

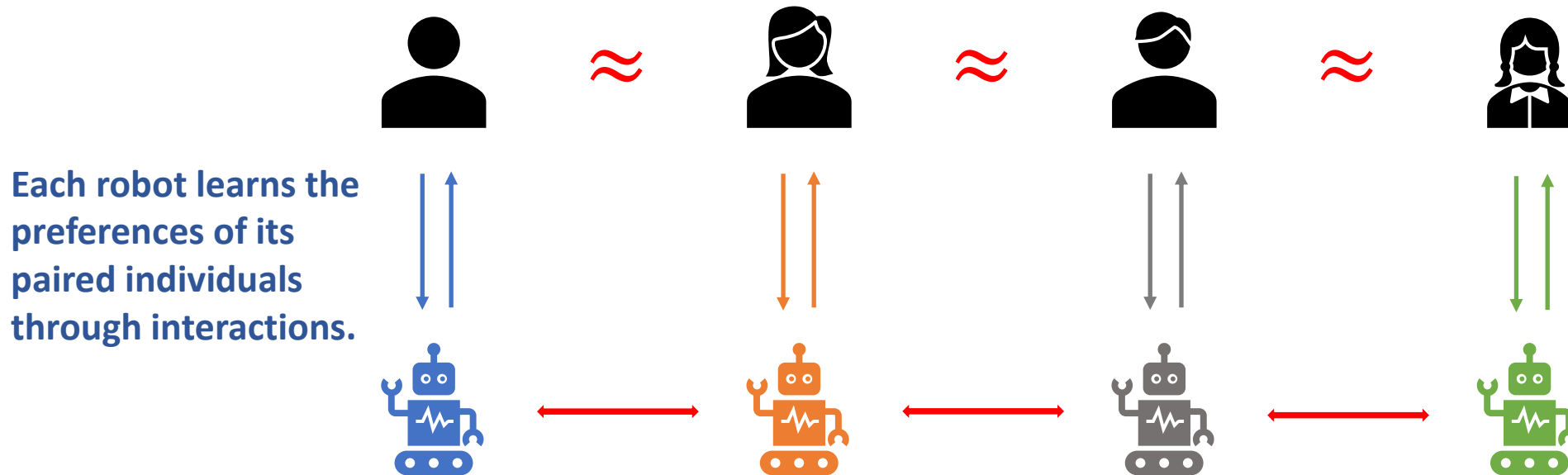
Provably Efficient Multi-Task Reinforcement Learning with Model Transfer

Chicheng Zhang and Zhi Wang

NeurIPS 2021



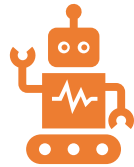
Heterogenous Multi-Task Online Reinforcement Learning (RL)



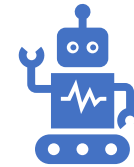
- A group of assistive robots deployed to provide personalized healthcare services (Kubota et al., 2020).
- Question: If the robots receive **similar yet nonidentical** feedback, how can they learn to perform their respective tasks faster in an **online** RL setting?

Multi-Player Episodic RL (MPERL)

- A set of M players (robots) concurrently interact with their respective environments, each represented as an Episodic MDP.



Alice



Bob



Action 1

$$R_p(s, 1) = 0.4$$



Action 2

$$R_p(s, 2) = 0.5$$



Action 3

$$R_p(s, 3) = 0.6$$



Action 1

$$R_q(s, 1) = 0.6$$



Action 2

$$R_q(s, 2) = 0.4$$



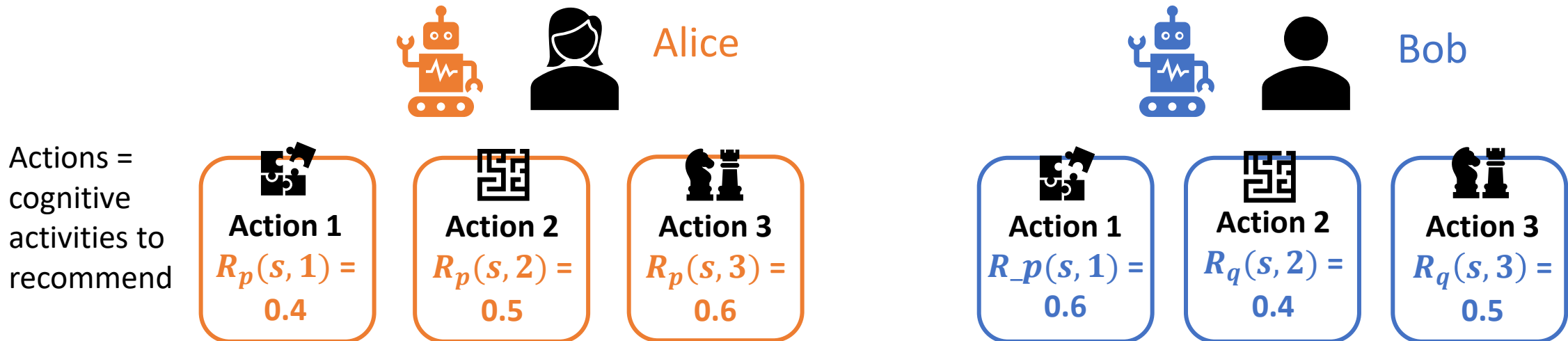
Action 3

$$R_q(s, 3) = 0.5$$

Actions =
cognitive
activities to
recommend

The ϵ -MPERL Problem

- A set of M players (robots) concurrently interact with their respective environments, each represented as an Episodic MDP.



$\forall p, q, s, a:$

$$\begin{aligned} |R_p(s, a) - R_q(s, a)| &\leq \epsilon \longrightarrow \epsilon: \text{dissimilarity parameter} \\ \|\mathbb{P}_p(\cdot | s, a) - \mathbb{P}_q(\cdot | s, a)\|_1 &\leq \epsilon/H \end{aligned}$$

The ϵ -MPERL Problem: formal setup

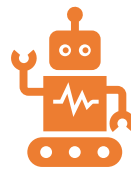
- M episodic, tabular, H -layered MDPs $(\mathcal{M}_p)_{p=1}^M$ with shared state-action spaces, and common initial distribution $\delta(s_0)$

- For episodes $k = 1, 2, \dots, K$:

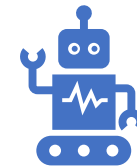
- For players $p = 1, 2, \dots, M$:

- Player p interacts with \mathcal{M}_p with policy $\pi^k(p)$ for one episode, obtaining trajectory τ_p^k

- All M trajectories $(\tau_p^k)_{p=1}^M$ are shared among the players



Alice



Bob

- Collective regret: $\text{Reg}(K) = \sum_{p=1}^M \sum_{k=1}^K V_p^*(s_0) - V_p^{\pi^k(p)}(s_0)$

Optimal value of player p Value of player p executing $\pi^k(p)$

Baseline: individual single-task learning

- Each player learns separately using a state-of-the-art online tabular RL algorithm, e.g., *Strong-Euler* (Simchowitz and Jamieson, 2019), achieving a collective regret of

- (Gap-independent bound) $\tilde{O}(M\sqrt{H^2SAK})$

- (Gap-dependent bound)

$$\tilde{O}\left(\sum_{p=1}^M\left(\sum_{(s,a)\in Z_{p,\text{opt}}}\frac{H^3\ln K}{\Delta_{p,\text{min}}}+\sum_{(s,a)\notin Z_{p,\text{opt}}}\frac{H^3\ln K}{\Delta_p(s,a)}\right)\right)$$

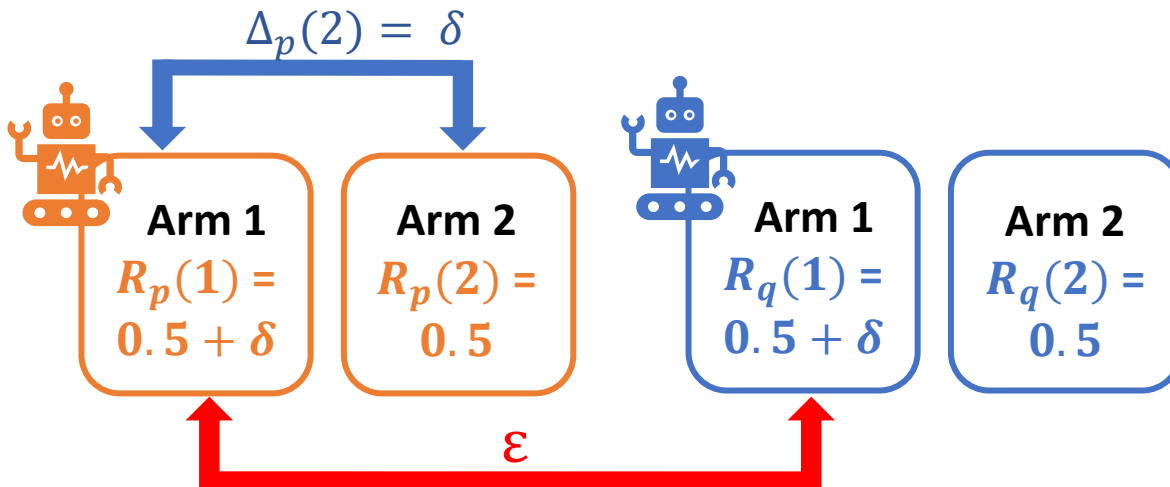
where $\Delta_p(s, a) := V_p^*(s) - Q_p^*(s, a)$, $Z_{p,\text{opt}} = \{(s, a) : \Delta_p(s, a) = 0\}$,

$$\Delta_{p,\text{min}} = \min_{(s,a)\notin Z_{p,\text{opt}}}\Delta_p(s, a)$$

- Can we do better with inter-task information sharing?

The benefit of multi-task learning

- (Wang, Zhang, Singh, Riek, Chaudhuri, 2021): in a multi-task multi-armed bandit setting, information sharing sometimes does not help, *information theoretically*.
- Example: For a fixed ε and $\delta < \varepsilon/4$, consider:



Claim: Any sublinear regret algorithm must have $\Omega\left(\frac{M \ln K}{\delta}\right)$ regret, no better than the individual single-task learning baseline.

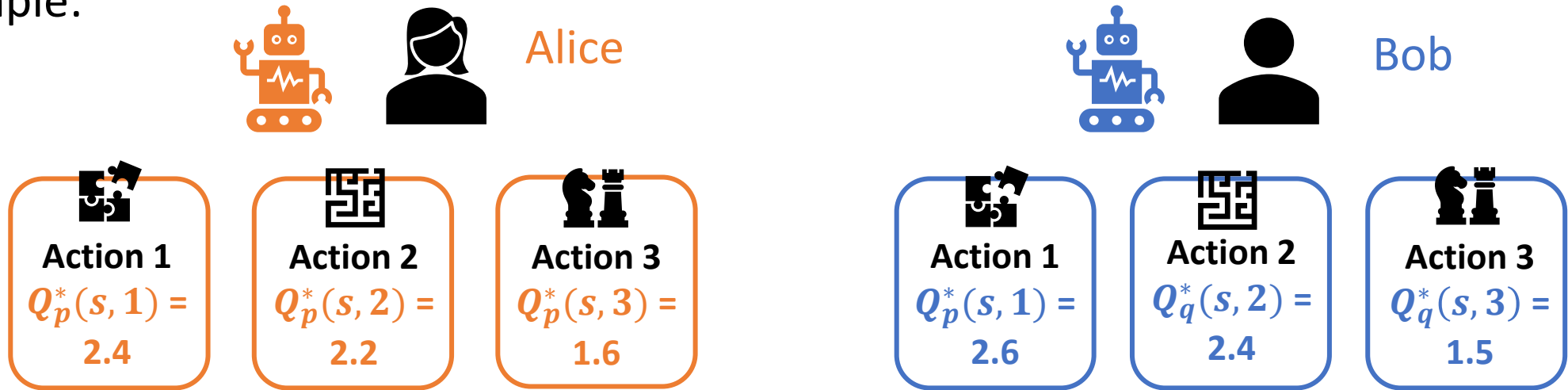
- Key observation: the benefit of multi-task learning depends on the interaction between ε and suboptimality gaps $\Delta_p(s, a)$

Key notion: subpar state-action pairs

- Subpar state-action pairs:

$$\mathcal{J}_\epsilon = \{(s, a) : \text{for some } p \in [M], \Delta_p(s, a) \geq \Omega(H\epsilon)\}$$

- Example:



- $(s, 3) \in \mathcal{J}_\epsilon; (s, 2) \notin \mathcal{J}_\epsilon$
- Subpar state-action pairs are those amenable for inter-task information sharing

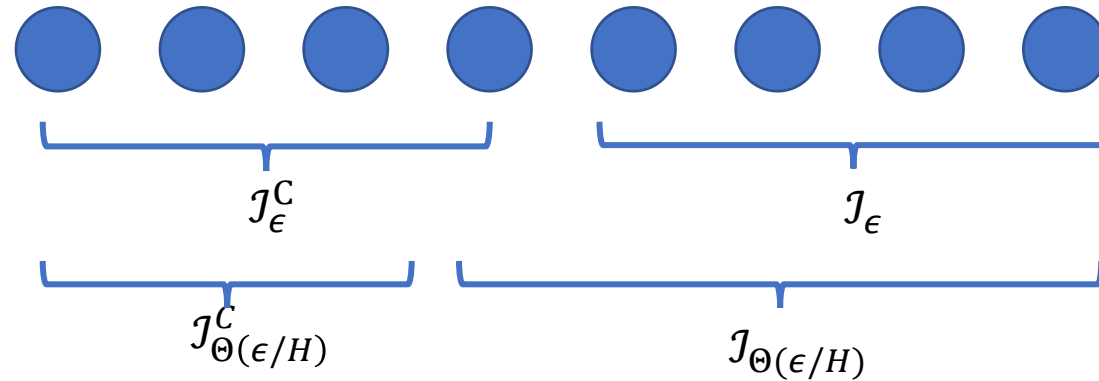
Our results

For ε -MPERL problems, assuming known ε :

- Our algorithm, Multi-Task-Euler(ε), achieves gap-dependent and gap-independent regret upper bounds
- We also show gap-dependent and gap-independent regret lower bounds, that nearly match the upper bounds for constant H

Our results: gap-independent bounds

State-action pairs



Individual
Strong-Euler

$$M \sqrt{H^2 |\mathcal{J}_\epsilon^C| K} \quad + \quad M \sqrt{H^2 |\mathcal{J}_\epsilon| K}$$

Multi-task-Euler(ϵ)

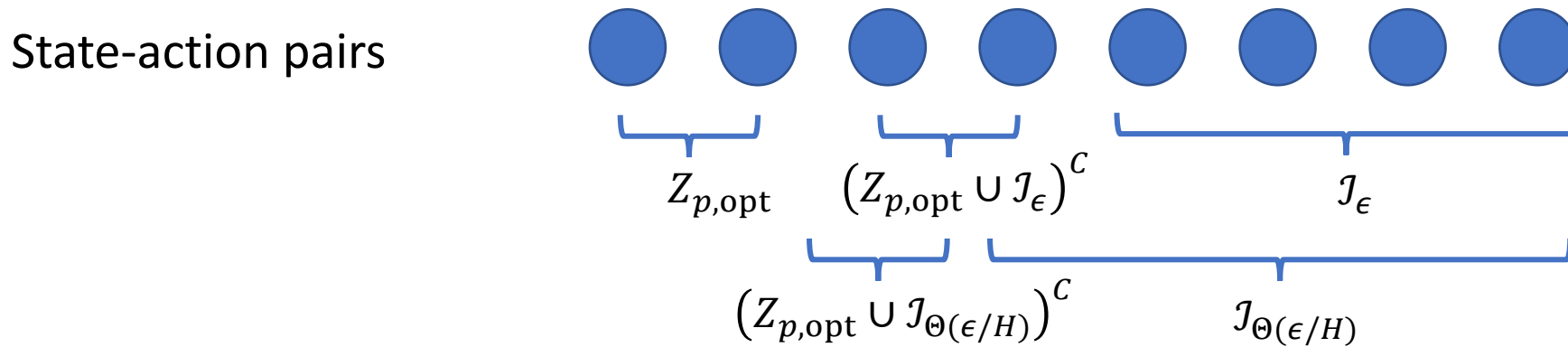
$$M \sqrt{H^2 |\mathcal{J}_\epsilon^C| K} \quad + \quad \sqrt{MH^2 |\mathcal{J}_\epsilon| K}$$

Lower bound

$$M \sqrt{H^2 |\mathcal{J}_{\Theta(\epsilon/H)}^C| K} \quad + \quad \sqrt{MH^2 |\mathcal{J}_{\Theta(\epsilon/H)}| K}$$

Our results: gap-dependent bounds

For player p 's contribution to the collective regret:



Individual Strong-Euler	$\sum_{s,a}$	$\frac{H^3 \ln K}{\Delta_{p,\min}}$	$\frac{H^3 \ln K}{\Delta_p(s, a)}$	$\frac{H^3 \ln K}{\Delta_p(s, a)}$
Multi-task-Euler(ϵ)	$\sum_{s,a}$	$\frac{H^3 \ln K}{\Delta_{p,\min}}$	$\frac{H^3 \ln K}{\Delta_p(s, a)}$	$\frac{1}{M} \cdot \frac{H^3 \ln K}{\Delta_p(s, a)}$
Lower bound	$\sum_{s,a}$		$\frac{H^2 \ln K}{\Delta_p(s, a)}$	$\frac{1}{M} \cdot \frac{H^2 \ln K}{\Delta_p(s, a)}$

Multi-task-Euler(ϵ): main ideas

- For each player p , Multi-Task-Euler(ϵ):
 1. Maintains two model estimates for \mathcal{M}_p : (1) an individual estimate $\widehat{\mathcal{M}}_p$ (2) an aggregate model estimate $\widehat{\mathcal{M}}$
 2. Performs a “heterogeneous” optimistic value iteration using both $\widehat{\mathcal{M}}_p$ and $\widehat{\mathcal{M}}$ to obtain \widehat{Q}_p , a tight upper confidence bound of Q_p^* , and executes its greedy policy
- Similar algorithmic idea of “model transfer” has appeared in prior works, e.g., (Taylor, Jong, & Stone, 2008), (Pazis & Parr, 2016)

Technical overview

- Upper bounds: a new surplus bound in the multi-task setting:

$$\widehat{Q}_p(s, a) - (R_p(s, a) + \langle \mathbb{P}_p(\cdot | s, a), \widehat{V}_p \rangle) \leq \tilde{O} \left(\min \left(\sqrt{\frac{1}{n_p(s, a)}}, \epsilon + \sqrt{\frac{1}{n(s, a)}} \right) \right),$$

and combine with the “clipping trick” (Simchowitz & Jamieson, 2019)

- Lower bounds: combine the multi-task bandit lower bounds (Wang, Zhang, Singh, Riek, Chaudhuri, 2021) with a standard bandit-to-RL conversion

Conclusion and open problems

- We study ε -MPERL, a new multi-task RL setting; this complements existing multi-task RL settings (e.g., Brunskill & Li, 2013, Liu, Guo, & Brunskill, 2016, Pazis & Parr, 2016)
- We give upper and lower bounds on the collective regret that are nearly matching for constant episode length H
- Open questions:
 - Improve the dependence on H in the collective regret bounds
 - Improve the dependence on $Z_{p,\text{opt}}$, similar to recent works (e.g., Xu, Ma, & Du, 2021)
 - Extensions to RL with function approximation

Thank you!