

Private Reinforcement Learning with PAC and Regret Guarantees

Giuseppe Vietri

University of Minnesota

Borja Balle

Deepmind

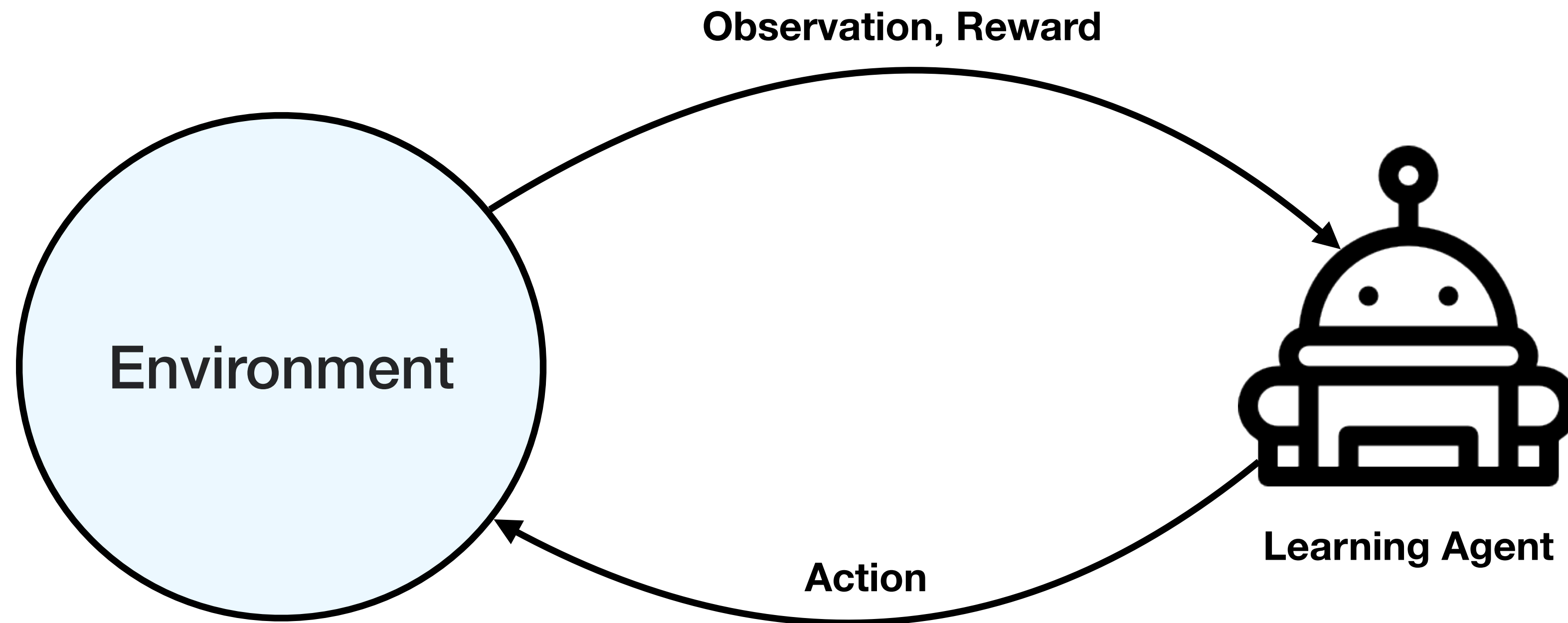
Akshay Krishnamurthy

Microsoft Research

Steven Wu

University of Minnesota

Reinforcement Learning



RL in healthcare

- Agent → Provider
- Environment → Patients
- Observations → Symptoms
- Actions → Treatments
- Reward → patient improves
- Privacy of patients?



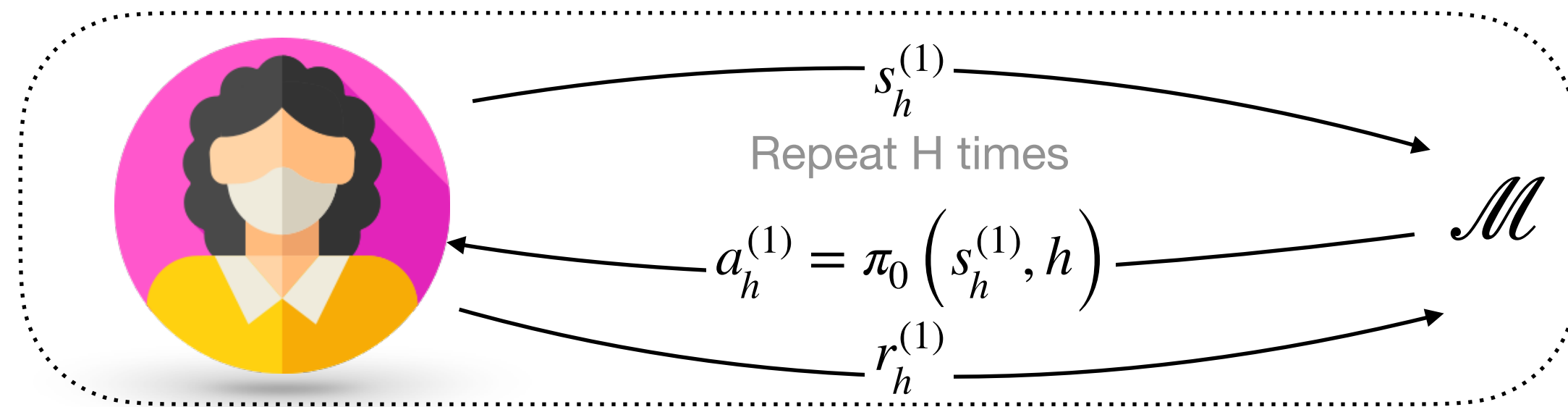
Episodic RL Protocol

Input: Learning Agent \mathcal{M} , User sequence $U = (\text{User 1}, \text{User 2}, \dots)$

Episode

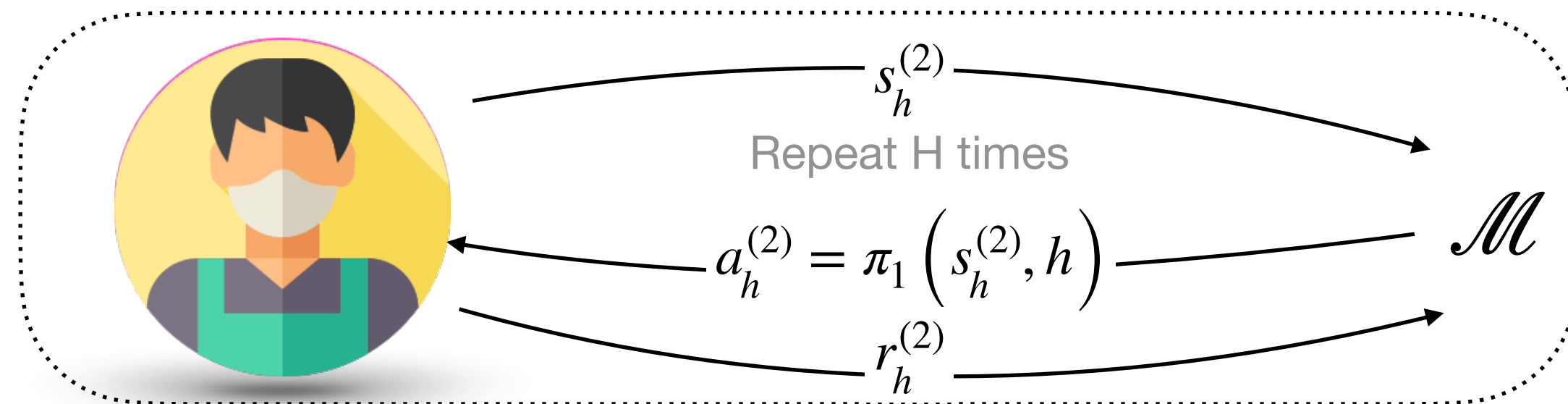
Initialize: π_0

1



Update: $\pi_1 \leftarrow \mathcal{M}$

2



Update: $\pi_2 \leftarrow \mathcal{M}$

⋮

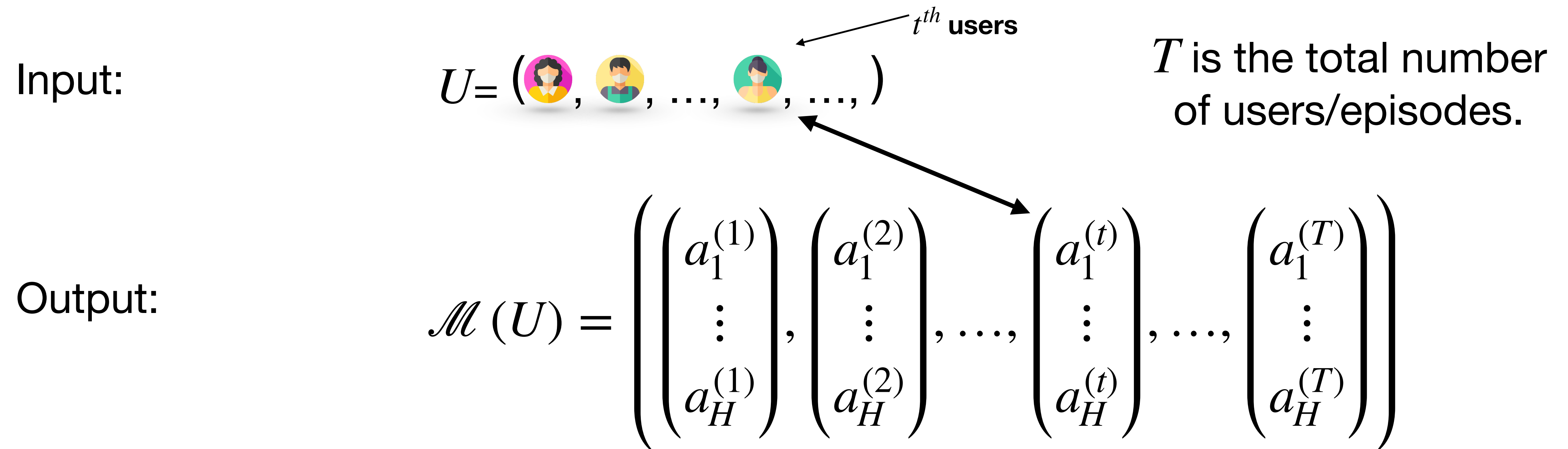
S = Number of States.

A = Number of Actions.

H = Time-steps per episode.

Main Results: A Privacy Formulation for RL

Let \mathcal{M} be a RL algorithm.



\mathcal{M} must satisfy ϵ -Joint Differential Privacy Under Continual Observation.

Main Results

Private

	PAC	Regret
Upper Bounds	$\tilde{O}\left(\left(\frac{SAH^4}{\alpha^2} + \frac{S^2AH^4}{\epsilon\alpha}\right)\right)$	$\tilde{O}\left(H^2\sqrt{SAT} + \frac{S^2AH^4}{\epsilon}\right)$
Lower Bound	$\tilde{\Omega}\left(\left(\frac{SAH^2}{\alpha^2} + \frac{SAH}{\epsilon\alpha}\right)\right)$	$\tilde{\Omega}\left(\sqrt{HSAT} + \frac{SAH}{\epsilon}\right)$

Prior work.
non-Private

	PAC	Regret
Upper Bounds	$\tilde{O}\left(\frac{SAH^4}{\alpha^2}\right)$ [Dann et al. 2017]	$\tilde{O}\left(\sqrt{SAHT}\right)$ [Azar et al. 2017]
Lower Bounds	$\tilde{\Omega}\left(\frac{SAH^3}{\alpha^2}\right)$ [Dann et al. 2017]	$\tilde{\Omega}\left(\sqrt{SAHT}\right)$ [Jaksch et al. 2010]

PUCB: A private RL Algorithm.

- **PUCB** is a private version of the **U**pper-**B**ound the **E**xpected next **V**alue algorithm (**UBEV**) [Dann et al., 2017].
- Compute an optimistic Q-value function using standard batch Q-learning updates with an optimism bonus.
- Keeps track of rewards and dynamics estimates.
- Follows a greedy policy: $\pi_t(s, h) = \arg \max_{a^*} Q(s, a^*, h)$

Q Learning

Without privacy:

For $t = 1 \dots, T$:

$$Q \leftarrow \text{OptimisticPlanning}(\hat{n}, \hat{r}, \hat{m})$$

$$s_1^{(t)} \sim \langle \text{uniform distribution over states} \rangle$$

For $h = 1 \dots, H$:

$$a_h^{(t)} = \arg \max_{a^*} Q(s_h^{(t)}, a^*, h)$$

$$r_h^{(t)} \sim R(s_h^{(t)}, a_h^{(t)}, h), s_{h+1}^{(t)} \sim P(\cdot | s_h^{(t)}, a_h^{(t)}, h)$$

Increment counters: $\hat{n}(s_h^{(t)}, a_h^{(t)}, h)$,

$$\hat{r}(s_h^{(t)}, a_h^{(t)}, h), \hat{m}(s_h^{(t)}, a_h^{(t)}, h, s_{h+1}^{(t)})$$

With Privacy

For $t = 1 \dots, T$:

$$Q \leftarrow \text{PrivateOptimisticPlanning}(\tilde{n}, \tilde{r}, \tilde{m})$$

$$s_1^{(t)} \sim \langle \text{uniform distribution over states} \rangle$$

For $h = 1 \dots, H$:

$$a_h^{(t)} = \arg \max_{a^*} Q(s_h^{(t)}, a^*, h)$$

$$r_h^{(t)} \sim R(s_h^{(t)}, a_h^{(t)}, h), s_{h+1}^{(t)} \sim P(\cdot | s_h^{(t)}, a_h^{(t)}, h)$$

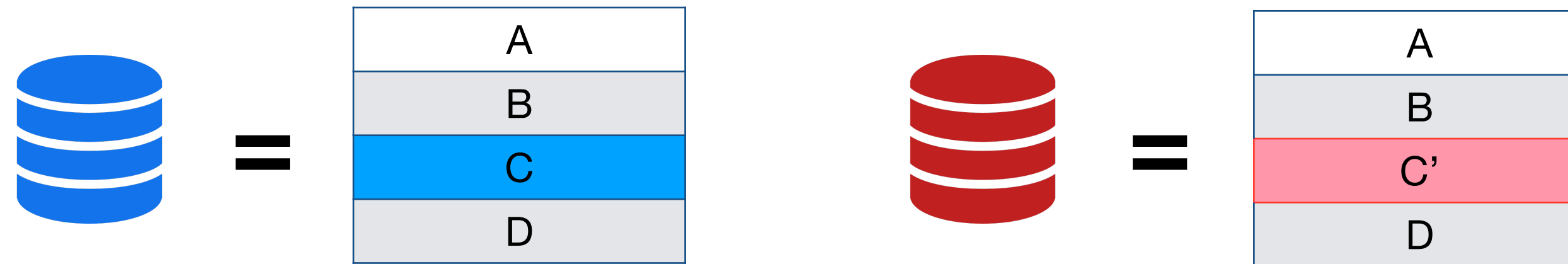
Increment private counters: $\tilde{n}(s_h^{(t)}, a_h^{(t)}, h)$,

$$\tilde{r}(s_h^{(t)}, a_h^{(t)}, h), \tilde{m}(s_h^{(t)}, a_h^{(t)}, h, s_{h+1}^{(t)})$$

Standard Differential Privacy (DP)

[Dwork et al., 2006]

Two datasets



are **neighbors** if they are different on only one row.

Definition: Mechanism M satisfies ϵ -differential privacy if, for all neighboring datasets and for all $r \in \text{range}(M)$

$$\Pr[M(\text{red database}) = r] \leq e^\epsilon \Pr[M(\text{blue database}) = r]$$

Why Is DP not Applicable?

- If the algorithm must satisfy **Differential Privacy** then for any two states s, s' , it holds that

$$\Pr [\pi_t (s, h) = a] \sim \Pr [\pi_t (s', h) = a]$$

Joint Differential Privacy Definition

$$U = (\text{User 1}, \text{User 2}, \dots, \text{User } t, \dots)$$

t-th users are different

$$U' = (\text{User 1}, \text{User 2}, \dots, \text{User } t', \dots)$$

U and U' are t -neighboring user sequences

Definition: A mechanism \mathcal{M} is ε -jointly differentially private if for all t , all t -neighboring user sequences U, U' and all future events $E \subseteq \mathcal{A}^{H \times [T-1]}$ we have

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

Event Counters

Total counters: $2SAH + S^2AH$

Use **Binary mechanism** from [Dwork et al., 2010] and [Chan et al., 2011].

If $\tilde{n}(s, a, h) \leftarrow$ Binary Mechanism With Privacy Parameter $\frac{\epsilon}{H}$

The **composition** of all counters satisfies ϵ -DP [Hsu et al., 2014].

$$\left| \tilde{n}(s, a, h) - \hat{n}(s, a, h) \right| \leq \frac{H}{\epsilon} \log(T)^{5/2} \log(2/\beta) := E_\epsilon$$

Balancing Exploration/Exploitation with Optimism.

Without privacy:

$$\widehat{Q}^+(s, a, h) = \widehat{Q}(s, a, h) + \widehat{\phi}(s, a, h)$$

Confidence due to
sampling error
[Dann et al. 2017]

Confidence term
due to privacy.
[This work]

With privacy:

$$\widetilde{Q}^+(s, a, h) = \widetilde{Q}(s, a, h) + \widetilde{\phi}(s, a, h) + \widetilde{\psi}(s, a, h)$$

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

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

JDP Proof: Billboard Lemma

Theorem: Algorithm PUCB satisfies ϵ -joint differential privacy.

- We use the billboard lemma from [Hsu et al., 2016]
- The billboard lemma: An algorithm is JDP if the output sent to each user is a function of the user's private data and a common signal computed using standard differential privacy.

- For example:

• The output for user t is: $a_h^{(t)} = \arg \max_{a^*} Q(s_h^{(t)}, a^*, h)$

$s_h^{(t)}$ is part of user t private data

The Q-function was computed using ϵ -differential privacy

PAC Upper Bound Proof: Optimism

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

$$\widehat{Q}^+(s, a, h) \leq \frac{\widetilde{r}(s, a, h) + E_\varepsilon + \sum_{s'} \widetilde{V}_{h+1}(s') (\widetilde{m}(s, a, h, s') + E_\varepsilon)}{\widetilde{n}(s, a, h) - E_\varepsilon} + \widehat{\phi}(s, a, h)$$

Case 1: If $\widetilde{n}(s, a, h) \geq 2E_\varepsilon$ the following holds: $\frac{1}{\widetilde{n}(s, a, h) - E_\varepsilon} \leq \left(\frac{1}{\widetilde{n}(s, a, h)} + \frac{2E_\varepsilon}{\widetilde{n}(s, a, h)^2} \right)$

$$\widehat{Q}^+(s, a, h) \leq \widetilde{Q}(s, a, h) + \left(\frac{1}{\widetilde{n}(s, a, h)} + \frac{2E_\varepsilon}{\widetilde{n}(s, a, h)^2} \right) (1 + SH)E_\varepsilon + \widehat{\phi}(s, a, h) = \widetilde{Q}^+(s, a, h)$$

← $\widetilde{\psi}(s, a, h)$

Case 2: If $\widetilde{n}(s, a, h) < 2E_\varepsilon$ then we make $\widetilde{Q}^+(s, a, h) = H$

PAC Lower Bound Proof

1. Lower bound for private-best-arm-identification problem.

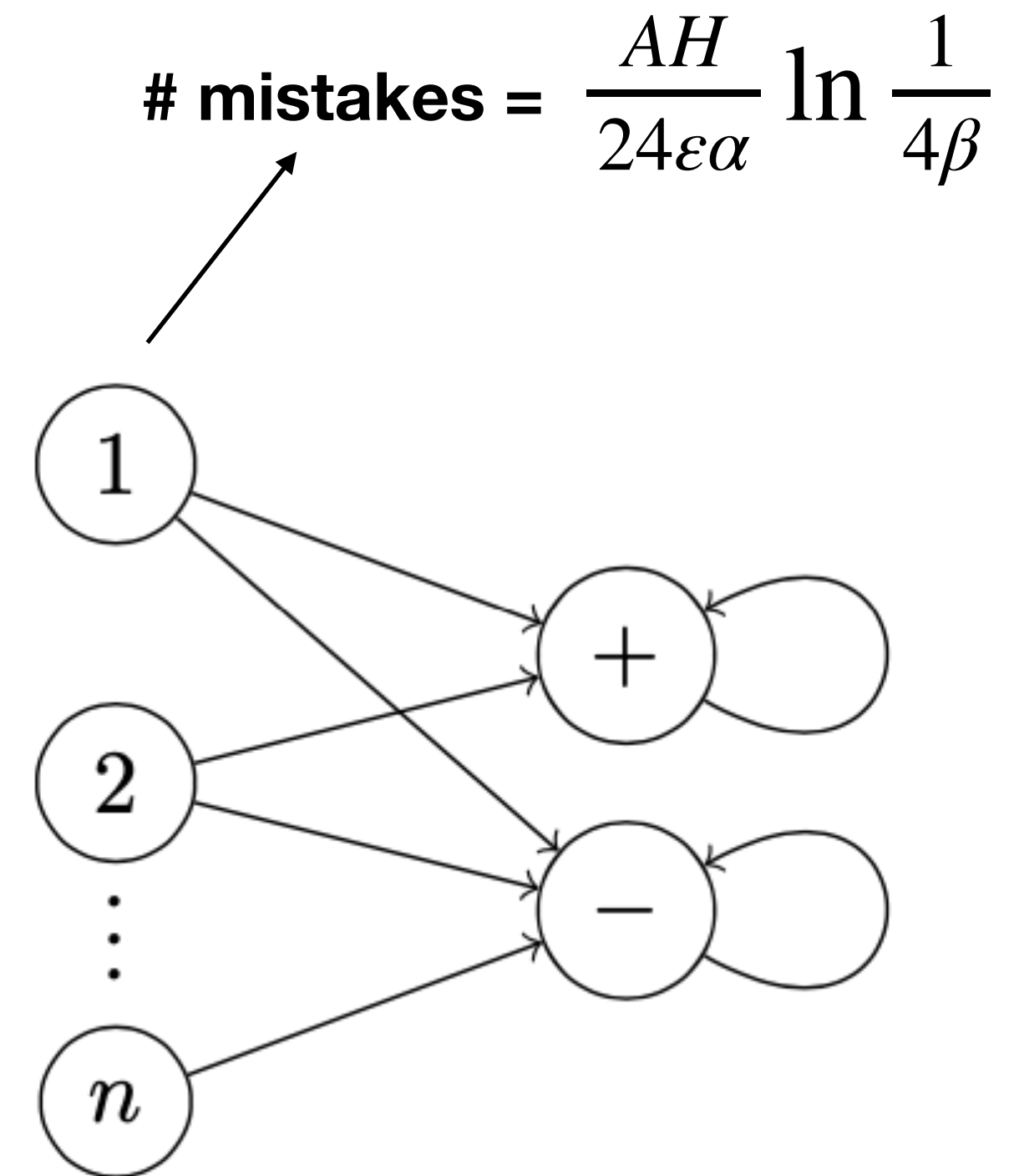
1. $\widetilde{\Omega} \left(\frac{A}{\varepsilon\alpha} \ln \frac{1}{4\beta} \right)$

2. We consider a simpler: *Public Initial State Setting*.

1. Each initial states $\{1, \dots, n\}$ is a private best-arm-identification MAB problem.

3. Therefore learner must make a total of at least $\frac{SAH}{24\varepsilon\alpha} \ln \frac{1}{4\beta}$ mistakes.

4. ε -JDP \implies ε -JDP in the public initial state setting.



Conclusion

- Introduced a meaningful formulation of privacy for RL.
- A private optimism based algorithm with PAC and regret Guarantees.
- First analysis of lower bounds for private RL.