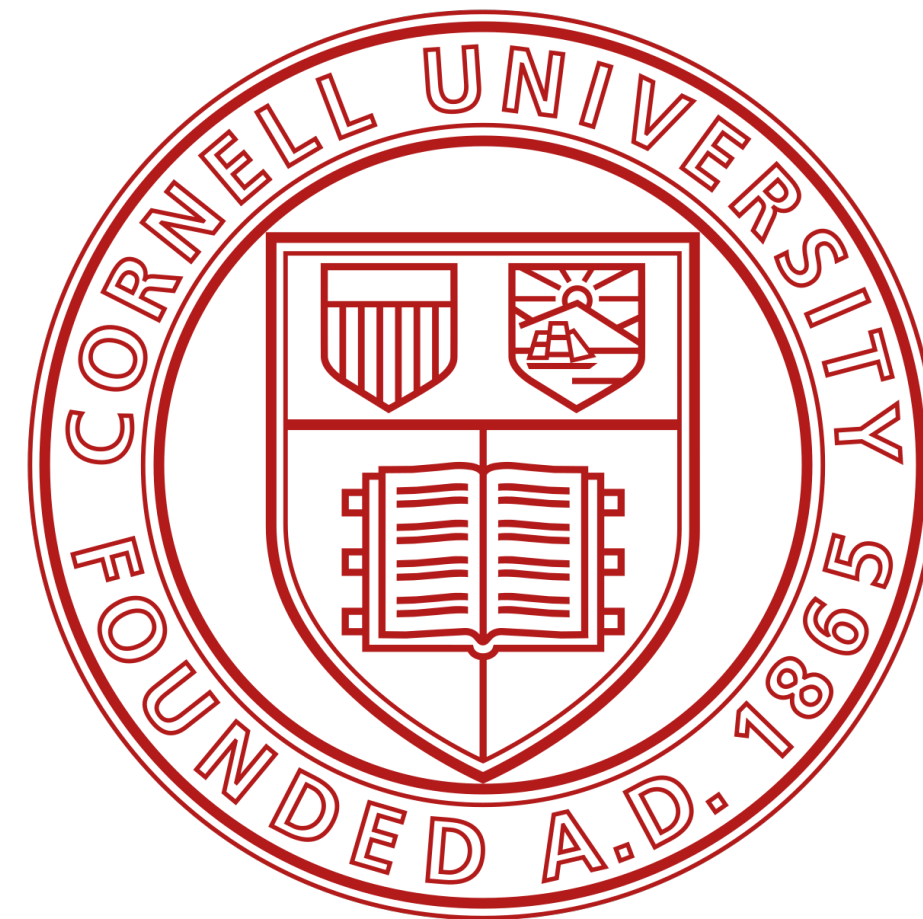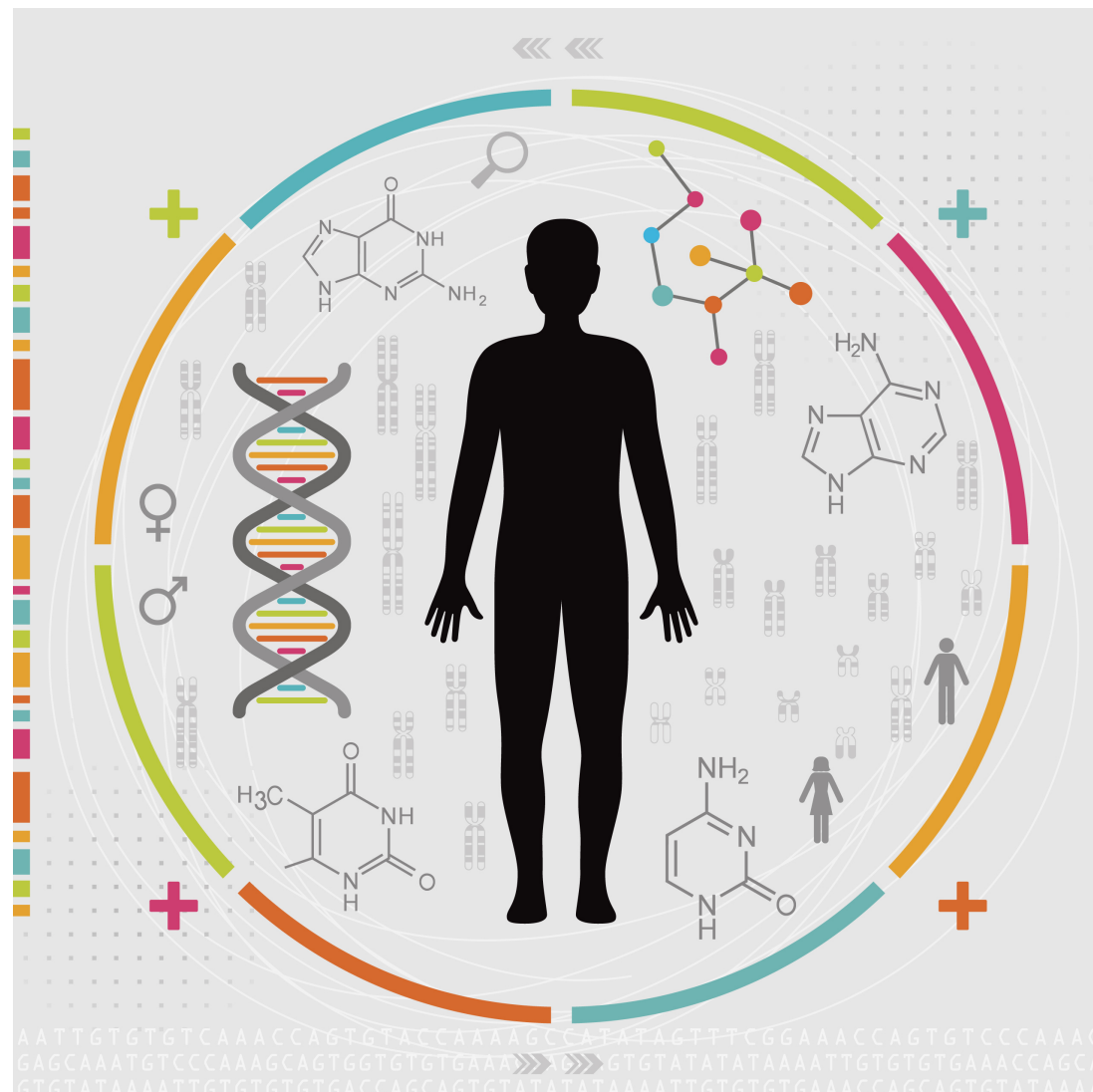# Towards Optimal Differentially Private Regret Bounds in Linear MDPs
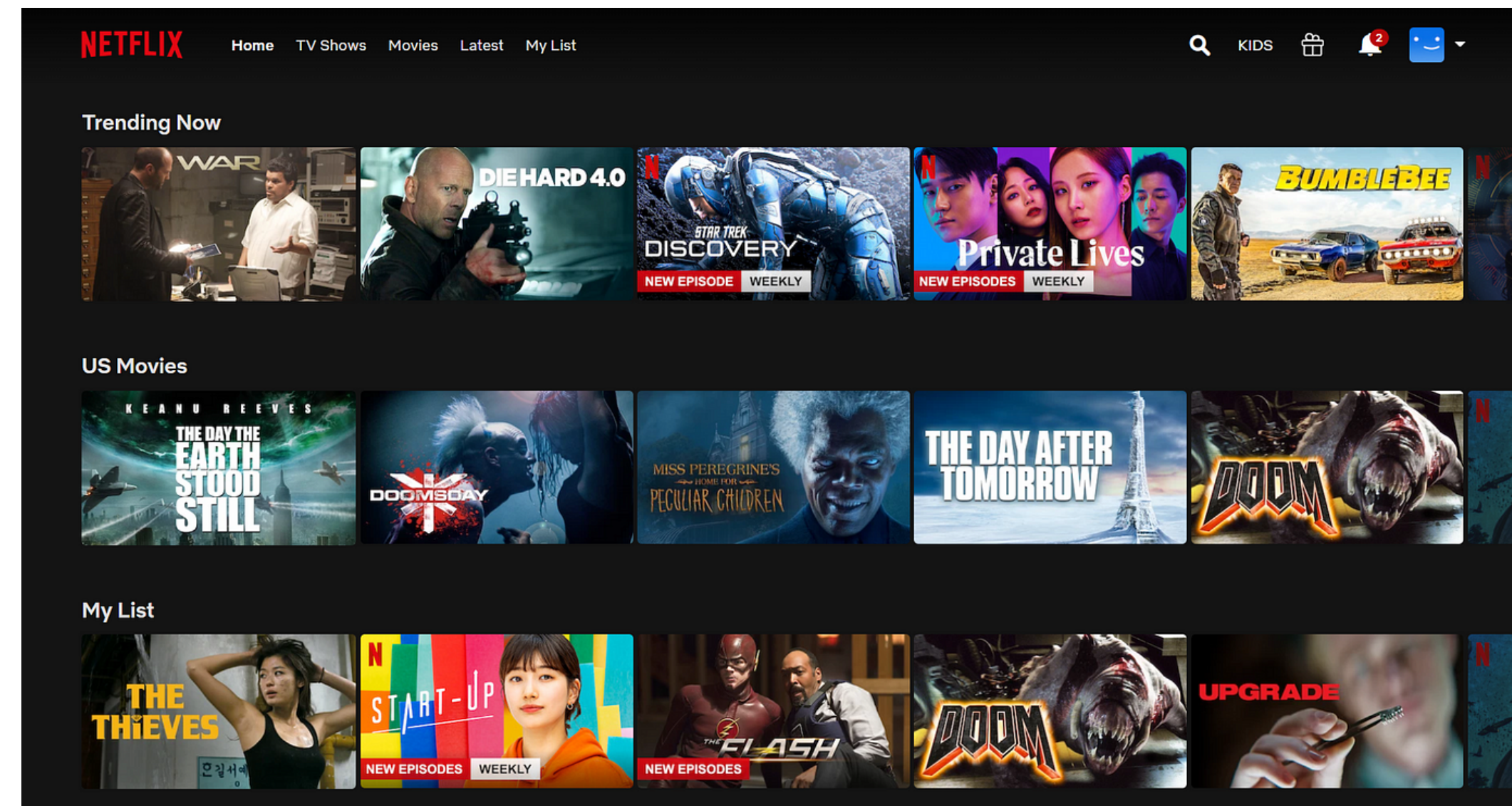


**Sharan Sahu, Statistics and Data Science, Cornell University**

# Recent Successes In Reinforcement Learning (RL)



Precision Medicine



Recommender Systems



Autonomous Driving

# Contextual Bandits

**User**



**Actions**

$$a_t^1$$
$$a_t^2$$
$$\vdots$$
$$a_t^k$$

Decision-Making Agent
$$\pi$$

Response
$$Y_i \in \{0,1\}$$

# Contextual Bandits

**User**

1

$s_t$

Decision-Making Agent

$\pi$

Actions

$a_t^1$
$a_t^2$
$\vdots$
$a_t^k$

Response
$Y_i \in \{0,1\}$

# Contextual Bandits

**User**

$s_t$

$\boxed{1}$

**Decision-Making Agent**
$\pi$

$\pi(a_t \mid s_t)$

$\boxed{2}$

**Actions**

$a_t^1$
$a_t^2$
$\vdots$
$a_t^k$

**Response**
$Y_i \in \{0,1\}$

# Contextual Bandits

**User**

$s_t$

We recommend
taking action $a_t$

**Decision-Making Agent**

$\pi$

$\pi(a_t \mid s_t)$

**Actions**

$a_t^1$
$a_t^2$
$\vdots$
$a_t^k$

**1**

**2**

**3**

**Response**
$Y_i \in \{0,1\}$

# Contextual Bandits

**User**



1

$s_t$

Decision-Making Agent
$\pi$

We recommend
taking action $a_t$

3

2

$\pi(a_t \mid s_t)$

4

**Actions**

$a_t^1$
$a_t^2$
$\vdots$
$a_t^k$

**Response**
$Y_i \in \{0,1\}$

1  Agent observes user state

2  Agent selects action based on user state

3  Selected action is provided to user

4  User provides a response for agent to learn

7

# Reasonable Policies May Use Sensitive Data

User 1

User 2

User 3

$s_0^{(1)}$

$s_0^{(2)}$

$s_0^{(3)}$

App

Good, Reasonable Policy

$\pi$

# Reasonable Policies May Use Sensitive Data



$s_0^{(1)}$

Name
Age
Weight
Health Habits
Physical Activity Levels
Health Issues
⋮

The policy has **access to information** that users may consider **sensitive** or **private**

# Neural Networks Can Memorize Personal Information From One Example

Anonymisation Fails
Single sample with personal features

John

DNN trained for
X-ray classification

Inserting the (memorised) unique feature
changes prediction

JOHN

**Hartley et al. 2023**

# We Must Incorporate Privacy-Preserving Mechanisms Into RL

**We require a mathematically rigorous framework that provides statistical guarantees for our (possibly randomized) mechanism:**

**Definition (Approximate Differential Privacy).** A mechanism $\mathcal{M}$ is $(\varepsilon, \delta)$-DP if for all neighboring datasets $\mathcal{U}, \mathcal{U}'$ that differ by one record and for all event $E$ in the output range

$$\mathbb{P}\left(\mathcal{M}(\mathcal{U}) \in E\right) \leq e^{\varepsilon}\mathbb{P}\left(\mathcal{M}(\mathcal{U}') \in E\right) + \delta$$

**Remark**: This is a relaxation of $\varepsilon$-DP as in many settings, achieving $\varepsilon$-DP is nearly impossible or comes at high utility cost

# Differential Privacy (DP)

Database $D_1$

+

Some other people's data

=

Database $D_2$

Mechanism $\mathscr{M}$

Mechanism $\mathscr{M}$

$\mathscr{M}\left(D_1\right)$

$\approx$

$\mathscr{M}\left(D_2\right)$

Mechanism $\mathscr{M}$ is differentially private if …

$\forall D_1, D_2$ that differ by at-most one record

$\mathscr{M}\left(D_1\right), \mathscr{M}\left(D_2\right)$ are indistinguishable

12

# There are a few issues with $(\varepsilon, \delta)$-DP

**Trusted Individual of the Central Agency**

**User $u_1$**

**Agent $\pi$**

$(\varepsilon, \delta)$-DP Mechanism $\mathscr{M}$

# There are a few issues with $(\varepsilon, \delta)$-DP

I trust this agent
with my sensitive
raw data $\mathscr{D}_{u_1}$

**User $u_1$**

**Agent $\pi$**

$(\varepsilon, \delta)$-DP

Mechanism $\mathscr{M}$

# There are a few issues with $(\varepsilon, \delta)$-DP

I trust this agent
with my sensitive
raw data $\mathscr{D}_{u_1}$

Query $Q(\mathscr{D}_{u_1})$

$(\varepsilon, \delta)$-DP

Mechanism $\mathscr{M}$

**User $u_1$**

**Agent $\pi$**

# There are a few issues with $(\varepsilon, \delta)$-DP

I trust this agent with my sensitive raw data $\mathscr{D}_{u_1}$

Query $Q(\mathscr{D}_{u_1})$

Noisy response $\mathscr{M}(Q(\mathscr{D}_{u_1}))$

$(\varepsilon, \delta)$-DP Mechanism $\mathscr{M}$

**User** $u_1$

**Agent** $\pi$

# There are a few issues with $(\varepsilon, \delta)$-DP



**Trusted Individual of the Central Agency**

I trust this agent with my sensitive raw data $\mathscr{D}_{u_1}$

Query $Q(\mathscr{D}_{u_1})$

$(\varepsilon, \delta)$-DP Mechanism $\mathscr{M}$

Agent $\pi$ recommends $a \sim \pi\left( \cdot \mid \mathscr{M}(Q(\mathscr{D}_{u_1})) \right)$

Noisy response $\mathscr{M}(Q(\mathscr{D}_{u_1}))$

**User $u_1$**

**Agent $\pi$**

# There are a few issues with $(\varepsilon, \delta)$-DP

I trust this agent with
my **new** sensitive
raw data $\mathscr{D}'_{u_1}$

Query $Q(\mathscr{D}'_{u_1})$

$(\varepsilon, \delta)$-DP
Mechanism $\mathscr{M}$

Agent $\pi$
recommends
$a \sim \pi\left( \, \cdot \mid \mathscr{M}(Q(\mathscr{D}'_{u_1})) \right)$

Noisy
response
$\mathscr{M}(Q(\mathscr{D}'_{u_1}))$

**User** $u_1$

**Agent** $\pi$

# There are a few issues with $(\varepsilon, \delta)$-DP



My recommendations didn't change even though I changed my data!

**Trusted Individual of the Central Agency**

Query $Q(\mathscr{D}'_{u_1})$

$(\varepsilon, \delta)$-DP Mechanism $\mathscr{M}$

**User** $u_1$

Agent $\pi$ recommends
$a \sim \pi\left( \cdot \mid \mathscr{M}(Q(\mathscr{D}'_{u_1})) \right)$

Noisy response $\mathscr{M}(Q(\mathscr{D}'_{u_1}))$

**Agent** $\pi$

# There are a few issues with $(\varepsilon, \delta)$-DP



Due to the $(\varepsilon, \delta)$-DP mechanism $\mathcal{M}$, any change in the dataset, including your own data, cannot change the output too much!

**Trusted Individual of the Central Agency**

Query $Q(\mathcal{D}'_{u_1})$

$(\varepsilon, \delta)$-DP Mechanism $\mathcal{M}$

Noisy response $\mathcal{M}(Q(\mathcal{D}'_{u_1}))$

**User** $u_1$

Agent $\pi$ recommends $a \sim \pi\left( \cdot \mid \mathcal{M}(Q(\mathcal{D}'_{u_1})) \right)$

**Agent** $\pi$

# There are a few issues with $(\varepsilon, \delta)$-DP



**Trusted Individual of the Central Agency**

I trust this agent with my sensitive raw data $\mathcal{D}_{u_2}$

Query $Q(\mathcal{D}_{u_2})$

$(\varepsilon, \delta)$-DP Mechanism $\mathcal{M}$

Agent $\pi$ recommends $a \sim \pi\left( \cdot \mid \mathcal{M}(Q(\mathcal{D}_{u_2})) \right)$

Noisy response $\mathcal{M}(Q(\mathcal{D}_{u_2}))$

**User** $u_2$

**Agent** $\pi$

# There are a few issues with $(\varepsilon, \delta)$-DP



Why are my recommendations the same as $u_1$? We are completely different people!

**Trusted Individual of the Central Agency**

**User** $u_2$

Agent $\pi$ recommends
$a \sim \pi\left( \cdot \mid \mathcal{M}(Q(\mathcal{D}_{u_2})) \right)$

**Agent** $\pi$

Query $Q(\mathcal{D}_{u_2})$

Noisy response
$\mathcal{M}(Q(\mathcal{D}_{u_2}))$

$(\varepsilon, \delta)$-DP
Mechanism $\mathcal{M}$

# There are a few issues with $(\varepsilon, \delta)$-DP



Giving vastly different recommendations between users would violate the guarantees of the $(\varepsilon, \delta)$-DP mechanism $\mathcal{M}$

**Trusted Individual of the Central Agency**

Query $Q(\mathcal{D}_{u_2})$

$(\varepsilon, \delta)$-DP Mechanism $\mathcal{M}$

Noisy response $\mathcal{M}(Q(\mathcal{D}_{u_2}))$

**User** $u_2$

Agent $\pi$ recommends $a \sim \pi \left( \cdot \mid \mathcal{M}(Q(\mathcal{D}_{u_2})) \right)$
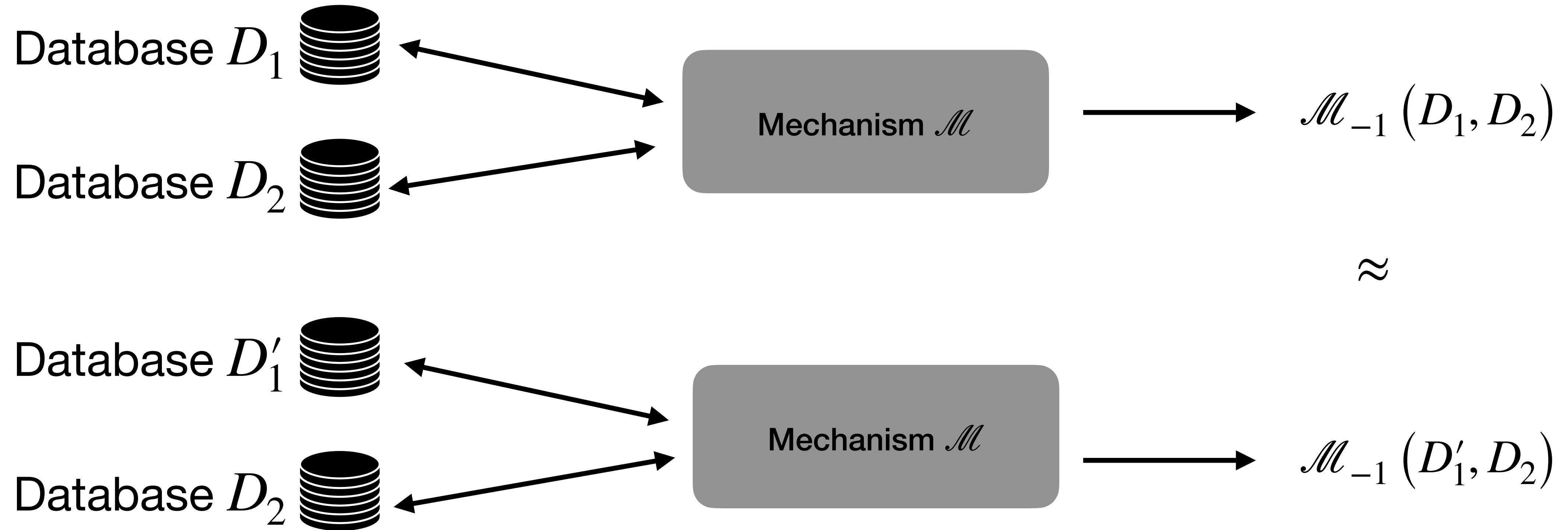
**Agent** $\pi$

# We need a further relaxation of DP …

One which works nicely with contextual bandit problems on a per-user level but does not sacrifice privacy on a per-decision or per-context level, ensuring that individual contexts do not overly influence the learned policy:

**Definition (Approximate Joint Differential Privacy).** A mechanism $\mathcal{M}$ is $(\varepsilon, \delta)$-JDP if for any $k \in [K]$, any user sequences $\mathcal{U}, \mathcal{U}'$ differing on the $k$-user and any $E \subset \mathcal{A}^{(K-1)H}$

$$\mathbb{P}\left(\mathcal{M}_{-k}(\mathcal{U}) \in E\right) \leq e^{\varepsilon}\mathbb{P}\left(\mathcal{M}_{-k}(\mathcal{U}') \in E\right)$$

# Joint Differential Privacy (JDP)

Database $D_1$

Database $D_2$

Mechanism $\mathscr{M}$

$$\mathscr{M}_{-1}\left(D_1, D_2\right)$$

$\approx$

Database $D_1'$

Database $D_2$

Mechanism $\mathscr{M}$

$$\mathscr{M}_{-1}\left(D_1', D_2\right)$$

Mechanism $\mathscr{M}$ is joint differentially private if ...

$\forall \left(D_1, D_2, \ldots, D_k\right), \left(D_1', D_2, \ldots, D_k\right)$ where only one party's data differs by at most one record

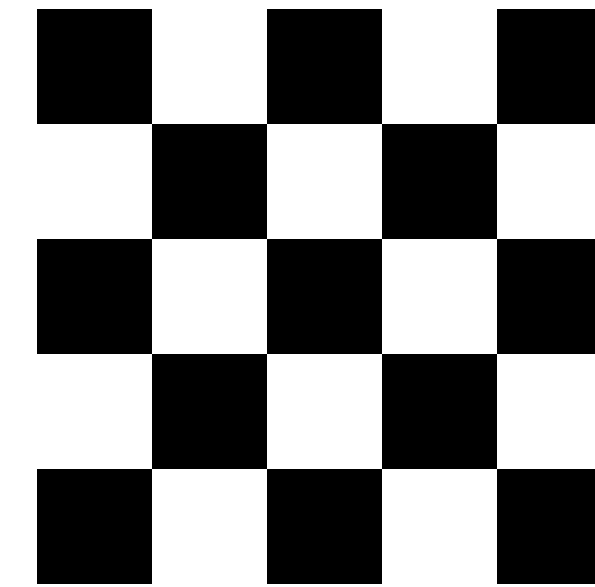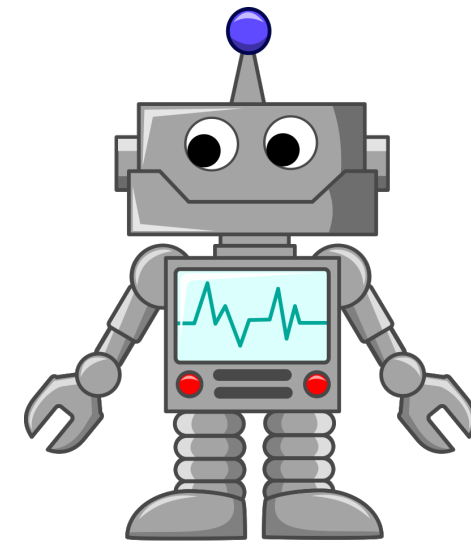$\mathscr{M}_{-1}\left(D_1, D_2\right), \mathscr{M}_{-1}\left(D_1', D_2\right)$ are indistinguishable

# In this talk:

Can we develop an efficient $(\varepsilon, \delta)$-JDP algorithm for sequential decision-making problems with **linear parametric representations**, and provide a novel algorithm with provably efficient guarantees for **privacy-preserving exploration**?
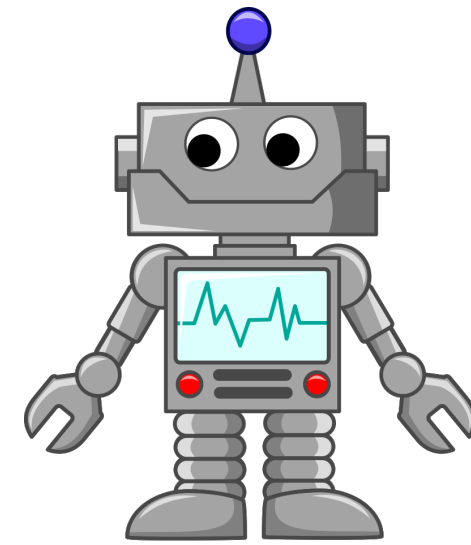
# **Outline**

1. Problem Setup + Previous Work and Motivation

2. Can we do better?

3. Our regret bound
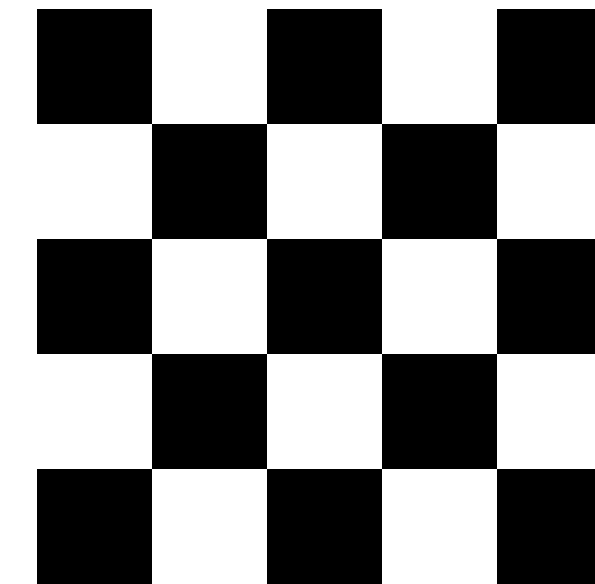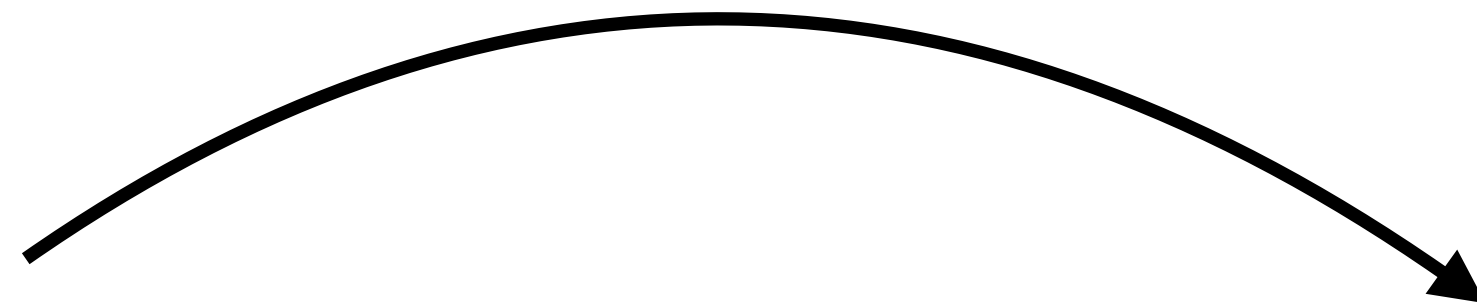
# Episodic Time-Inhomogeneous Finite-Horizon MDPs

# Episodic Time-Inhomogeneous Finite-Horizon MDPs
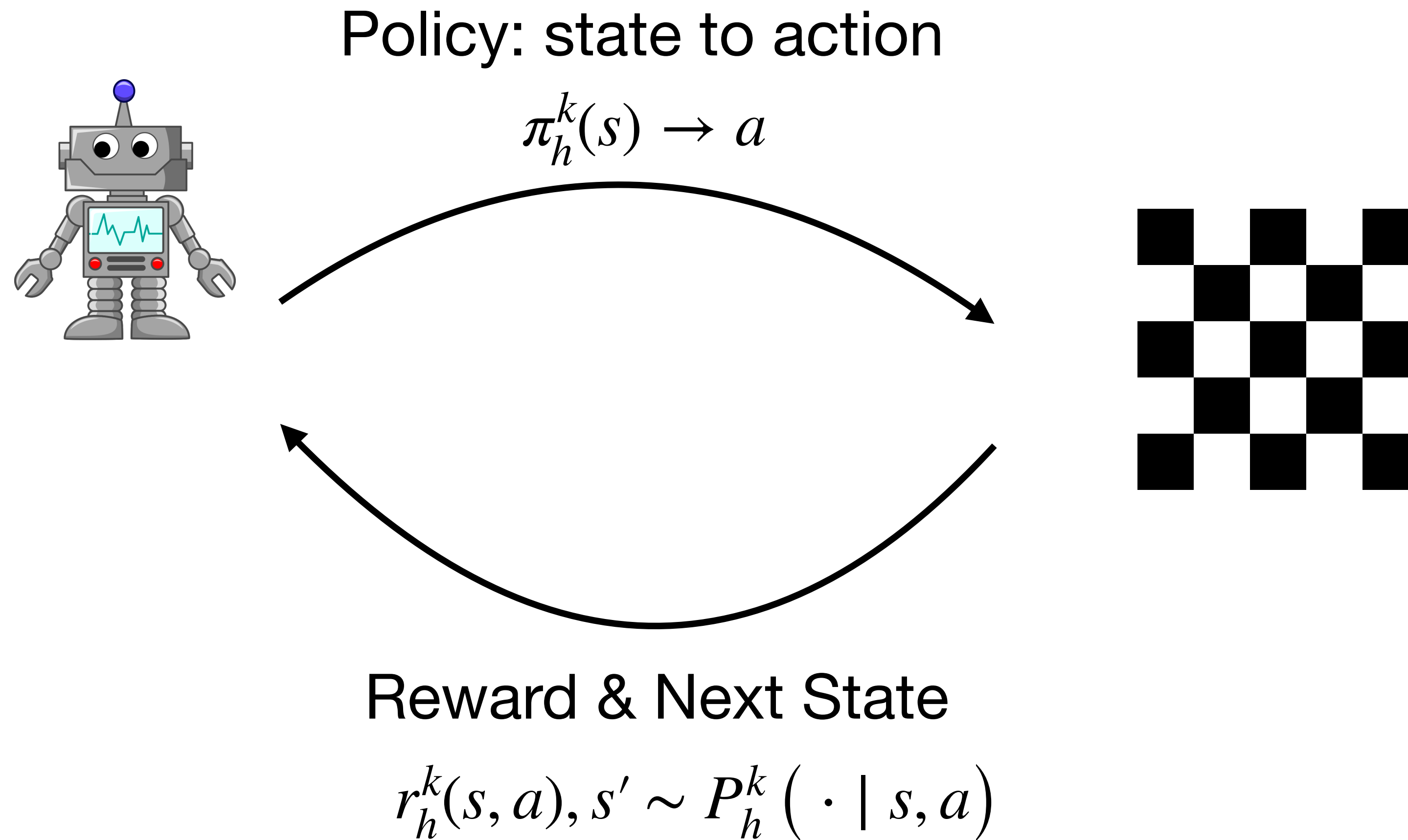
Policy: state to action

$$\pi_h^k(s) \to a$$

# Episodic Time-Inhomogeneous Finite-Horizon MDPs

Policy: state to action

$$\pi_h^k(s) \to a$$

Reward & Next State

$$r_h^k(s, a), s' \sim P_h^k \left( \cdot \mid s, a \right)$$

# Episodic Time-Inhomogeneous Finite-Horizon MDPs

Policy: state to action

$$\pi_h^k(s) \to a$$

$H$ Times

Reward & Next State

$$r_h^k(s, a), s' \sim P_h^k\left(\cdot \mid s, a\right)$$

# Episodic Time-Inhomogeneous Finite-Horizon MDPs

Policy: state to action

$$\pi_h^k(s) \to a$$

$$\forall k \in [K],$$
play $H$ times

Reward & Next State

$$r_h^k(s, a), s' \sim P_h^k \left( \, \cdot \mid s, a \right)$$

# Episodic Time-Inhomogeneous Finite-Horizon MDPs

Policy: state to action

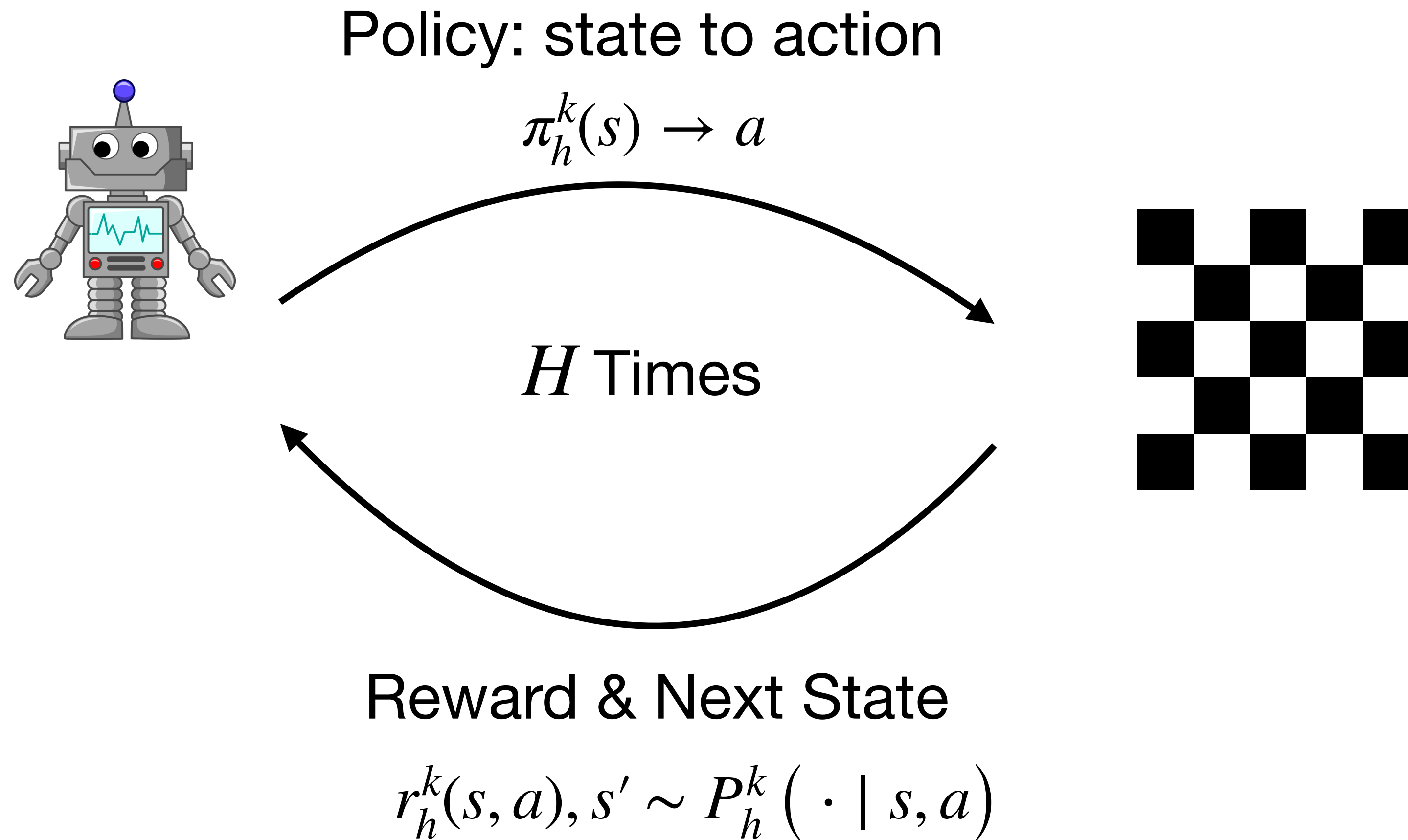$$\pi_h^k(s) \to a$$

$$\tau^k = \left\{ s_h^k, a_h^k \right\}_{h=1}^H$$

$$\forall k \in [K],$$
play $H$ times

Reward & Next State

$$r_h^k(s, a), s' \sim P_h^k \left( \cdot \mid s, a \right)$$

33

# Episodic Time-Inhomogeneous Finite-Horizon MDPs

Policy: state to action

$$\pi_h^k(s) \rightarrow a$$

$$\tau^k = \left\{ s_h^k, a_h^k \right\}_{h=1}^{H}$$

$$\forall k \in [K],$$
play $H$ times

Reward & Next State

$$r_h^k(s, a), s' \sim P_h^k \left( \cdot \mid s, a \right)$$

**Finite-Horizon MDP:** $\mathcal{M} = \left\{ \mathcal{S}, \mathcal{A}, \left\{ r_h \right\}_{h=1}^{H}, \left\{ P_h \right\}_{h=1}^{H}, H \right\}$  $H < \infty$

# Formal RL Problem Setting

**Setting:** Episodic inhomogeneous finite horizon MDP $\mathcal{M} = \{\mathcal{S}, \mathcal{A}, \{\mathbb{P}_h\}_h, \{r_h\}_h, H\}$ where $\mathcal{S}, \mathcal{A}$ are the states and actions, respectively, $H \in \mathbb{Z}$ is the length of each episode, $\mathbb{P}_h : \mathcal{S} \times \mathcal{A} \to \Delta(\mathcal{S})$, $r_h : \mathcal{S} \times \mathcal{A} \to [0,1]$ are the time-dependent transition probability and deterministic reward function. $\mathcal{S}$ is measurable and possibly uncountable, and $\mathcal{A}$ is finite. In this setting, the policy is time-dependent and we denote this $\pi = \{\pi_1, \cdots, \pi_H\}$

$$\text{Regret}(K) = \sum_{k=1}^{K} \left[ V_1^* \left( s_1^k \right) - V_1^{\pi_k} \left( s_1^k \right) \right]$$

# Linear MDP

$(s,a)$ $\bullet$ $=$ $(s,a)$ $\phi(s,a)$

$r_h \in \mathbb{R}^{SA}$  $\quad$  $\Phi \in \mathbb{R}^{SA \times d}$ $\quad$ $\theta_h \in \mathbb{R}^d$

$$\exists \theta_h, \phi^\star : \forall s, a, h, \; r_h(s,a) = \phi^\star(s,a)^\top \theta_h$$

# Linear MDP



$$\exists \mu_h, \phi^\star : \forall s, a, h, s', P_h(s' \mid s, a) = \phi^\star(s, a)^\top \mu_h(s')$$

# LSVI-UCB

**Algorithm 1** Least-Squares Value Iteration with UCB (LSVI-UCB)

1: **for** episode $k = 1, \ldots, K$ **do**
2:      Receive the initial state $x_1^k$.
3:      **for** step $h = H, \ldots, 1$ **do**
4:           $\Lambda_h \leftarrow \sum_{\tau=1}^{k-1} \phi(x_h^\tau, a_h^\tau) \phi(x_h^\tau, a_h^\tau)^\top + \lambda \cdot \mathbf{I}$.
5:           $\mathbf{w}_h \leftarrow \Lambda_h^{-1} \sum_{\tau=1}^{k-1} \phi(x_h^\tau, a_h^\tau)[r_h(x_h^\tau, a_h^\tau) + \max_a Q_{h+1}(x_{h+1}^\tau, a)]$.
6:           $Q_h(\cdot, \cdot) \leftarrow \min\{\mathbf{w}_h^\top \phi(\cdot, \cdot) + \beta[\phi(\cdot, \cdot)^\top \Lambda_h^{-1} \phi(\cdot, \cdot)]^{1/2}, H\}$.
7:      **for** step $h = 1, \ldots, H$ **do**
8:           Take action $a_h^k \leftarrow \mathrm{argmax}_{a \in \mathcal{A}} Q_h(x_h^k, a)$, and observe $x_{h+1}^k$.

**LSVI-UCB Algorithm [Jin et al. 2020]**

38

# LSVI-UCB

**Algorithm 1** Least-Squares Value Iteration with UCB (LSVI-UCB)

1: **for** episode $k = 1, \ldots, K$ **do**
2:     Receive the initial state $x_1^k$.
3:     **for** step $h = H, \ldots, 1$ **do**
4:         $\Lambda_h \leftarrow \sum_{\tau=1}^{k-1} \phi(x_h^\tau, a_h^\tau)\phi(x_h^\tau, a_h^\tau)^\top + \lambda \cdot \mathbf{I}$.
5:         $\mathbf{w}_h \leftarrow \Lambda_h^{-1} \sum_{\tau=1}^{k-1} \phi(x_h^\tau, a_h^\tau)[r_h(x_h^\tau, a_h^\tau) + \max_a Q_{h+1}(x_{h+1}^\tau, a)]$.
6:         $Q_h(\cdot, \cdot) \leftarrow \min\{\mathbf{w}_h^\top \phi(\cdot, \cdot) + \beta[\phi(\cdot, \cdot)^\top \Lambda_h^{-1} \phi(\cdot, \cdot)]^{1/2}, H\}$.
7:     **for** step $h = 1, \ldots, H$ **do**
8:         Take action $a_h^k \leftarrow \mathrm{argmax}_{a \in \mathcal{A}} Q_h(x_h^k, a)$, and observe $x_{h+1}^k$.

**LSVI-UCB Algorithm [Jin et al. 2020]**

# LSVI-UCB

**Algorithm 1** Least-Squares Value Iteration with UCB (LSVI-UCB)

1: **for** episode $k = 1, \ldots, K$ **do**
2:      Receive the initial state $x_1^k$.
3:      **for** step $h = H, \ldots, 1$ **do**
4:          $\Lambda_h \leftarrow \sum_{\tau=1}^{k-1} \phi(x_h^\tau, a_h^\tau) \phi(x_h^\tau, a_h^\tau)^\top + \lambda \cdot \mathbf{I}$.
5:          $\mathbf{w}_h \leftarrow \Lambda_h^{-1} \sum_{\tau=1}^{k-1} \phi(x_h^\tau, a_h^\tau)[r_h(x_h^\tau, a_h^\tau) + \max_a Q_{h+1}(x_{h+1}^\tau, a)]$.
6:          $Q_h(\cdot, \cdot) \leftarrow \min\{\mathbf{w}_h^\top \phi(\cdot, \cdot) + \beta[\phi(\cdot, \cdot)^\top \Lambda_h^{-1} \phi(\cdot, \cdot)]^{1/2}, H\}$.
7:      **for** step $h = 1, \ldots, H$ **do**
8:          Take action $a_h^k \leftarrow \mathrm{argmax}_{a \in \mathcal{A}} Q_h(x_h^k, a)$, and observe $x_{h+1}^k$.

**Think of $\phi(s, a)$ as one-hot vector, then $\Lambda_h$ is capturing something similar to visitation counts which uses trajectory information with possibly private data**

**LSVI-UCB Algorithm [Jin et al. 2020]**

# LSVI-UCB

**Algorithm 1** Least-Squares Value Iteration with UCB (LSVI-UCB)

1: **for** episode $k = 1, \ldots, K$ **do**
2:      Receive the initial state $x_1^k$.
3:      **for** step $h = H, \ldots, 1$ **do**
4:          $\Lambda_h \leftarrow \sum_{\tau=1}^{k-1} \phi(x_h^\tau, a_h^\tau) \phi(x_h^\tau, a_h^\tau)^\top + \lambda \cdot \mathbf{I}$.
5:          $\mathbf{w}_h \leftarrow \Lambda_h^{-1} \sum_{\tau=1}^{k-1} \phi(x_h^\tau, a_h^\tau) [r_h(x_h^\tau, a_h^\tau) + \max_a Q_{h+1}(x_{h+1}^\tau, a)]$.
6:          $Q_h(\cdot, \cdot) \leftarrow \min\{\mathbf{w}_h^\top \phi(\cdot, \cdot) + \beta[\phi(\cdot, \cdot)^\top \Lambda_h^{-1} \phi(\cdot, \cdot)]^{1/2}, H\}$.
7:      **for** step $h = 1, \ldots, H$ **do**
8:          Take action $a_h^k \leftarrow \mathrm{argmax}_{a \in \mathcal{A}} Q_h(x_h^k, a)$, and observe $x_{h+1}^k$.

**LSVI-UCB Algorithm [Jin et al. 2020]**

# LSVI-UCB

**Algorithm 1** Least-Squares Value Iteration with UCB (LSVI-UCB)

1: **for** episode $k = 1, \ldots, K$ **do**
2:      Receive the initial state $x_1^k$.
3:      **for** step $h = H, \ldots, 1$ **do**
4:          $\Lambda_h \leftarrow \sum_{\tau=1}^{k-1} \phi(x_h^\tau, a_h^\tau)\phi(x_h^\tau, a_h^\tau)^\top + \lambda \cdot \mathbf{I}$.
5:          $\mathbf{w}_h \leftarrow \Lambda_h^{-1} \sum_{\tau=1}^{k-1} \phi(x_h^\tau, a_h^\tau)[r_h(x_h^\tau, a_h^\tau) + \max_a Q_{h+1}(x_{h+1}^\tau, a)]$.
6:          $Q_h(\cdot, \cdot) \leftarrow \min\{\mathbf{w}_h^\top \phi(\cdot, \cdot) + \beta[\phi(\cdot, \cdot)^\top \Lambda_h^{-1} \phi(\cdot, \cdot)]^{1/2}, H\}$.
7:      **for** step $h = 1, \ldots, H$ **do**
8:          Take action $a_h^k \leftarrow \text{argmax}_{a \in \mathcal{A}} Q_h(x_h^k, a)$, and observe $x_{h+1}^k$.

**These are parameter estimates for the feature regressors which allow us to calculate the Q-function due to Linear MDPs. Again, this can leak information about trajectories taken by the policy**

**LSVI-UCB Algorithm [Jin et al. 2020]**

42

# Differential Privacy Techniques

Theorem (Gaussian Mechanism). $M_G(x, f(\,\cdot\,), \varepsilon, \delta) = f(x) + (Y_1, \ldots, Y_k)$ where $Y_i \sim \mathcal{N}\left(0, \sigma^2\right)$

with $\sigma^2 = \dfrac{\Delta_2(f)\sqrt{2\log(2/\delta)}}{\varepsilon}$ is $(\varepsilon, \delta)$-DP

Theorem (Billboard Lemma). If you have a mechanism $\mathcal{M}_{DP}$ that is $(\varepsilon, \delta)$-DP, then any function $f_i : \mathcal{U}_i \times \mathcal{R} \to \mathcal{R}_i$ that depends on user i's data and the output of the mechanism satisfies $(\varepsilon, \delta)$-JDP

# Previous Work

[Luyo et al. 2021]. Fix any privacy level $\varepsilon, \delta \in (0,1)$. For any $p \in (0,1)$, their algorithm is $(\varepsilon, \delta)$-JDP and, with probability at least $1 - p$, its regret is bounded as follows:

$$R(K) = \tilde{O}(\sqrt{d^3 H^4 K} + H^{11/5} d^{8/5} K^{3/5} / \varepsilon^{2/5})$$

## Techniques Used

$$\tilde{\Lambda}_h = \Lambda_h + \mathcal{N}\left(0, \mathcal{O}\left(\frac{1}{\varepsilon}\sqrt{BH}\log(1/\delta)\right)\right)$$

$$\tilde{u}_h = u_h + \mathcal{N}\left(0, \mathcal{O}\left(\frac{1}{\varepsilon}\sqrt{H^2 B}\log(1/\delta)\right)\right)$$

**Static Batching** to reduce the number of policy switches to $\mathcal{O}(\text{poly}(K))$

# Previous Work

[Ngo et al. 2022]. Fix any privacy level $\varepsilon, \delta \in (0,1)$. For any $p \in (0,1)$, their algorithm is $(\varepsilon, \delta)$-JDP and, with probability at least $1 - p$, its regret is bounded as follows:

$$R(K) = \tilde{O}(\sqrt{d^3 H^4 K} + H^3 d^{5/4} K^{1/2}/\varepsilon^{1/2})$$

## Techniques Used

Same techniques as previous work but instead of a static batching schedule, they use **Adaptive Batching** to reduce the number of policy switches to $\mathcal{O}(\log(K))$

# LSVI-UCB

---

**Algorithm 1** Least-Squares Value Iteration with UCB (LSVI-UCB)

---

1: **for** episode $k = 1, \ldots, K$ **do**
2:     Receive the initial state $x_1^k$.
3:     **for** step $h = H, \ldots, 1$ **do**
4:         $\Lambda_h \leftarrow \sum_{\tau=1}^{k-1} \phi(x_h^\tau, a_h^\tau) \phi(x_h^\tau, a_h^\tau)^\top + \lambda \cdot \mathbf{I}$.
5:         $\mathbf{w}_h \leftarrow \Lambda_h^{-1} \sum_{\tau=1}^{k-1} \phi(x_h^\tau, a_h^\tau)[r_h(x_h^\tau, a_h^\tau) + \max_a Q_{h+1}(x_{h+1}^\tau, a)]$.
6:         $Q_h(\cdot, \cdot) \leftarrow \min\{\mathbf{w}_h^\top \phi(\cdot, \cdot) + \beta[\phi(\cdot, \cdot)^\top \Lambda_h^{-1} \phi(\cdot, \cdot)]^{1/2}, H\}$.
7:     **for** step $h = 1, \ldots, H$ **do**
8:         Take action $a_h^k \leftarrow \mathrm{argmax}_{a \in \mathcal{A}} Q_h(x_h^k, a)$, and observe $x_{h+1}^k$.

---

**LSVI-UCB Algorithm [Jin et al. 2020]**

Achieves regret $R(K) = \tilde{\mathcal{O}}\left(H^2 \sqrt{d^3 K}\right)$

46

# Motivating Work: LSVI-UCB+

[Hu et al. 2022]. Set $\lambda = 1/\left(H^2\sqrt{d}\right)$. Then, with probability at least $1 - 10\delta$, the regret of LSVI-UCB+ is upper bounded by

$$R(K) = \tilde{\mathcal{O}}\left(d\sqrt{H^3 K}\right)$$

## Techniques Used

Instead of solving a **ridge regression** problem, we solve a **weighted ridge regression** problem using estimated weights from data. This allows us to use a self-normalized martingale argument using **Azuma-Bernstein** rather than **Azuma-Hoeffding** to get a bonus that improves our regret

# Motivating Work: JDP In Tabular MDPs

[Qiao and Wang. 2023]. For any privacy budget $\varepsilon > 0$, failure probability $0 < \beta < 1$, and any privatizer where the private counts are close to the true counts with high probability, with probability at least $1 - \beta$, their algorithm is $(\varepsilon, \delta)$-JDP and achieves regret upper bounded by:

$$R(K) = \tilde{\mathcal{O}}\left(\sqrt{H^3 SAK} + S^2 AH^3/\varepsilon\right)$$

## Techniques Used

In previous work, since we would use a **Hoeffding-bound** that only depends on the counts, it is sufficient to **privatize the counts** loosely using Gaussian noise with sufficient variance component. However, to use a **Bernstein-bound**, we need to **carefully privatize the counts** to ensure that we can upper bound the variance term in a Bernstein-bound

**Can we design a $(\varepsilon, \delta)$-JDP algorithm that is near minimax optimal for non-private learning and improves the cost of privacy using more refined privatization and concentration techniques?**

$$R(K) = \tilde{O}(\sqrt{d^3 H^4 K} + H^3 d^{5/4} K^{1/2}/\varepsilon^{1/2})$$

**Can we design a $(\varepsilon, \delta)$-JDP algorithm that is near minimax optimal for non-private learning and improves the cost of privacy using more refined privatization and concentration techniques?**

$$R(K) = \tilde{O}(\sqrt{d^3 H^4 K} + H^3 d^{5/4} K^{1/2}/\varepsilon^{1/2})$$

**Non-private learning regret: We can do better using LSVI-UCB+**

**Can we design a $(\varepsilon, \delta)$-JDP algorithm that is near minimax optimal for non-private learning and improves the cost of privacy using more refined privatization and concentration techniques?**

$$R(K) = \tilde{O}(\sqrt{d^3 H^4 K} + H^3 d^{5/4} K^{1/2}/\varepsilon^{1/2})$$

**Non-private learning regret: We can do better using LSVI-UCB+**

**Cost of privacy: can we improve this**

$$\mathcal{O}\left(\text{poly}\left(HdK\right)/\varepsilon\right)$$

# Our Work

Fix any privacy level $\varepsilon, \delta \in (0,1)$. For any $p \in (0,1)$, their algorithm is $(\varepsilon, \delta)$-JDP and, with probability at least $1 - p$, its regret is bounded as follows:

$$R(K) = \tilde{O}\left( d\sqrt{H^3 K} + \frac{H^{19/8} d^{15/8} K^{3/4}}{\varepsilon} \right)$$

Compared to Luyo et al. (2021) and Ngo et al. (2022), this regret bound achieves tighter dependence on $H, d$ for the non-private terms and tighter dependence on $H, \varepsilon$ for the private terms

# Proof Sketch

1. Identify terms in the non-private algorithm that are used for estimating

2. Privatize them by (cleverly) adding noise to the terms

3. Prove the utility of the privatized terms (how close are they to the non-private terms)

4. Use the private terms in place of the non-private terms and use your standard LSVI-UCB techniques (i.e. self-normalized martingale concentrations, uniform covering arguments, elliptical potentials, and utility of the privatized terms)

# Questions?