performance over **both the training and evaluation regimes.**
- **Conclusion: Exploration-Conscious RL and specifically, the Surrogate approach, can easily help to improve variety of RL algorithms.**

---

# Summary

- We define Exploration Conscious RL and analyze its properties.
- Exploration Conscious RL can improve performance over **both the training and evaluation regimes**.
- **Conclusion:** Exploration-Conscious RL and specially, the Surrogate approach, can easily help to improve variety of algorithms.

**SEE YOU AT POSTER**

**#90**

