

# Private Reinforcement Learning with PAC and Regret Guarantees

Giuseppe Vietri<sup>1</sup>, Borja Balle<sup>2</sup>, Akshay Krishnamurthy<sup>3</sup>, Zhiwei Steven Wu<sup>4</sup>

<sup>1</sup>University of Minnesota, <sup>2</sup>Deepmind, <sup>3</sup>Microsoft Research, <sup>4</sup>Carnegie Mellon University



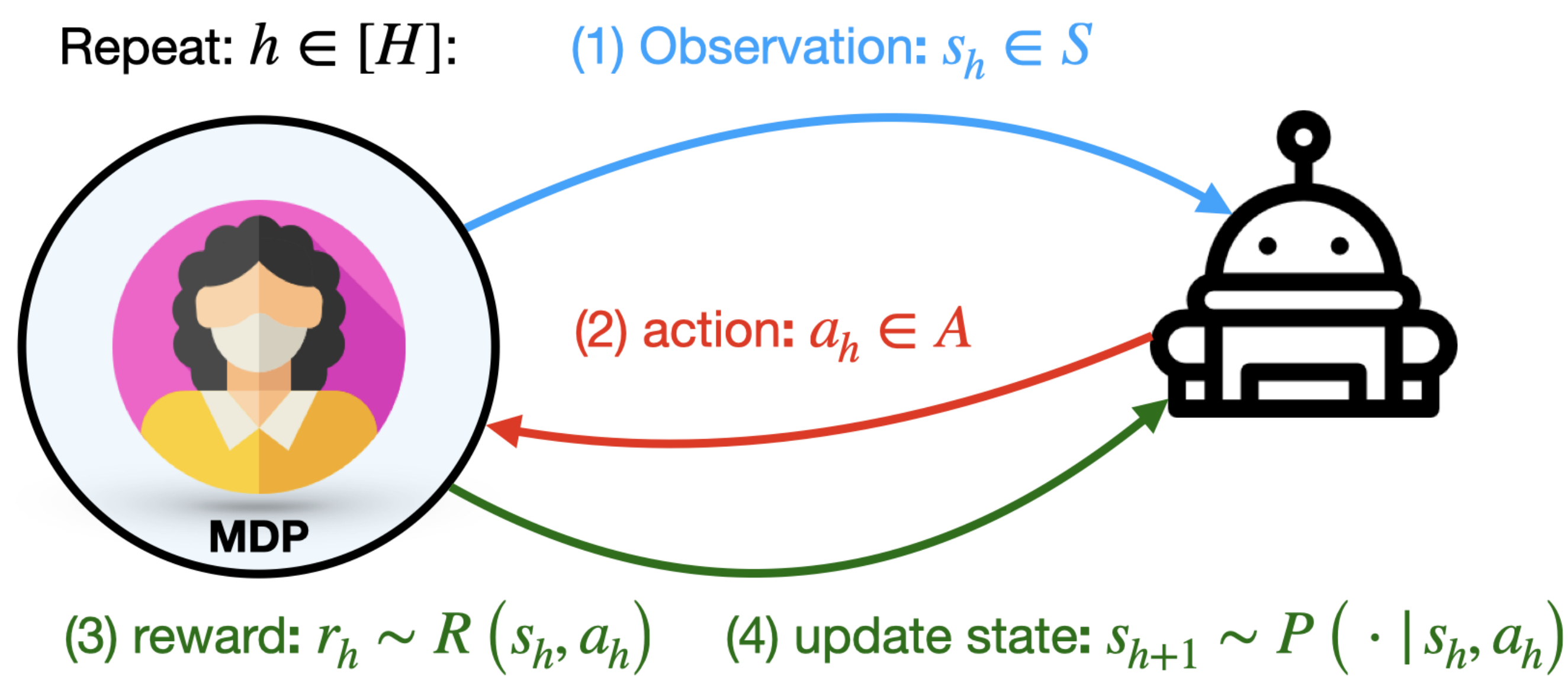
Microsoft  
Research



Google  
DeepMind

## Online Learning MDP

1. An agent interacts with a fix horizon MDP over a sequence of episodes:
2.  $M = (S, A, R, P, p_0, H)$  // Environment
3. **Privacy in RL**: we assume that sequence of states and rewards is sensitive data belonging to users.



## Contributions

1. Definition of Joint Differential Privacy (**JDP**) in Reinforcement Learning (**RL**).
2. Almost optimal  $\varepsilon$ -**JDP RL** algorithm with bounded *sample complexity* and bounded *regret*.
3. Lower bound for any **RL** algorithm that satisfies **JDP**.

## Results

1. Let  $\alpha$  be a target accuracy and  $\varepsilon$  the privacy parameter.
  - *How many episodes does it take to learn an  $\alpha$ -optimal policy?*
  - A smaller  $\varepsilon$  means the algorithm is more private.
2. Upper bound sample complexity:

$$\tilde{O}\left(\frac{SAH^4}{\alpha^2} + \frac{S^2AH^4}{\varepsilon\alpha}\right) \quad (1)$$

3. Lower bound sample complexity:

$$\tilde{\Omega}\left(\frac{SAH^2}{\alpha^2} + \frac{SAH}{\varepsilon\alpha}\right) \quad (2)$$

## Privacy Definition

- We represent a user as a tree of depth  $H$  encoding all possible state and reward paths.
- Let  $\mathcal{M}$  be a **RL** algorithm:
  - The INPUT is a sequence of  $T$  users.
  - The OUTPUT is a sequence of  $T$  actions and a final policy  $\pi_T$ .

$$\mathcal{M}(\text{user}_1, \dots, \text{user}_T) = (a^{(1)}, \dots, a^{(T)}, \pi_T) \quad a^{(t)} = (a_1^{(t)}, \dots, a_h^{(t)})$$

- $U$  and  $U'$  are  $t$ -neighboring if they differ only on user at position  $t$ :

$$U = (\text{user}_1, \dots, \text{user}_t, \dots, \text{user}_T), \quad U' = (\text{user}_1, \dots, \text{user}'_t, \dots, \text{user}_T)$$

**Definition** (KPRU14). A mechanism  $\mathcal{M}$  is  $\varepsilon$ -**JDP** if for all  $t$ , all  $t$ -neighboring user sequences  $U, U'$  and all events  $E \subseteq A^{H \times [T-1]} \times \Pi$  we have

$$\Pr[\mathcal{M}_{-t}(U) \in E] \leq e^\varepsilon \Pr[\mathcal{M}_{-t}(U') \in E] \quad (3)$$

## PUCB

- **PUCB** is a **JDP** version of UBEV [DLB 2017].
- Non-private event counters:  $\hat{n}_t(s, a, h), \hat{m}_t(s, a, h, s'), \hat{r}_t(s, a, h)$ .
- Private event counters:  $\tilde{n}_t(s, a, h), \tilde{m}_t(s, a, h, s'), \tilde{r}_t(s, a, h)$ .
- Use Binary mechanism (**BM**) from [Dwork et al., 2010] and [Chan et al., 2011]. For any  $t \in [T]$ :

$$|\hat{n}_t(s, a, h) - \tilde{n}_t(s, a, h)| \leq \frac{H}{\varepsilon} \log(T)^{5/2} \log(SAH/\beta) := E_\varepsilon \quad (4)$$

- Balance exploration/exploitation with optimism:

– Compute  $\tilde{Q}_t^+$  with **DP**:

$$\tilde{Q}_t^+(s, a, h) = (\text{Reward}) + (\text{Future reward}) + (\text{Bonus}) \quad (5)$$

– Greedy policy:  $\pi_t(s, h) = \arg \max_{a^*} \tilde{Q}_t^+(s, a^*, h)$ .

## Private Optimistic Q-function

$$\tilde{Q}^+(s, a, h) = \frac{\tilde{r}(s, a, h) + \sum_{s'} \tilde{V}_{h+1}(s') \tilde{m}(s, a, h, s')}{\tilde{n}(s, a, h)} + \widetilde{\text{conf}}(s, a, h)$$

Immediate  
reward

Expected  
future reward

Bonus term

$$\widetilde{\text{conf}}(s, a, h) = \tilde{\phi}(s, a, h) + \tilde{\psi}(s, a, h)$$

$$\tilde{\phi}(s, a, h) = (1 + H) \sqrt{\frac{T/\beta}{\tilde{n}(s, a, h)}}$$

sampling error  
[Dann et al. 2017]

$$\tilde{\psi}(s, a, h) = (1 + SH) \left( \frac{3E_\varepsilon}{\tilde{n}(s, a, h)} + \frac{2E_\varepsilon^2}{\tilde{n}(s, a, h)^2} \right)$$

privacy error  
[This work].

## PUCB Sample Complexity Analysis

- Construct  $\tilde{Q}_t^+$  with DP counters and show that:

$$\tilde{Q}_t^+(s, a, h) \geq Q^*(s, a, h) \quad (6)$$

- Optimality gap decomposition:

$$V^* - V^{\pi_t} := \Delta_t \leq \sum_{(s, a, h) \in (S, A, [H])} w_t(s, a, h) \widetilde{\text{conf}}_t(s, a, h) \quad (7)$$

- Bounding number of episodes where  $\Delta_t > \alpha$ .

## Lower Bound Analysis

- Consider hard-MDP construction.
- Lower bound for **DP** best-arm-identification:
- We show that finding the  $\alpha$ -optimal arm takes  $\tilde{\Omega}\left(\frac{A}{\varepsilon\alpha} \ln \frac{1}{4\beta}\right)$  tries.
- We consider a Public Initial State (**PIS**) Setting:
  - And do a reduction to **DP** best-arm-identification.
  - We show that the learner makes at least  $\frac{SAH}{24\varepsilon\alpha} \ln \frac{1}{4\beta}$  mistakes in the PIS setting.
  - If algorithm satisfies  $\varepsilon$ -**JDP**  $\implies$  satisfies  $\varepsilon$ -**JDP** in the PIS setting.

## Acknowledgements

Giuseppe Vietri has been supported by the GAANN fellowship from the U.S. Department of Education. We want to thank Matthew Joseph, whose comments improved our definition of joint-differential-privacy.