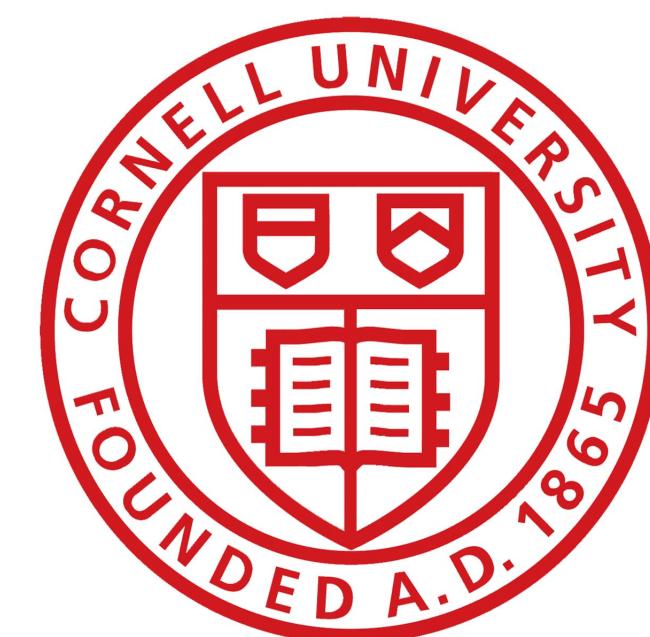


Towards Optimal Differentially Private Regret Bounds in Linear MDPs

Sharan Sahu | Stats and Data Sci. PhD | Cornell University

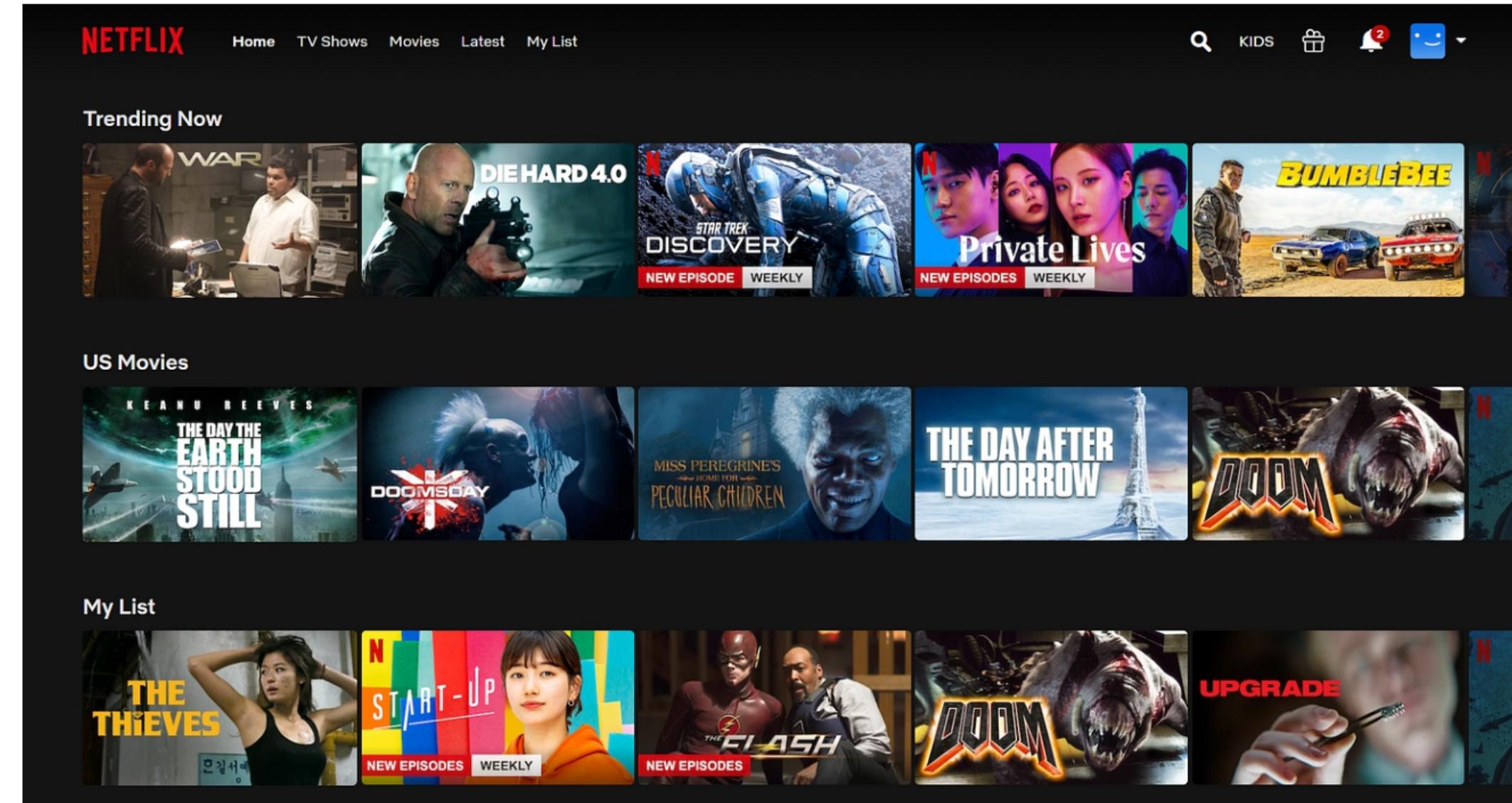


Cornell University®

Recent Successes of Reinforcement Learning (RL)



Precision Medicine

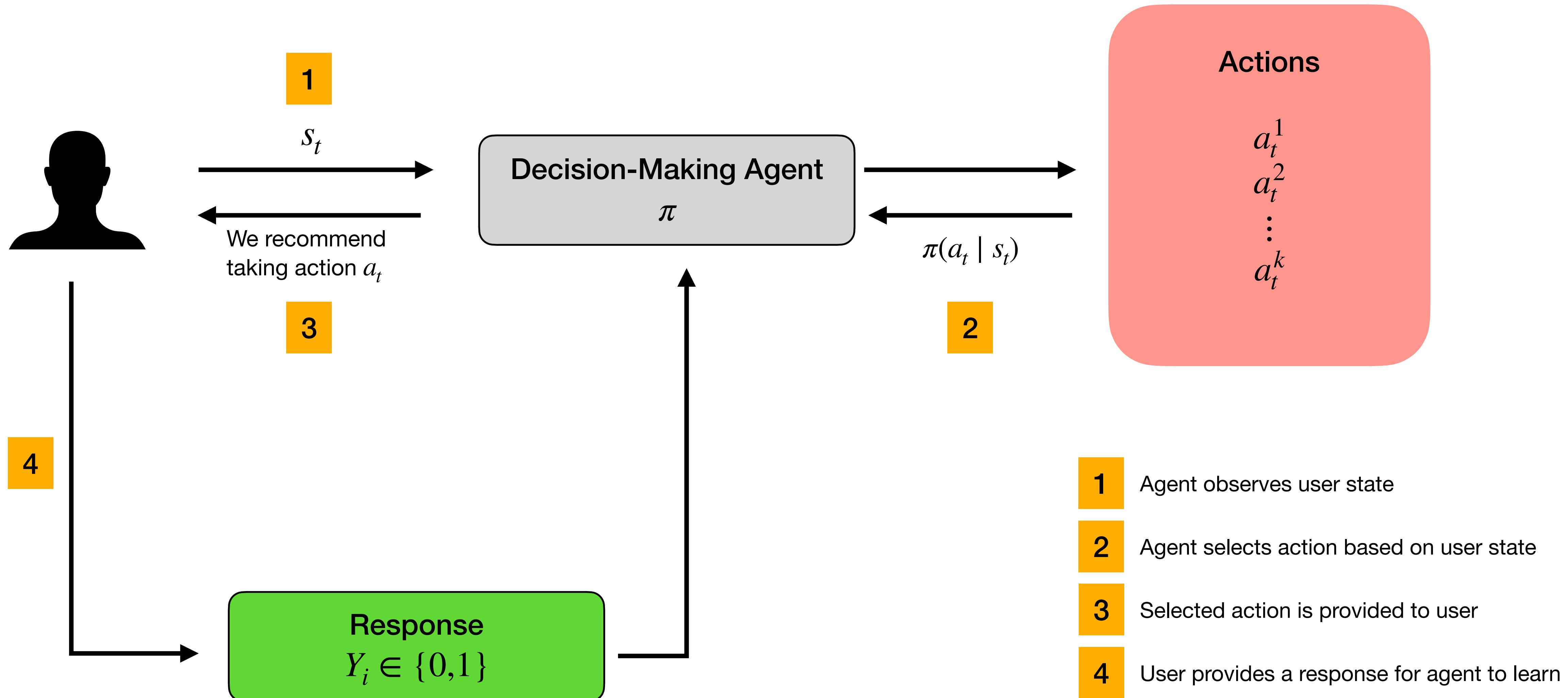


Recommender Systems

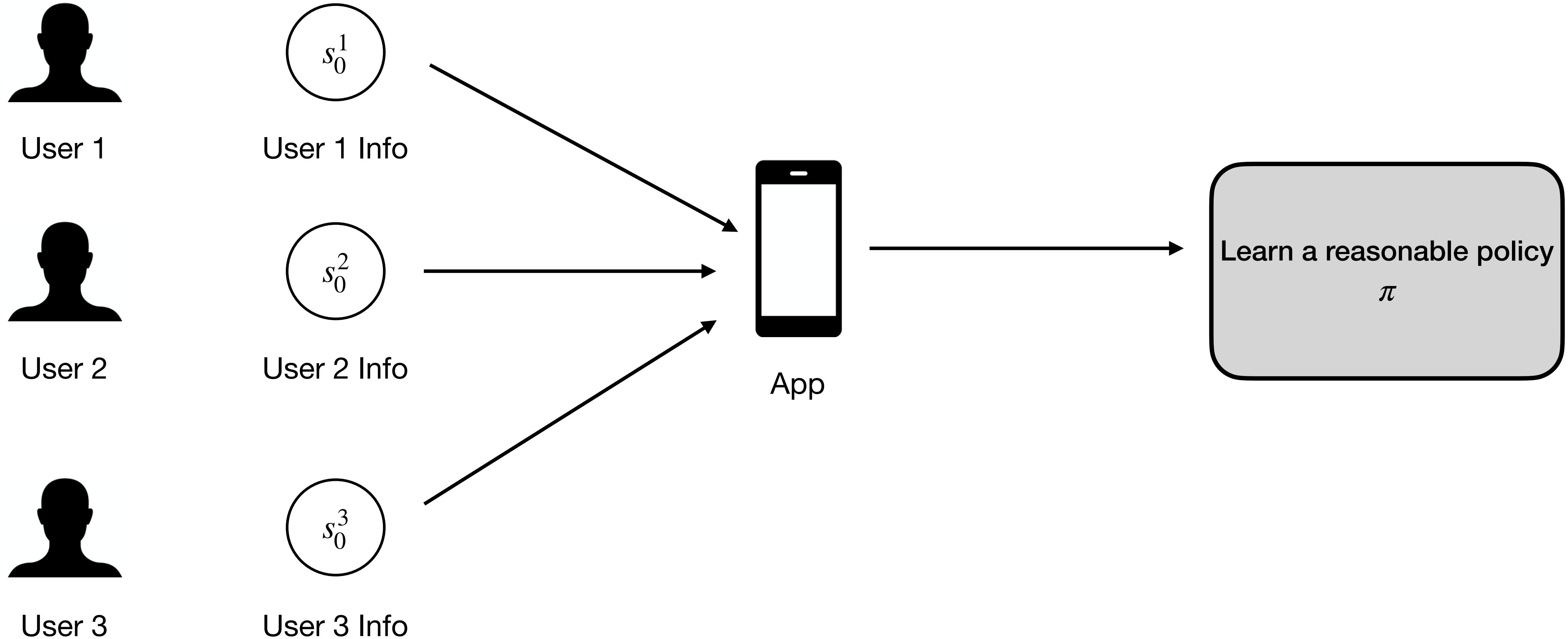


Autonomous Driving

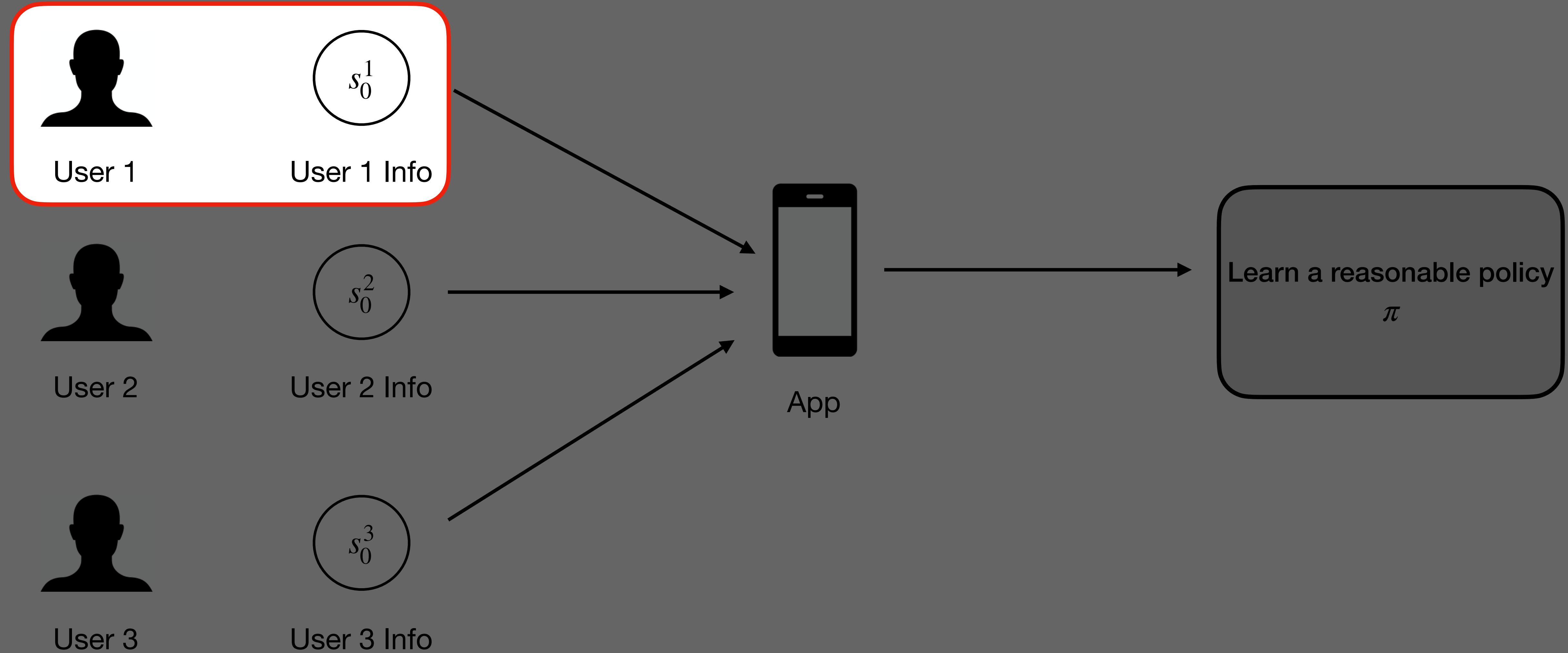
Contextual Bandits



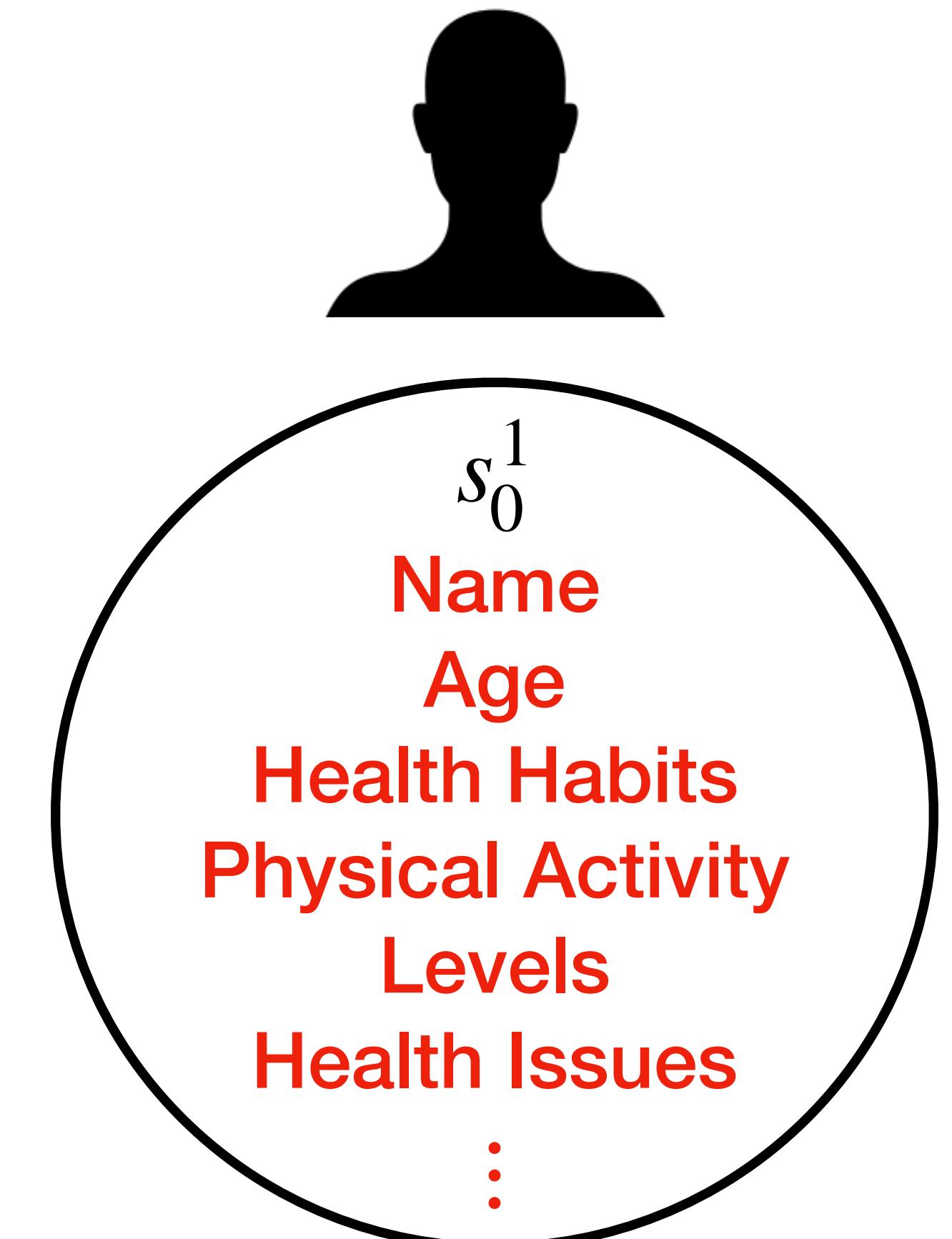
Reasonable Policies May Use Sensitive Data



Reasonable Policies May Use Sensitive Data



Reasonable Policies May Use Sensitive Data



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

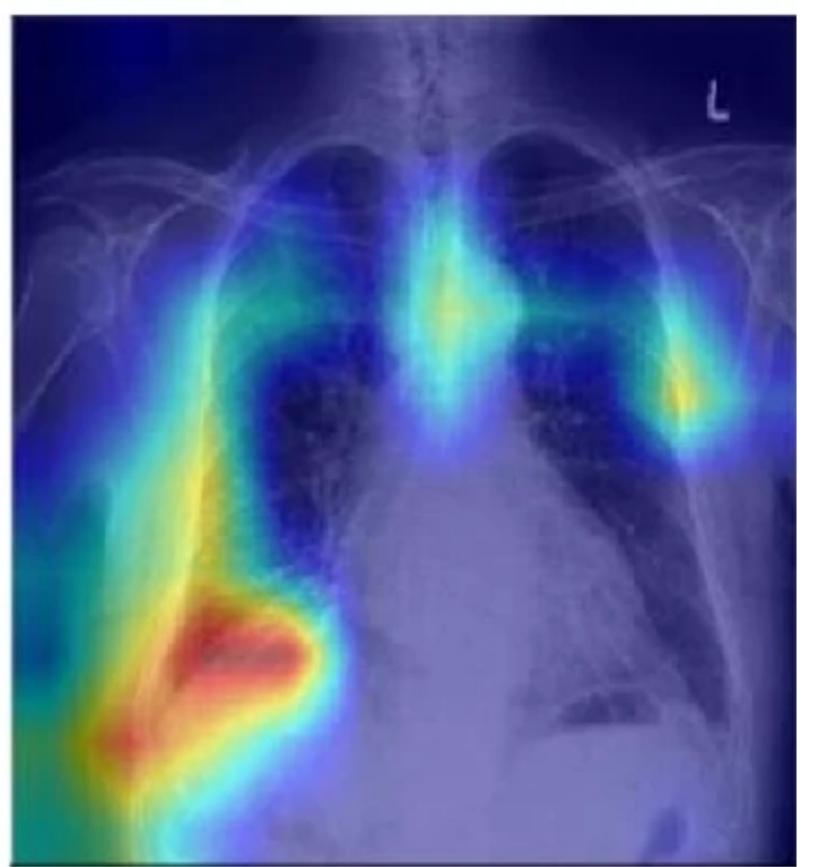
Neural Networks Can Memorize Personal Information From One Example

Anonymisation Fails
Single sample with personal features



DNN trained for
X-ray classification

Inserting the (memorised) unique feature
changes prediction



We Must Incorporate Privacy-Preserving Mechanisms Into RL

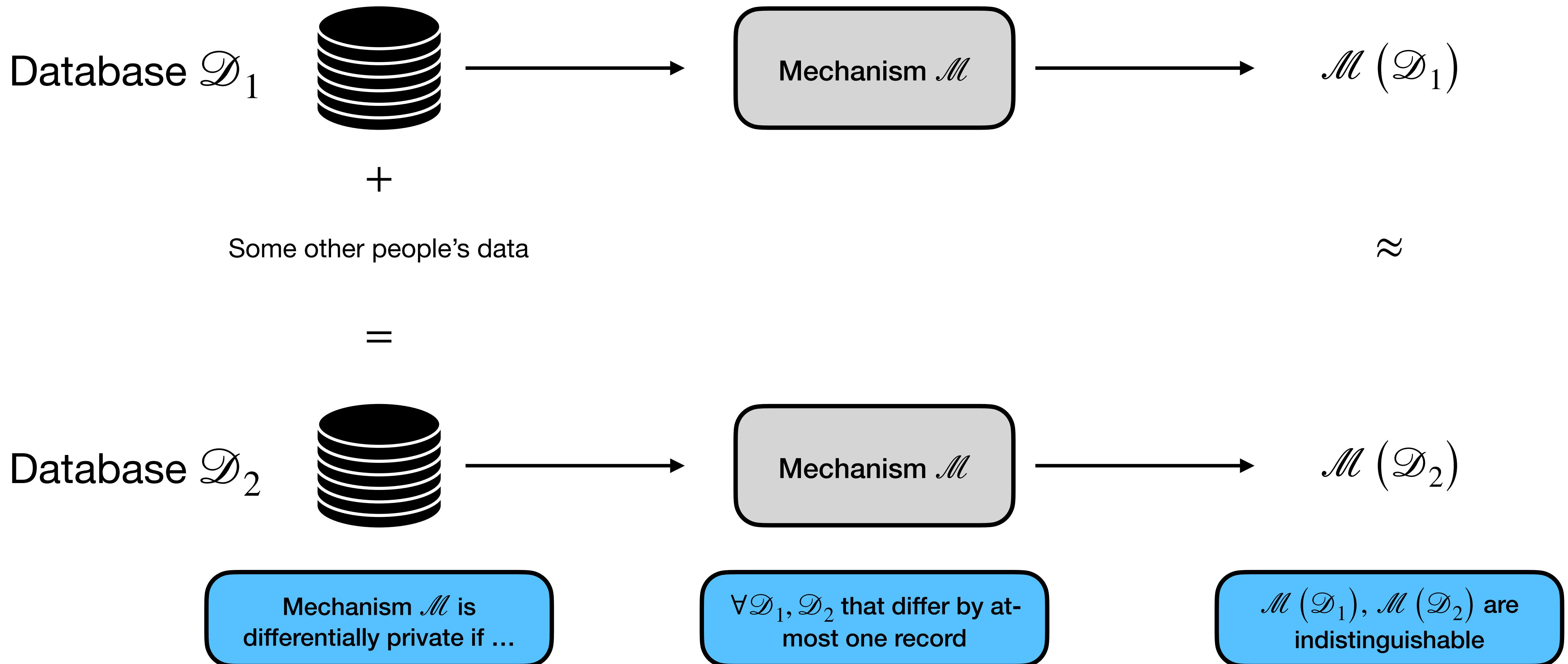
We require a mathematically rigorous framework that provides statistical guarantees for our (possibly randomized) mechanism

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

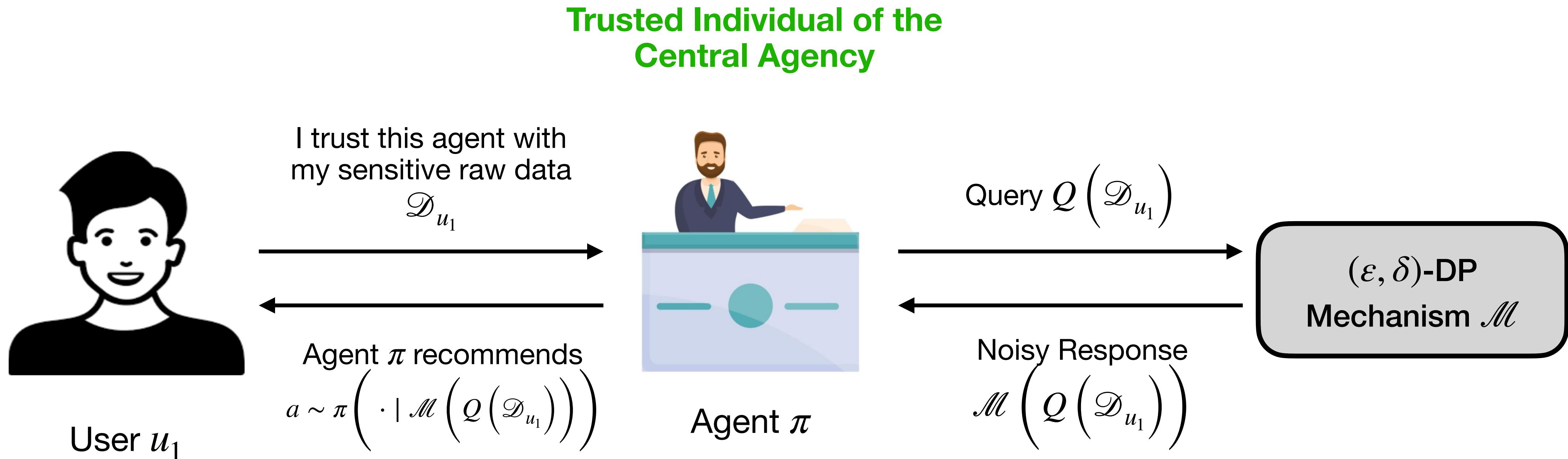
$$\mathbb{P}(\mathcal{M}(\mathcal{U}) \in \mathcal{E}) \leq e^\varepsilon \mathbb{P}(\mathcal{M}(\mathcal{U}') \in \mathcal{E}) + \delta$$

Remark: This is a relaxation of ε -DP as in many settings, achieving ε -DP is near impossible or comes at high utility cost.

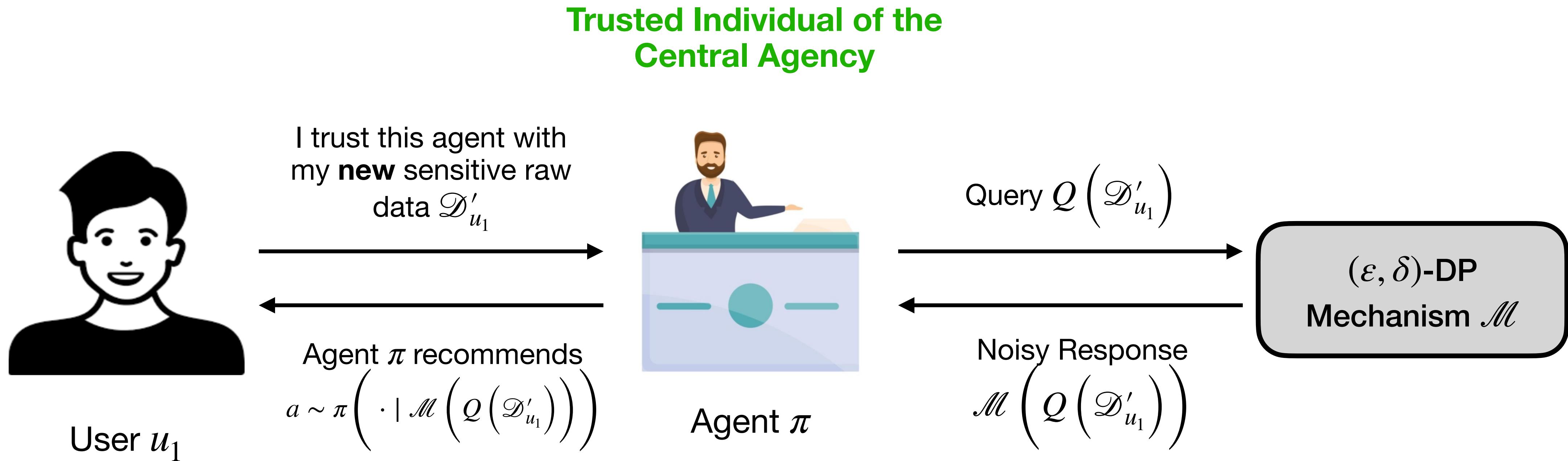
Differential Privacy (DP)



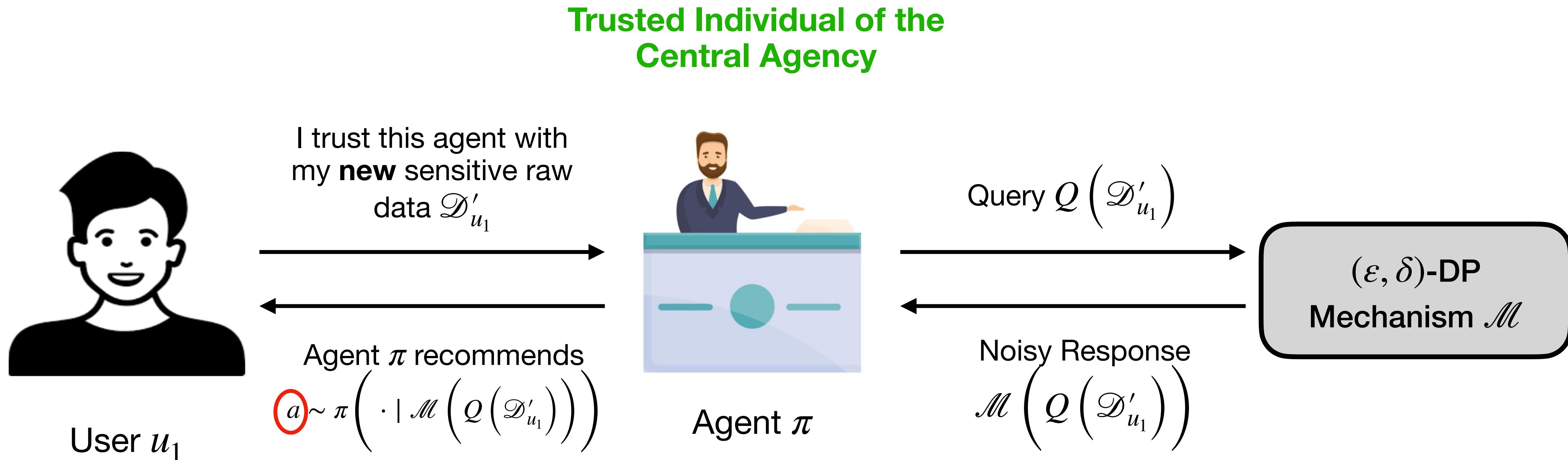
There are a few issues with (ε, δ) -DP



There are a few issues with (ε, δ) -DP

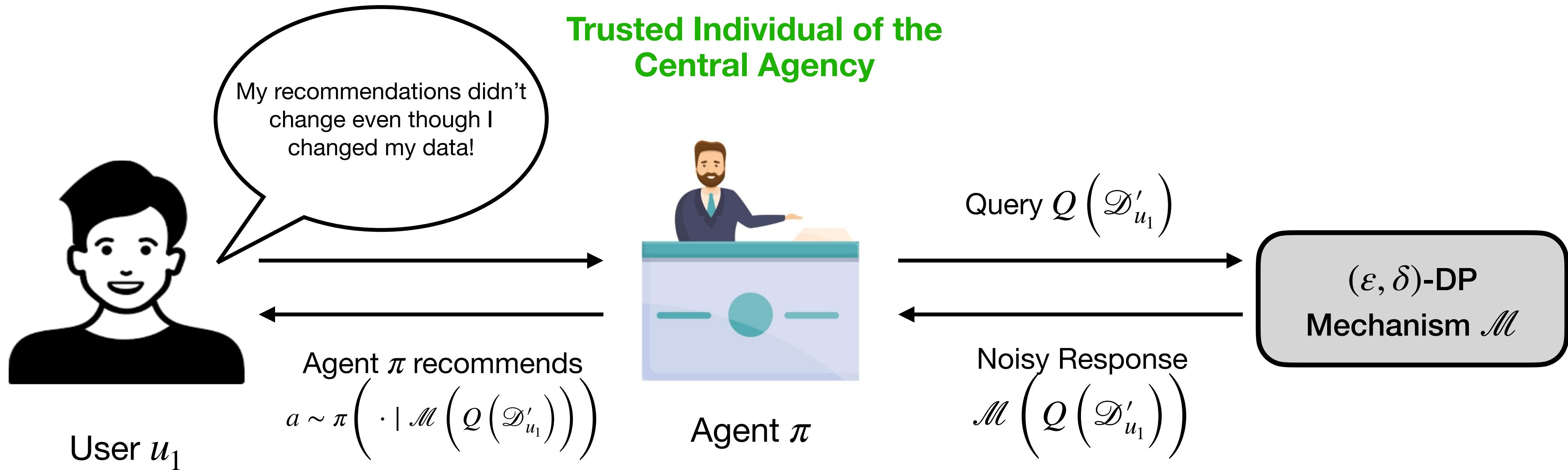


There are a few issues with (ε, δ) -DP

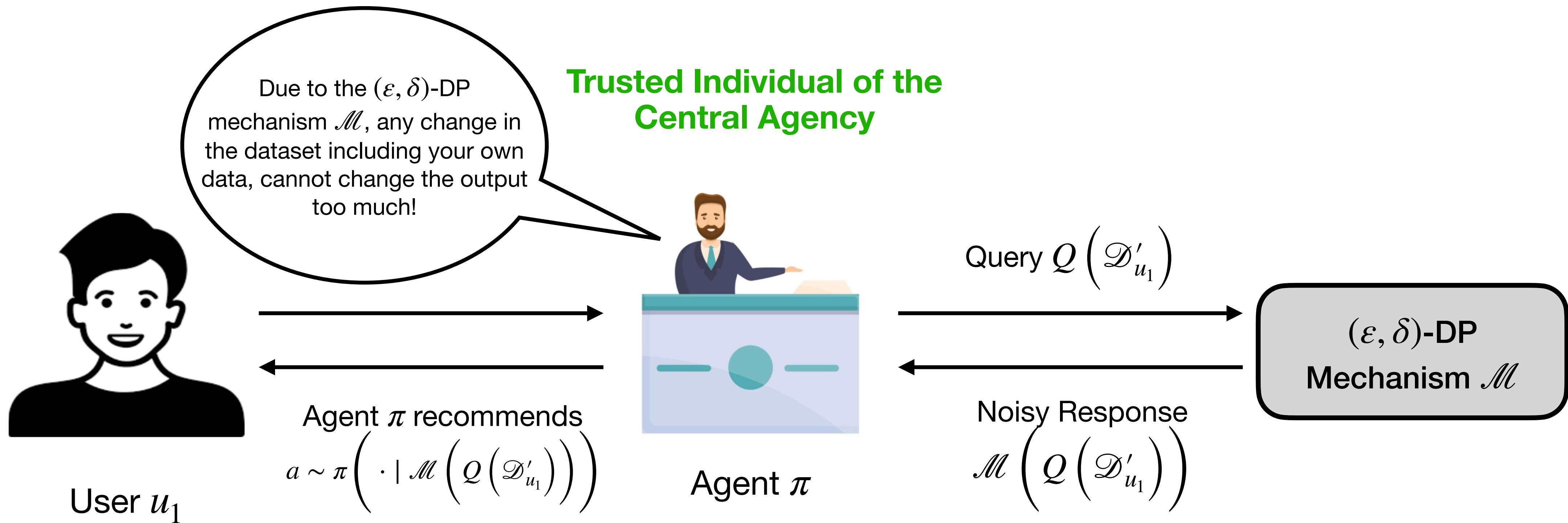


This is the same exact action recommended with the old data!

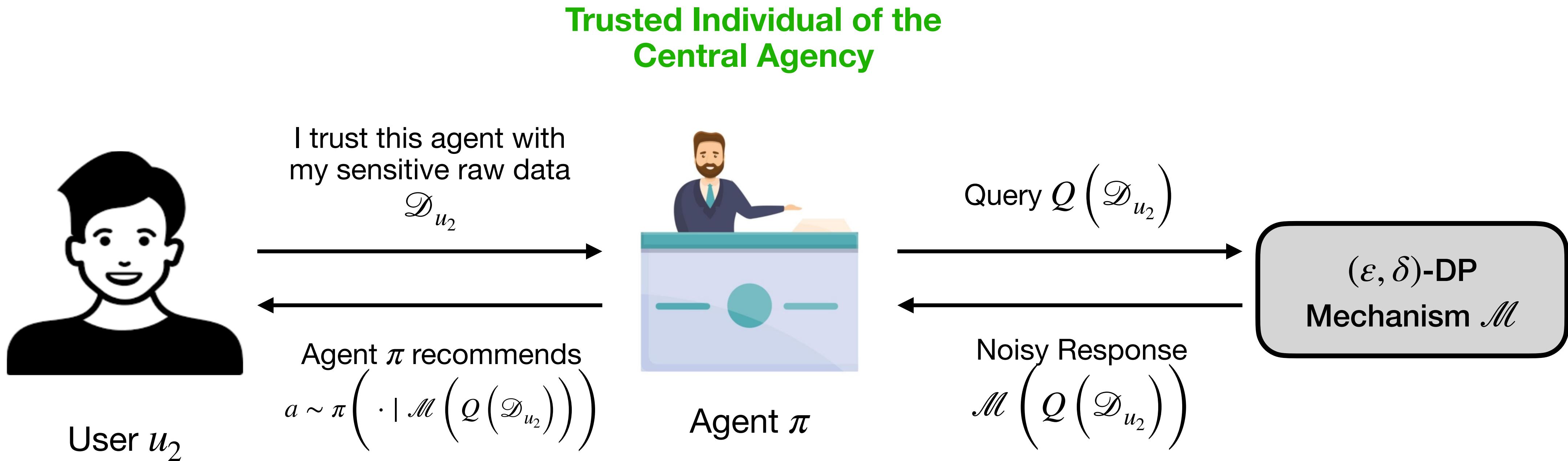
There are a few issues with (ε, δ) -DP



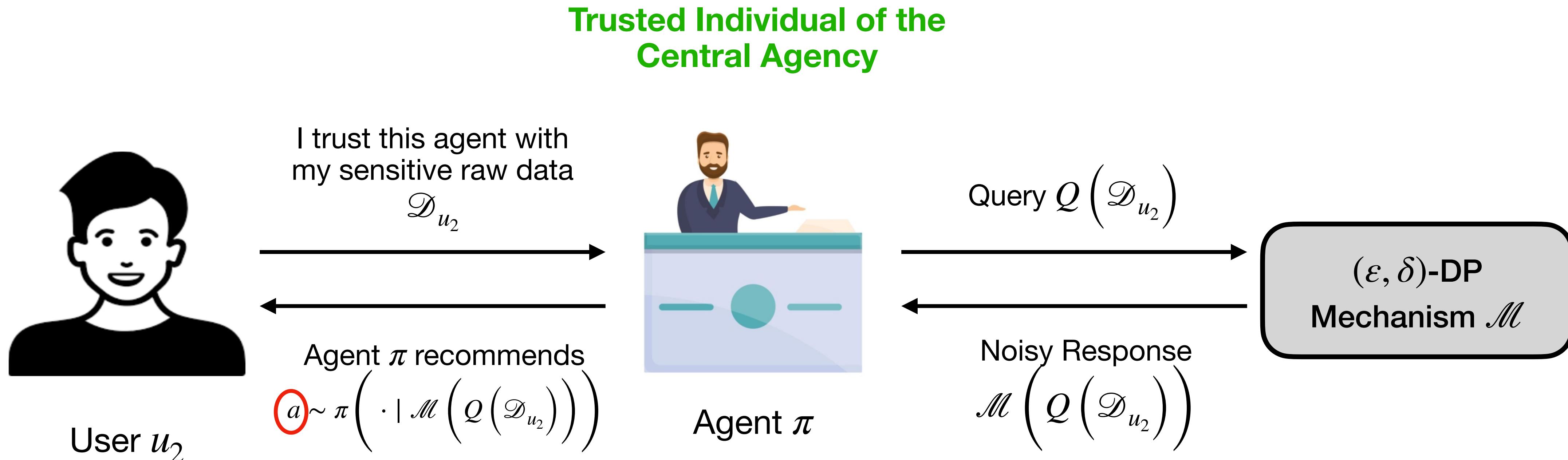
There are a few issues with (ε, δ) -DP



There are a few issues with (ε, δ) -DP

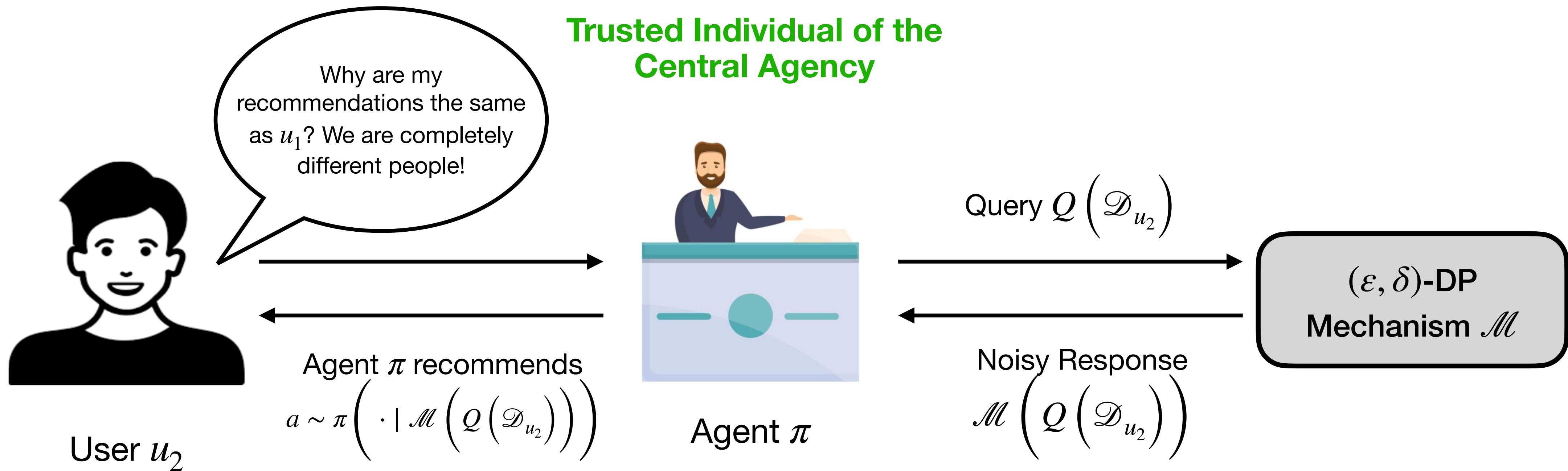


There are a few issues with (ε, δ) -DP

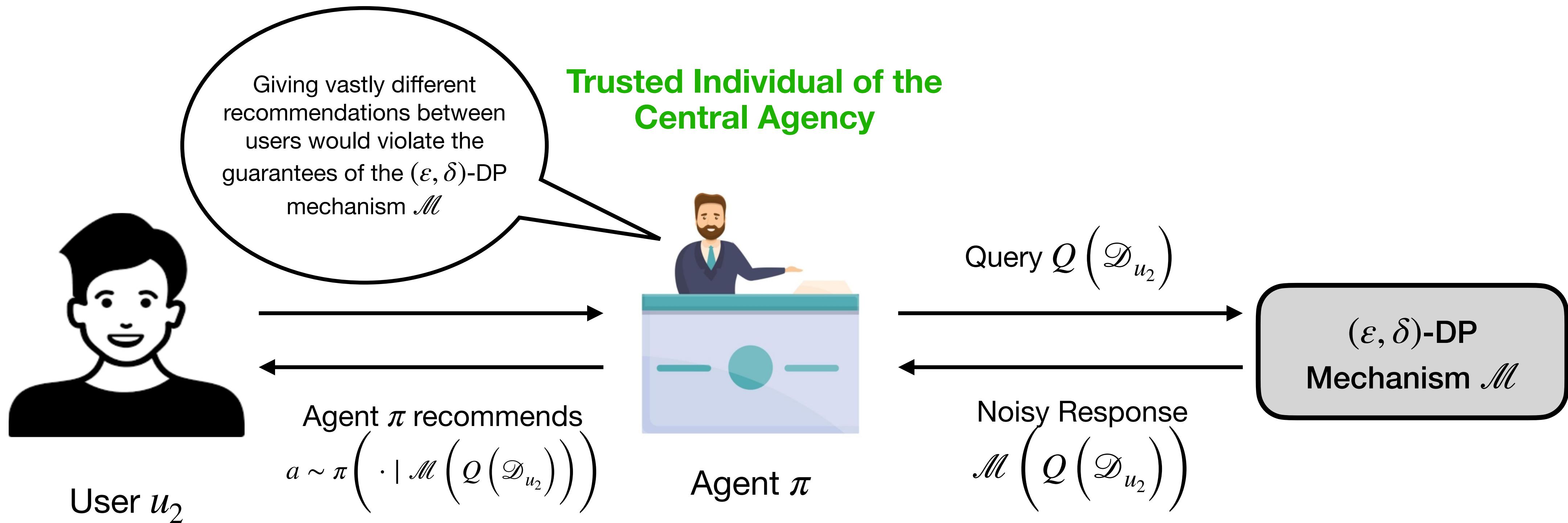


This is the same exact action recommended to the other user!

There are a few issues with (ε, δ) -DP



There are a few issues with (ε, δ) -DP



We need a further relaxation of DP ...

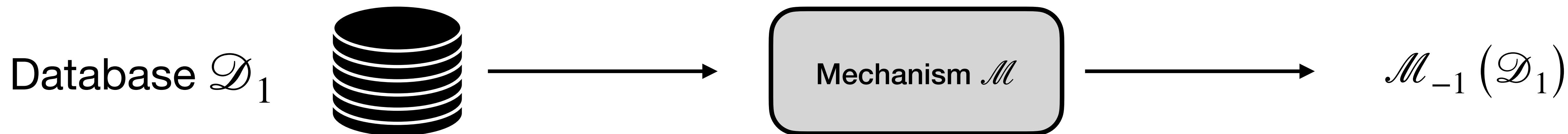
One that works nicely with contextual bandit problems on a per-user level but does not sacrifice privacy on a per-decision or per-context level

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

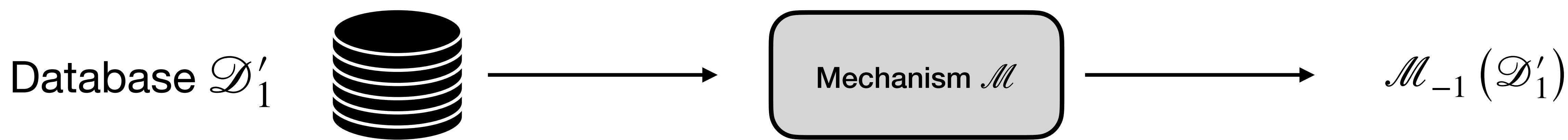
$$\mathbb{P}(\mathcal{M}_{-k}(\mathcal{U}) \in \mathcal{E}) \leq e^\varepsilon \mathbb{P}(\mathcal{M}_{-k}(\mathcal{U}') \in \mathcal{E}) + \delta$$

Remark: JDP allows for **better utility** than standard DP in some contexts since it permits individual outputs to depend more heavily on the individual's own data

Joint Differential Privacy (JDP)



\approx



Mechanism \mathcal{M} is joint differentially private if ...

$\forall \mathcal{D}_1, \mathcal{D}'_1$ where only the data only differs by at most party one's data

$\mathcal{M}^{-1}(\mathcal{D}_1), \mathcal{M}^{-1}(\mathcal{D}'_1)$ are indistinguishable

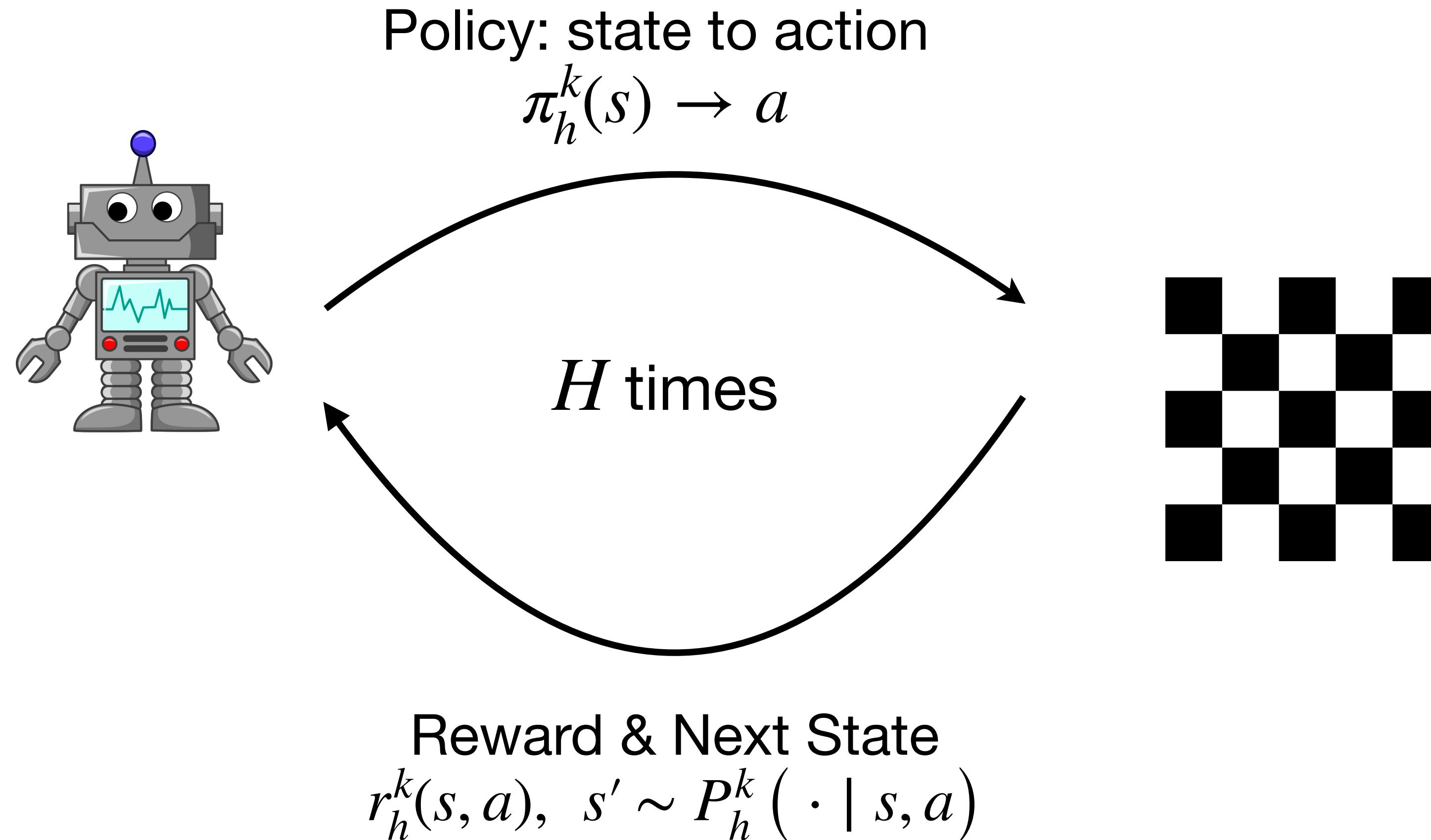
In this talk:

Can we develop an efficient (ε, δ) -JDP algorithm for sequential decision-making problems with **linear parametric representations**, and provide a novel algorithm with provably efficient guarantees for **privacy-preserving exploration**?

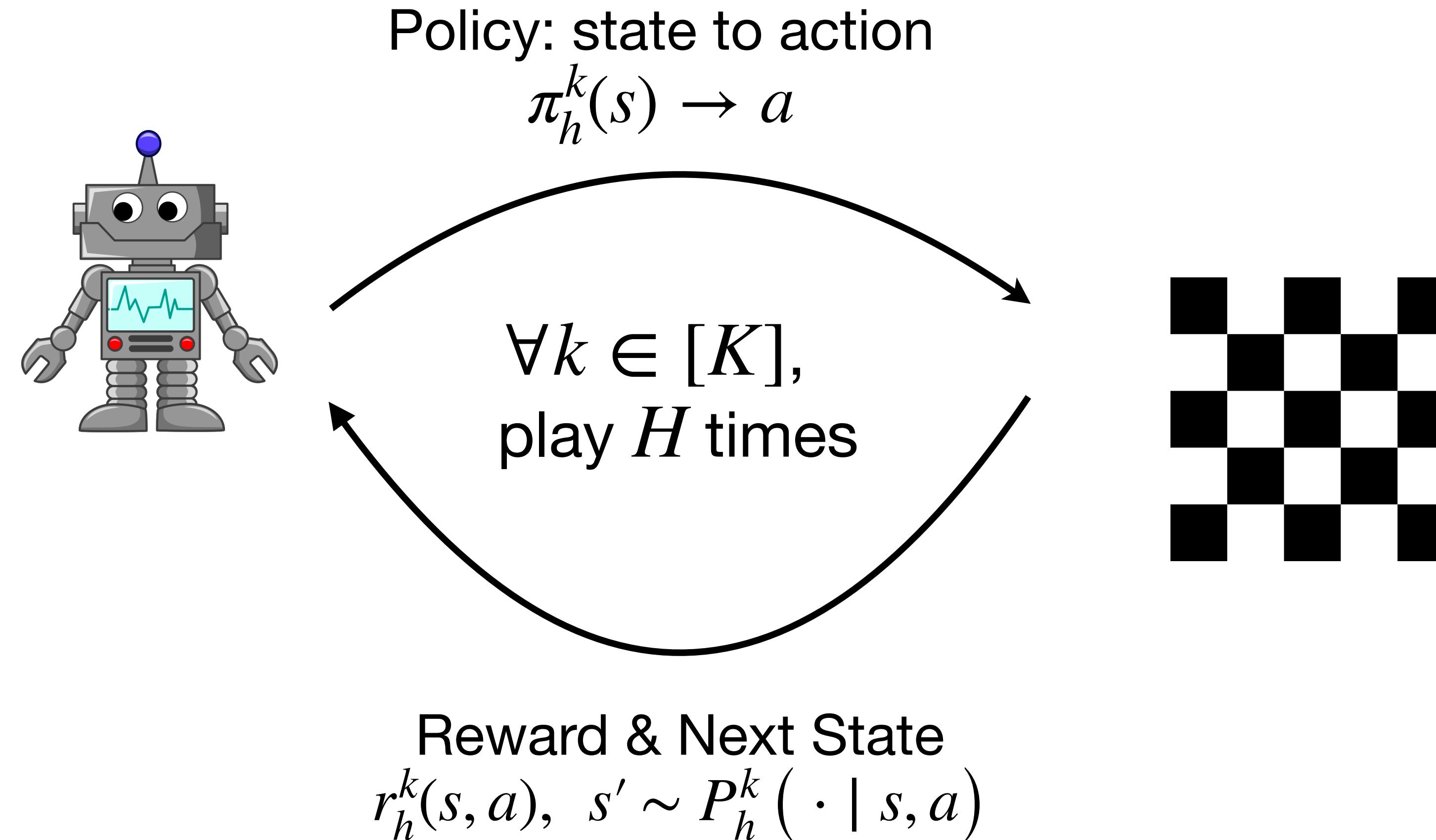
Outline

1. Problem Setup & Previous Work and Motivation
2. Can we do better?
3. Our regret bound with proof sketch

Episodic Time-Inhomogeneous Finite-Horizon MDPs

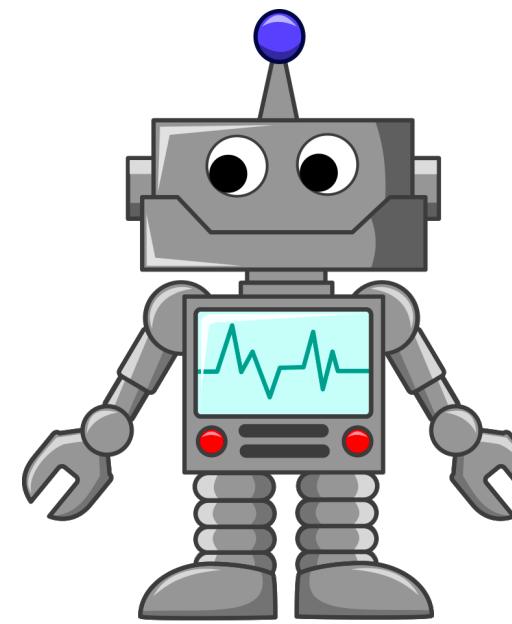


Episodic Time-Inhomogeneous Finite-Horizon MDPs



Episodic Time-Inhomogeneous Finite-Horizon MDPs

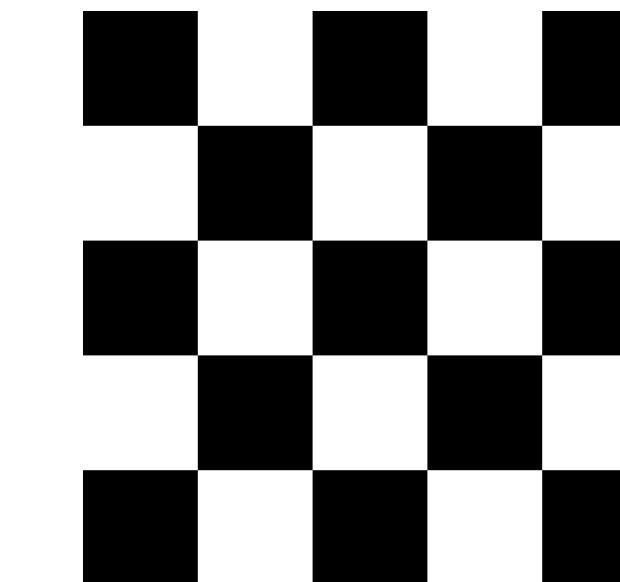
$$\tau^k = \{s_h^k, a_h^k\}_{h=1}^H$$



Policy: state to action

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

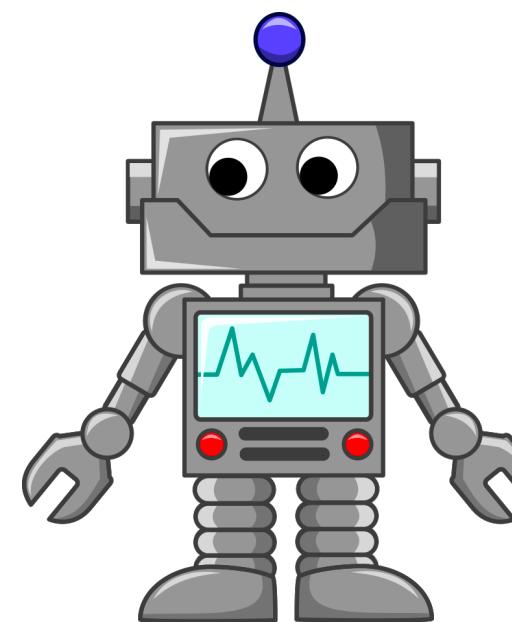
$\forall k \in [K]$,
play H times



Reward & Next State
 $r_h^k(s, a), s' \sim P_h^k(\cdot | s, a)$

Episodic Time-Inhomogeneous Finite-Horizon MDPs

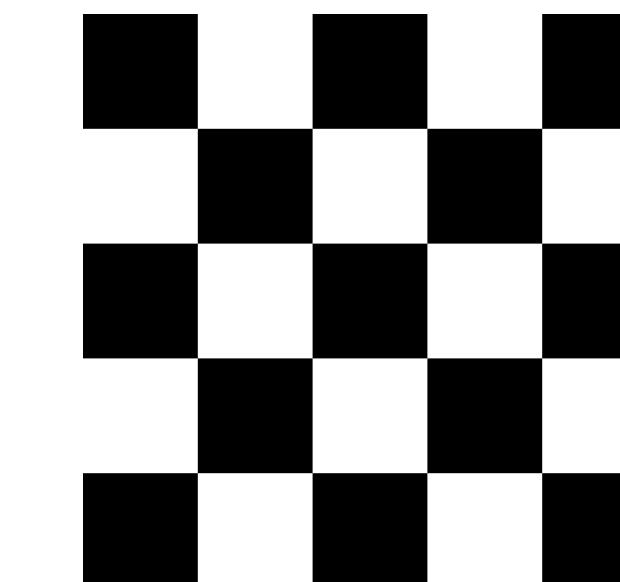
$$\tau^k = \{s_h^k, a_h^k\}_{h=1}^H$$



Policy: state to action

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

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



Reward & Next State
 $r_h^k(s, a), s' \sim P_h^k(\cdot | s, a)$

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

Formal Reinforcement Learning Problem Setting

Let $\mathcal{M} = \left\{ \mathcal{S}, \mathcal{A}, \left\{ r_h \right\}_{h=1}^H, \left\{ \mathcal{P}_h \right\}_h^H, H \right\}$ be an episodic inhomogeneous finite-horizon Markov Decision Process (MDP) where \mathcal{S}, \mathcal{A} are the states and actions, respectively, and $H \in \mathbb{Z}$ is the length of each episode. We call $\mathcal{P}_h : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ the state-transition probability and $r_h : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ the reward function.

Formal Reinforcement Learning Problem Setting

Let $\mathcal{M} = \left\{ \mathcal{S}, \mathcal{A}, \left\{ r_h \right\}_{h=1}^H, \left\{ \mathcal{P}_h \right\}_h^H, H \right\}$ be an episodic inhomogeneous finite-horizon Markov Decision Process (MDP) where \mathcal{S}, \mathcal{A} are the states and actions, respectively, and $H \in \mathbb{Z}$ is the length of each episode. We call $\mathcal{P}_h : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ the state-transition probability and $r_h : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ the reward function.

$$V_h^\pi(s) = \mathbb{E} \left[\sum_{t=h}^H r_t(s_t, a_t) \mid s_h = s, a_t \sim \pi_t(s_t) \right]$$

Value Function (State-value)

$$Q_h^\pi(s, a) = \mathbb{E} \left[\sum_{t=h}^H r(s_t, a_t) \mid s_h = s, a_h = a, a_t \sim \pi_t(s_t) \right]$$

Q-function (Action-value)

Formal Reinforcement Learning Problem Setting

Useful Identities For Later (Bellman Equations)

Let $\mathcal{M} = \{$
(MDP) where
call \mathcal{P}

$$\begin{aligned} Q^\pi(s, a) &= \mathbb{E} \left[\sum_{h=0}^H r(s_h, a_h) \mid s_0 = s, a_0 = a, a_h \sim \pi(\cdot \mid s_h) \right] \\ &= r(s_0, a_0) + \sum_{(s', a') \in \mathcal{S} \times \mathcal{A}} \mathcal{P}(s' \mid s_0, a_0) \pi(a' \mid s') r(s', a') \\ &= r(s_0, a_0) + \mathbb{E}_{s' \sim \mathcal{P}(\cdot \mid s_0, a_0)} V^\pi(s') \end{aligned}$$

$$V^*(s) = \max_{\pi \in \Pi} V^\pi(s) = \max_{a \in \mathcal{A}} Q^*(s, a)$$

ess
We

ate-value)

on (Action-value)

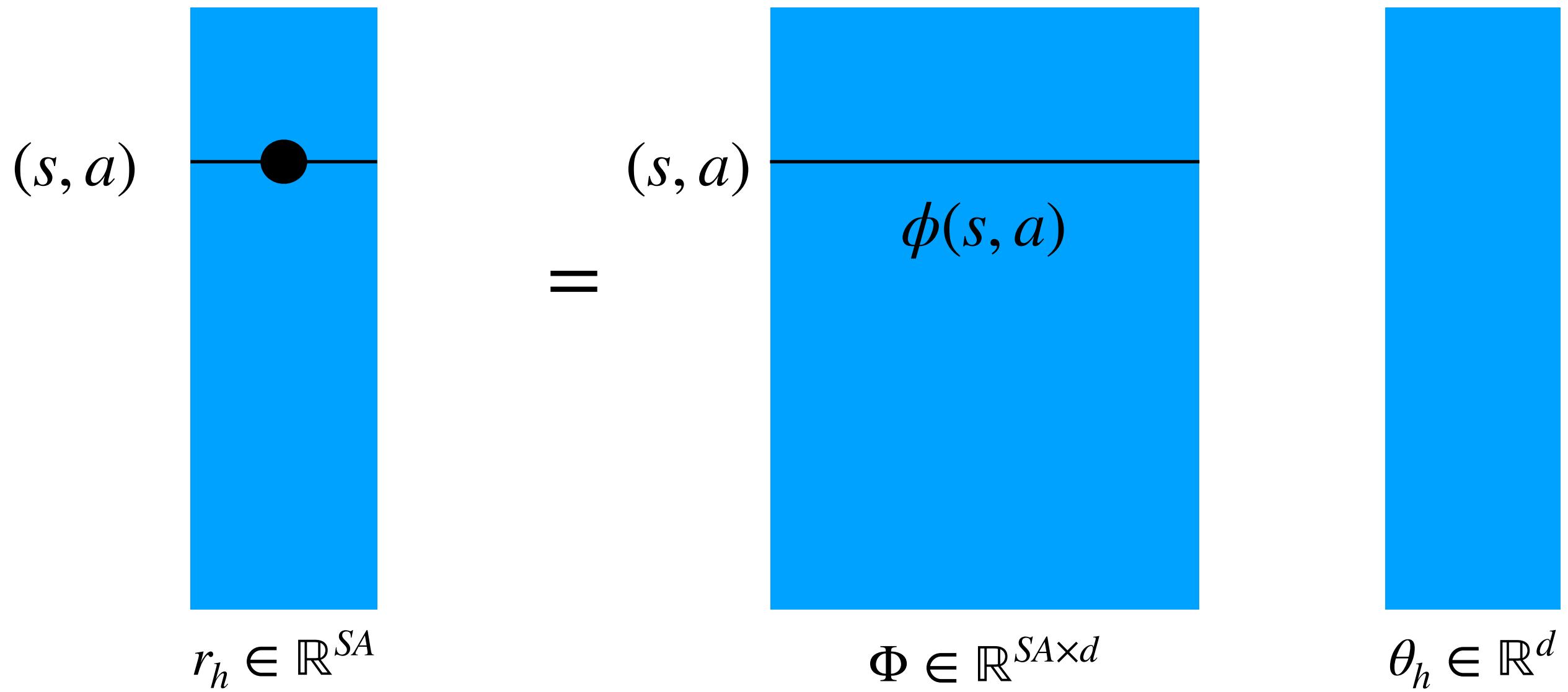
Formal Reinforcement Learning Problem Setting

Let $\mathcal{M} = \left\{ \mathcal{S}, \mathcal{A}, \left\{ r_h \right\}_{h=1}^H, \left\{ \mathcal{P}_h \right\}_h^H, H \right\}$ be an episodic inhomogeneous finite-horizon Markov Decision Process (MDP) where where \mathcal{S}, \mathcal{A} are the states and actions, respectively, and $H \in \mathbb{Z}$ is the length of each episode. We call $\mathcal{P}_h : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ the state-transition probability and $r_h : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ the reward function.

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

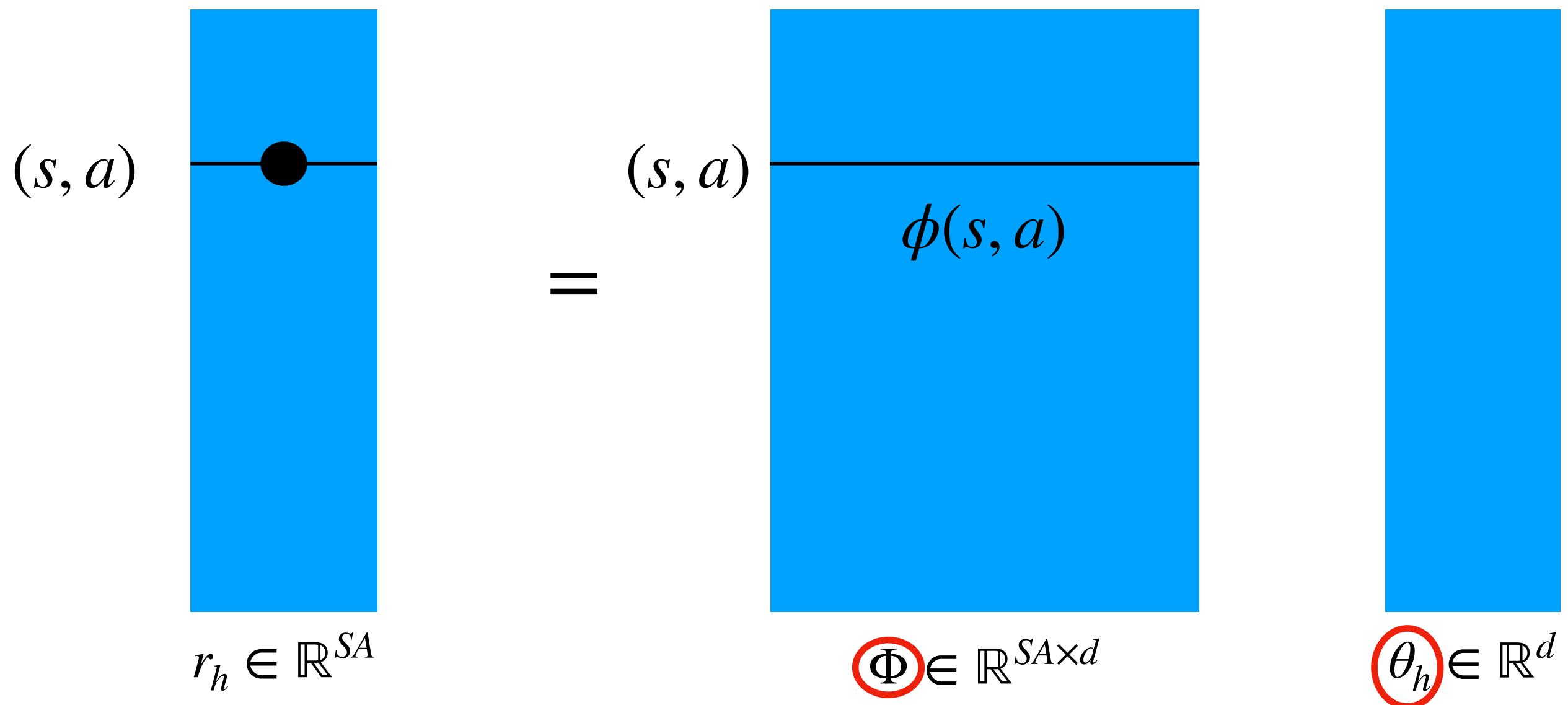
Regret

Linear MDP



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

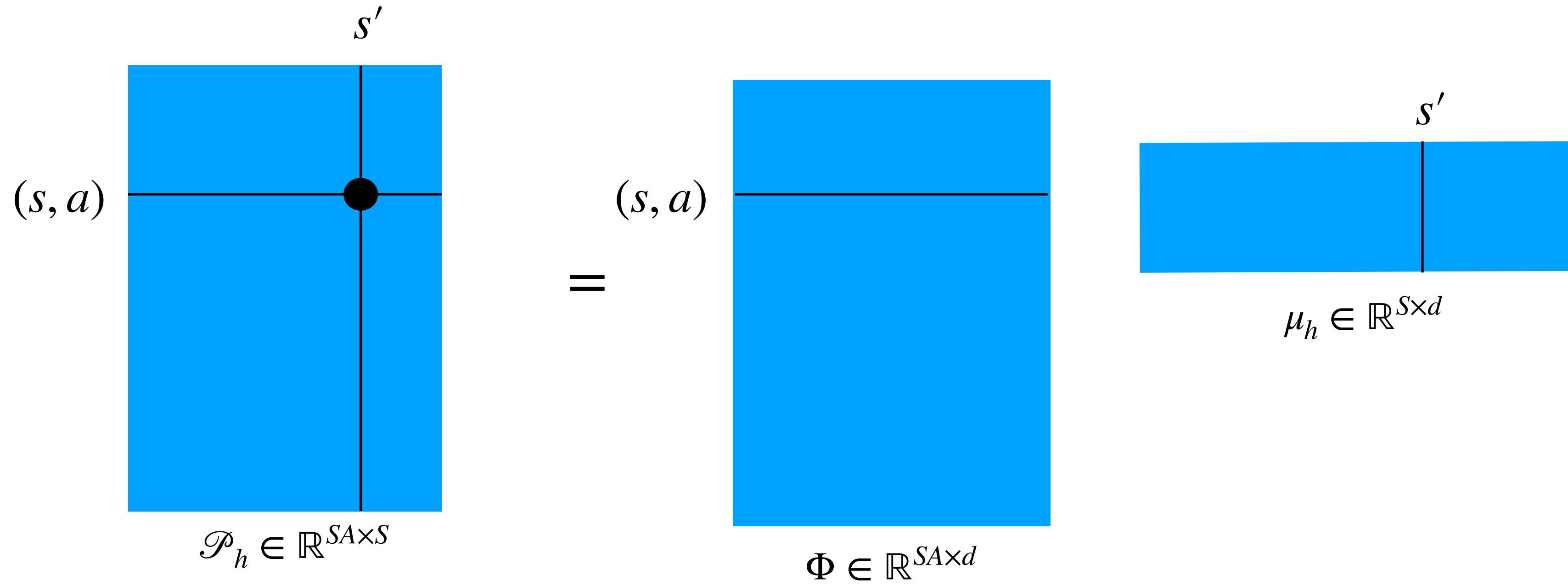
Linear MDP



**These quantities
are known**

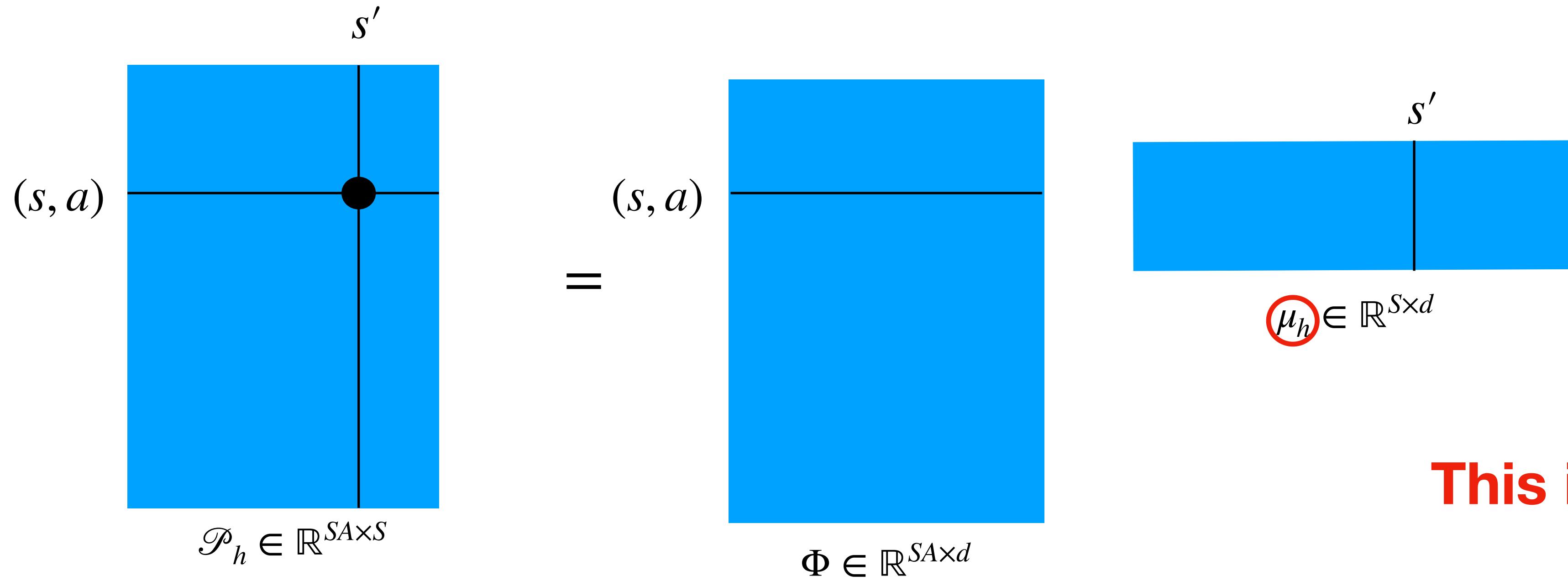
$$\exists \theta_h, \phi : \forall s, a, r_h(s, a) = \theta_h^\top \phi(s, a)$$

Linear MDP



$$\exists \mu_h, \phi : \forall s, a, h, s', \mathcal{P}_h(s' | s, a) = \mu_h(s')^\top \phi(s, a)$$

Linear MDP



$$\exists \mu_h, \phi : \forall s, a, h, s', \mathcal{P}_h(s' | s, a) = \mu_h(s')^\top \phi(s, a)$$

Linear MDP \implies Action-Value function is linear

$$\begin{aligned} Q_h(s, a) &= r_h(s, a) + \mathbb{E}_{s' \sim \mathcal{P}_h(\cdot | s, a)} V_{h+1}(s') \\ &= \theta_h^\top \phi(s, a) + \mathcal{P}_h(s, a)^\top V_{h+1} \\ &= \theta_h^\top \phi(s, a) + (\mu_h V_{h+1})^\top \phi(s, a) \\ &= (\theta_h + \mu_h V_{h+1})^\top \phi(s, a) \\ &= w_h^\top \phi(s, a) \end{aligned}$$

Linear MDP \implies Action-Value function is linear

$$\begin{aligned} Q_h(s, a) &= r_h(s, a) + \mathbb{E}_{s' \sim \mathcal{P}_h(\cdot | s, a)} V_{h+1}(s') \\ &= \theta_h^\top \phi(s, a) + \mathcal{P}_h(s, a)^\top V_{h+1} \\ &= \theta_h^\top \phi(s, a) + (\mu_h V_{h+1})^\top \phi(s, a) \\ &= (\theta_h + \mu_h V_{h+1})^\top \phi(s, a) \\ &= \textcolor{red}{w_h^\top} \phi(s, a) \end{aligned}$$

If we learn this,
we can estimate
Q and thus the
optimal policy!

Learning The Transition Dynamics With Ridge Regression

At each time step h , we try to solve

$$\hat{w}_h = \operatorname{argmin}_w \sum_{i=1}^K \left(w^\top \phi(s_h^{(i)}, a_h^{(i)}) - y_h^{(i)} \right)^2 + \lambda \|w\|_2^2$$

With $\lambda > 0$ and the target labels being:

$$y_h^{(i)} = r_h^{(i)} + \max_{a' \in \mathcal{A}} Q_{h+1} \left(s_{h+1}^{(i)}, a' \right)$$

Learning The Transition Dynamics With Ridge Regression

From the solution of ridge regression, we find that

$$w_h^k = \Lambda_h^{-1} \sum_{i=1}^{k-1} \phi(s_h^i, a_h^i) \left(r_h(s_h^i, a_h^i) + \max_{a' \in \mathcal{A}} Q_h(s_h^i, a') \right)$$

$$\text{where } \Lambda_h = \sum_{i=1}^{k-1} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^\top + \lambda I$$

LSVI-UCB

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

- 1: **for** episode $k = 1, \dots, K$ **do**
- 2: Receive the initial state x_1^k .
- 3: **for** step $h = H, \dots, 1$ **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 $h = 1, \dots, H$ **do**
- 8: Take action $a_h^k \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} Q_h(x_h^k, a)$, and observe x_{h+1}^k .

LSVI-UCB

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

- 1: **for** episode $k = 1, \dots, K$ **do**
- 2: Receive the initial state x_1^k .
- 3: **for** step $h = H, \dots, 1$ **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 $h = 1, \dots, H$ **do**
- 8: Take action $a_h^k \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} Q_h(x_h^k, a)$, and observe x_{h+1}^k .

LSVI-UCB Bonus

The agent is learning from limited data. Like a regression confidence interval, we want to hedge against uncertainty in our estimate of \hat{w}_h . For any new (s, a) , the uncertainty in prediction is proportional to

$$\beta \|\phi(s, a)\|_{\Lambda_h^{-1}} = \beta \sqrt{\phi(s, a)^\top \Lambda_h^{-1} \phi(s, a)}$$

where $\beta = dH \sqrt{d \log \left(\frac{dKH}{\delta} \right)}$

LSVI-UCB Bonus

The agent is learning from limited data. Like a regression confidence interval, we want to hedge against uncertainty in our estimate of \hat{w}_h . For any new (s, a) , the uncertainty in prediction is proportional to

$$\beta \|\phi(s, a)\|_{\Lambda_h^{-1}} = \beta \sqrt{\phi(s, a)^\top \Lambda_h^{-1} \phi(s, a)}$$

$$\text{where } \beta = dH \sqrt{d \log \left(\frac{dKH}{\delta} \right)}$$

Intuition: Another way to think about this is that this is a carefully curated bonus given to our agent that promotes exploration by taking actions that are less certain. It ensures that with high probability $Q_h^k(s, a)$ is an upper confidence bound of the true Q function $Q_h^*(s, a) \quad \forall (s, a)$

LSVI-UCB Bonus

The agent is learning from limited data. Like a regression confidence interval, we want to hedge against uncertainty in our estimate of \hat{w}_h . For any new (s, a) , the uncertainty in prediction is proportional to

$$\beta \|\phi(s, a)\|_{\Lambda_h^{-1}} = \beta \sqrt{\phi(s, a)^\top \Lambda_h^{-1} \phi(s, a)}$$

where $\beta = dH \sqrt{d \log \left(\frac{dKH}{\delta} \right)}$

**How do we get
this?**

Intuition: Another way to think about this is that this is a carefully curated bonus given to our agent that promotes exploration by taking actions that are less certain. It ensures that with high probability $Q_h^k(s, a)$ is an upper confidence bound of the true Q function $Q_h^*(s, a) \quad \forall (s, a)$

LSVI-UCB Bonus Proof Sketch

Lemma A (Theorem 1 in Abbasi-Yadkori et al. 2011). Let $\{\mathcal{F}_t\}_{t=0}^\infty$ be a filtration. Let $\{\eta_t\}_{t=1}^\infty$ be a real-valued stochastic process such that η_t is \mathcal{F}_t measurable and η_t is conditionally R -sub-Gaussian for some $R \geq 0$ i.e.

$$\forall \lambda \in \mathbb{R}, \quad \mathbb{E}[e^{\lambda \eta_t} | \mathcal{F}_{t-1}] \leq \exp\left(\frac{\lambda^2 R^2}{2}\right)$$

Let $\{x_t\}_{t=1}^\infty$ be an \mathbb{R}^d -valued stochastic process such that x_t is \mathcal{F}_{t-1} -measurable. Assume that Z is a $d \times d$ positive definite matrix. For any $k \geq 0$, define

$$Z_k = Z + \sum_{s=1}^t X_s X_s^\top$$

Then, for any $\delta > 0$, with probability at least $1 - \delta$, for all $t \geq 0$

$$\left\| \sum_{i=1}^k x_i \eta_i \right\|_{Z_k^{-1}}^2 \leq 2R^2 \log\left(\frac{\det(Z_k)^{1/2} \det(Z)^{-1/2}}{\delta}\right)$$

This is a self-normalizing martingale bound using **Azuma-Hoeffding**

LSVI-UCB Bonus Proof Sketch

Lemma B (Lemma D.6 of Jin et al. 2020). Let \mathcal{V} denote a class of functions mapping from \mathcal{S} to \mathbb{R} with the following parametric form

$$V(\cdot) = \min \left\{ \max_a \left[w^\top \phi(\cdot, a) + \beta \sqrt{\phi(\cdot, a)^\top \Lambda^{-1} \phi(\cdot, a)} \right], H \right\}$$

where the parameters (w, β, Λ) satisfy $\|w\| \leq L$, $\beta \in [0, B]$, and the minimum eigenvalue satisfies $\lambda_{\min}(\Lambda) \geq \lambda$. Assume $\|\phi(s, a)\| \leq 1$ for all (s, a) pairs, and let \mathcal{N}_ε be the ε -covering number of \mathcal{V} with respect to distance

$$\text{dist}(V, V') = \sup_s |V(s) - V'(s)|$$

Then,

$$\log \mathcal{N}_\varepsilon \leq d \log \left(1 + \frac{4L}{\varepsilon} \right) + d^2 \log \left(1 + \frac{8\sqrt{dB^2}}{\lambda \varepsilon^2} \right).$$

LSVI-UCB Regret

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

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

$$\leq \sum_{k=1}^K \sum_{h=1}^H \zeta_k^h + 2\beta \sum_{k=1}^K \sum_{h=1}^H \sqrt{(\phi_k^h)^\top (\Lambda_k^h)^{-1} \phi_k^h}$$

$$\underset{\sim}{<} \widetilde{\mathcal{O}}(d^3 H^4 K)$$

Privacy Concerns In LSVI-UCB

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

- 1: **for** episode $k = 1, \dots, K$ **do**
- 2: Receive the initial state x_1^k .
- 3: **for** step $h = H, \dots, 1$ **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 $h = 1, \dots, H$ **do**
- 8: Take action $a_h^k \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} Q_h(x_h^k, a)$, and observe x_{h+1}^k .

Privacy Concerns In LSVI-UCB

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

```
1: for episode  $k = 1, \dots, K$  do
2:   Receive the initial state  $x_1^k$ .
3:   for step  $h = H, \dots, 1$  do
4:      $\Lambda_h \leftarrow \sum_{\tau=1}^{k-1} \phi(x_h^\tau, a_h^\tau) \phi(x_h^\tau, a_h^\tau)^\top + \lambda \cdot \mathbf{I}$ . Think of  $\phi(s, a)$  as a one-hot vector, then  $\Lambda_h$  is capturing something similar to visitation counts which uses trajectory information with possibly private data
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  $h = 1, \dots, H$  do
8:     Take action  $a_h^k \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} Q_h(x_h^k, a)$ , and observe  $x_{h+1}^k$ .
```

Privacy Concerns In LSVI-UCB

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

```
1: for episode  $k = 1, \dots, K$  do
2:   Receive the initial state  $x_1^k$ .
3:   for step  $h = H, \dots, 1$  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  $h = 1, \dots, H$  do
8:     Take action  $a_h^k \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} Q_h(x_h^k, a)$ , 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. These can also leak information about trajectories taken by the policy

\mathbf{w}_h

We need to privatize these terms!

Differential Privacy Tools

Lemma C (Lemma 1.7 of Bun and Steiner. 2016). Let $\mathcal{A} : \mathcal{X}^n \rightarrow \mathcal{Y}$ and $\mathcal{A}' : \mathcal{X}^n \times \mathcal{Y} \rightarrow \mathcal{Z}$ be (possibly randomized) mechanisms. For $\delta > 0$, suppose \mathcal{A} satisfies (ϵ_1, δ) -DP, and for each $y \in \mathcal{Y}$, $\mathcal{A}'(\cdot, y)$ satisfies (ϵ_2, δ) -DP. Define the composed mechanism

$$\mathcal{A}''(x) = \mathcal{A}'(x, \mathcal{A}(x))$$

Then, \mathcal{A}'' satisfies $(\epsilon_1 + \epsilon_2, 2\delta)$ -DP.

Lemma D (Theorem 3.22 of Dwork and Roth. 2014). Let $f : \mathbb{N}^{\mathcal{X}} \rightarrow \mathbb{R}^d$ be an arbitrary d-dimensional function with

$$\Delta(f) = \max_{\mathcal{U} \sim \mathcal{U}'} ||f(\mathcal{U}) - f(\mathcal{U}')||_2$$

where $\mathcal{U} \sim \mathcal{U}'$ are neighboring datasets. The Gaussian mechanism $\mathcal{M}_{\text{Gauss}}$ with noise level σ is given by

$$\mathcal{M}_{\text{Gauss}}(\mathcal{U}) = f(\mathcal{U}) + \mathcal{N}(0, \sigma^2 I_d)$$

For all $0 < \delta, \epsilon < 1$, a Gaussian Mechanism with noise parameter $\sigma = \frac{\Delta}{\epsilon} \sqrt{2 \log(1.25/\delta)}$ satisfies (ϵ, δ) -DP.

Differential Privacy Tools

Lemma E (Billboard Lemma of Hsu et al. 2013). Suppose that a randomized mechanism $\mathcal{A} : \mathcal{X}^n \rightarrow \mathcal{Y}$ is (ϵ, δ) -DP. Let $U \in \mathcal{U}$ be a dataset containing n users. Then, consider any set of functions $f_i : \mathcal{U}_i \times \mathcal{Y} \rightarrow \mathcal{Y}_i$ for $i \in [n]$ where \mathcal{U}_i is the portion of the dataset containing user i 's data. Then, the composition $\left\{ f_i(\Pi_i(U), \mathcal{A}(U)) \right\}_{i \in [n]}$ is (ϵ, δ) -JDP where $\Pi_i : \mathcal{U} \rightarrow \mathcal{U}_i$ is the canonical projection to the i -th user's data

Remark: This lemma tells us that if we construct an (ϵ, δ) -DP algorithm and we have a function f_i that operates on user i 's data, then that mechanism is indeed (ϵ, δ) .

Previous Work

[Theorem 8 of Luyo et al. 2021]. Fix any privacy level $\varepsilon, \delta \in (0,1)$. For any $p \in (0,1)$, their algorithm is (ε, δ) -JDP and, with probability at least $1 - p$, its regret is bounded as follows:

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

Approach: Privatize Λ_h, w_h

$$\tilde{\Lambda}_h = \Lambda_h + \mathcal{N} \left(0, \mathcal{O} \left(K^{1/5} d^{3/10} H^{2/5} \varepsilon^{-4/5} \log(1/\delta) \right) \right)$$

$$\tilde{w}_{b+1,h} = \tilde{\Lambda}_{b+1,h}^{-1} \sum_{i=1}^{b+1} \phi(s_{i,h}, a_{i,h}) [r_h(s_{i,h}, a_{i,h}) + V_{b+1,h+1}(s_{i,h+1})] + \mathcal{N} \left(0, \tilde{\Lambda}_{b+1,h}^{-1} \cdot \mathcal{O} \left(\frac{1}{\varepsilon} H^2 B \log \left(\frac{1}{\delta} \right) \right) \cdot \tilde{\Lambda}_{b+1,h}^{-1} \right)$$

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

Previous Work

[Theorem 16 of Ngo et al. 2022]. Fix any privacy level $\varepsilon, \delta \in (0,1)$. For any $p \in (0,1)$, their algorithm is (ε, δ) -JDP and, with probability at least $1 - p$, its regret is bounded as follows:

$$\mathcal{R}(K) \leq \widetilde{\mathcal{O}} \left(\sqrt{d^3 H^4 K} + H^3 d^{5/4} K^{1/2} / \varepsilon^{1/2} \right)$$

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

Motivating Work: LSVI-UCB++

[**Theorem 5.1 of He et al. 2023**]. For any linear MDP \mathcal{M} with K sufficiently large, if we set the parameters $\lambda = 1/H^2$ and the confidence intervals $\beta, \bar{\beta}, \tilde{\beta}$ as

$$\begin{aligned}\beta &= \mathcal{O} \left(H\sqrt{d\lambda} + \sqrt{d \log^2 (1 + dKH/(\delta\lambda))} \right) \\ \bar{\beta} &= \mathcal{O} \left(H\sqrt{d\lambda} + \sqrt{d^3 H^2 \log^2 (dHK/(\delta\lambda))} \right) \\ \tilde{\beta} &= \mathcal{O} \left(H^2 \sqrt{d\lambda} + \sqrt{d^3 H^4 \log^2 (dHK/(\delta\lambda))} \right)\end{aligned}$$

then with high probability of at least $1 - 7\delta$, the regret of LSVI-UCB++ is upper bounded as follows

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

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

[Theorem 4.1 of [Qiao and Wang. 2023](#)]. For any privacy budget $\epsilon > 0$, 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 (ϵ, δ) -JDP and achieves regret upper bounded by:

$$\mathcal{R}(K) \leq \widetilde{O} \left(\sqrt{H^3 SAK} + S^2 A H^3 / \epsilon \right)$$

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 a sufficient variance component. However, to use a **Bernstein-bound**, we need to **carefully privatize the bounds** to ensure that we can **upper bound the variance term** in a Bernstein-bound

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

$$\mathcal{R}(K) \leq \tilde{\sigma} \left(\sqrt{d^3 H^4 K} + H^3 d^{5/4} K^{1/2} / \epsilon^{1/2} \right)$$

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

Cost of privacy: can
we improve this to
 $\mathcal{O}(\text{poly}(dHK)/\epsilon)$

Yes We Can: DP-LSVI-UCB++

[**Theorem 3.2 of Sahu. 2025**]. For any linear MDP \mathcal{M} with K sufficiently large, if we set the parameters

$$\begin{aligned}\lambda_{\widetilde{\Lambda}} &= \mathcal{O}\left(\sqrt{dHK} \left(2 + \left(\frac{\log(5H/\delta)}{d}\right)^{2/3}\right)\right) \\ L &= \mathcal{O}\left(H\sqrt{dHK \log(dKH/\delta)}\right)\end{aligned}$$

and the confidence intervals $\hat{\beta}, \check{\beta}, \bar{\beta}$ as

$$\begin{aligned}\hat{\beta} = \check{\beta} &= \mathcal{O}\left(HL\sqrt{d\lambda_{\widetilde{\Lambda}}} + \sqrt{d^3H^2 \log^2(dH^3KL^2/(\delta\lambda_{\widetilde{\Lambda}}))}\right) \\ \bar{\beta} &= \mathcal{O}\left(H^2L^2\sqrt{d\lambda_{\widetilde{\Lambda}}} + \sqrt{d^3H^4 \log^2(dH^4KL^2/(\delta\lambda_{\widetilde{\Lambda}}))}\right)\end{aligned}$$

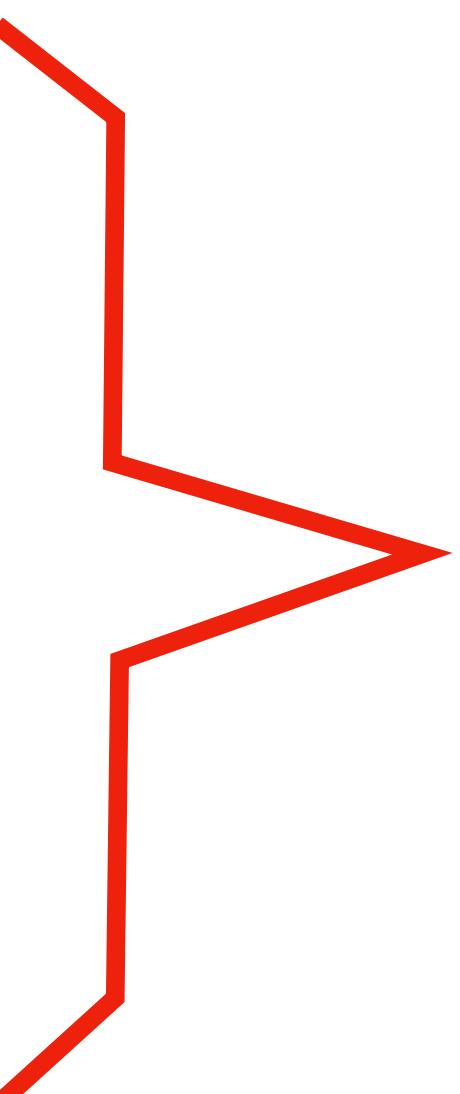
then with high probability of at least $1 - 7\delta$, the regret of DP-LSVI-UCB++ is upper bounded as follows

$$\mathcal{R}(K) \leq \widetilde{\mathcal{O}}\left(d\sqrt{H^3K} + H^{15/4}d^{7/6}K^{1/2}/\epsilon\right)$$

Algorithm 1 LSVI-UCB++

Require: Regularization parameter $\lambda > 0$, confidence radius $\beta, \bar{\beta}, \tilde{\beta}$

- 1: Initialize $k_{\text{last}} = 0$ and for each stage $h \in [H]$ set $\Sigma_{0,h}, \Sigma_{1,h} \leftarrow \lambda I$
- 2: For each stage $h \in [H]$ and state-action $(s, a) \in \mathcal{S} \times \mathcal{A}$, set $Q_{0,h}(s, a) \leftarrow H, \check{Q}_{0,h}(s, a) \leftarrow 0$
- 3: **for** episodes $k = 1, \dots, K$ **do**
- 4: Received the initial state s_1^k .
- 5: **for** stage $h = H, \dots, 1$ **do**
- 6: $\hat{\mathbf{w}}_{k,h} = \Sigma_{k,h}^{-1} \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) V_{k,h+1}(s_{h+1}^i)$
- 7: $\check{\mathbf{w}}_{k,h} = \Sigma_{k,h}^{-1} \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) \check{V}_{k,h+1}(s_{h+1}^i)$
- 8: **if** there exists a stage $h' \in [H]$ such that $\det(\Sigma_{k,h'}) \geq 2 \det(\Sigma_{k_{\text{last}}, h'})$ **then**
- 9: $Q_{k,h}(s, a) = \min \left\{ r_h(s, a) + \hat{\mathbf{w}}_{k,h}^\top \phi(s, a) + \beta \sqrt{\phi(s, a)^\top \Sigma_{k,h}^{-1} \phi(s, a)}, Q_{k-1,h}(s, a), H \right\}$
- 10: $\check{Q}_{k,h}(s, a) = \max \left\{ r_h(s, a) + \check{\mathbf{w}}_{k,h}^\top \phi(s, a) - \bar{\beta} \sqrt{\phi(s, a)^\top \Sigma_{k,h}^{-1} \phi(s, a)}, \check{Q}_{k-1,h}(s, a), 0 \right\}$
- 11: Set the last updating episode $k_{\text{last}} = k$
- 12: **else**
- 13: $Q_{k,h}(s, a) = Q_{k-1,h}(s, a)$
- 14: $\check{Q}_{k,h}(s, a) = \check{Q}_{k-1,h}(s, a)$
- 15: **end if**
- 16: $V_{k,h}(s) = \max_a Q_{k,h}(s, a)$
- 17: $\check{V}_{k,h}(s) = \max_a \check{Q}_{k,h}(s, a)$
- 18: **end for**
- 19: **for** stage $h = 1, \dots, H$ **do**
- 20: Take action $a_h^k \leftarrow \text{argmax}_a Q_{k,h}(s_h^k, a)$
- 21: Set the estimated variance $\sigma_{k,h}$ as in (4.1)
- 22: $\bar{\sigma}_{k,h} \leftarrow \max \left\{ \sigma_{k,h}, H, 2d^3 H^2 \|\phi(s_h^k, a_h^k)\|_{\Sigma_{k,h}^{-1}}^{1/2} \right\}$
- 23: $\Sigma_{k+1,h} = \Sigma_{k,h} + \bar{\sigma}_{k,h}^{-2} \phi(s_h^k, a_h^k) \phi(s_h^k, a_h^k)^\top$
- 24: Receive next state s_{h+1}^k
- 25: **end for**
- 26: **end for**



We take these terms from LSVI-UCB++ and privatize them with sufficient noise using a Gaussian mechanism

DP-LSVI-UCB++ Privacy Guarantee

[Theorem 3.1 of [Sahu. 2025](#)]. DP-LSVI-UCB++ satisfies (ϵ, δ) -JDP where $\epsilon = \rho + 2\sqrt{\rho \log(1/\delta)}$ for $\rho > 0$

Proof Sketch:

1. Compute the l_2 -sensitivity of each privatized estimator by considering neighboring sequences $\mathcal{U}, \mathcal{U}'$
2. Use a Gaussian mechanism with sufficient noise to ensure each is $\rho/4KH$ -zCDP
3. By Advanced Composition, it is ρ -zCDP
4. By a conversion from zCDP to DP, DP-LSVI-UCB++ is (ϵ, δ) -DP
5. Use the Billboard Lemma to conclude that DP-LSVI-UCB++ is (ϵ, δ) -JDP.

DP-LSVI-UCB++ Privacy Guarantee

[Theorem 3.1 of [Sahu. 2025](#)]. DP-LSVI-UCB++ satisfies (ϵ, δ) -JDP where $\epsilon = \rho + 2\sqrt{\rho \log(1/\delta)}$ for $\rho > 0$

Proof Sketch:

1. Compute the l_2 -sensitivity of each privatized estimator by considering neighboring sequences $\mathcal{U}, \mathcal{U}'$
2. Use a Gaussian mechanism with sufficient noise to ensure each is $\rho/4KH$ -zCDP
3. By Advanced Composition, it is ρ -zCDP
4. By a conversion from zCDP to DP, DP-LSVI-UCB++ is (ϵ, δ) -DP
5. Use the Billboard Lemma to conclude that DP-LSVI-UCB++ is (ϵ, δ) -JDP.

DP-LSVI-LUCB++ Privacy Guarantee

[Theorem 3.]

0

Definition (Zero-Concentrated Differential Privacy (zCDP)).

A randomized mechanism $\mathcal{M} : \mathcal{X}^n \rightarrow \mathcal{Y}$ is (ξ, ρ) -zCDP if

$\forall x, x' \in \mathcal{X}^n$ differing on a single entry and all $\alpha \in (1, \infty)$

Proof S

$$D_\alpha (\mathcal{M}(x) || \mathcal{M}(x')) \leq \xi + \rho\alpha$$

where $D_\alpha (\mathcal{M}(x) || \mathcal{M}(x'))$ is the α -Rényi divergence between distribution $\mathcal{M}(x)$ and $\mathcal{M}(x')$. Equivalently,

$$\mathbb{E} [e^{(\alpha-1)Z}] \leq e^{(\alpha-1)(\xi+\rho\alpha)}$$

DP-LSVI-UCB++ Regret Proof Sketch

1. Use the privatized terms and prove their utility i.e. how close are they to the non-privatized terms. Since we used a Gaussian mechanism, it is sufficient to use a Gaussian concentration inequality on the matrix operator norm or the l_2 -norm
2. Use the private terms in place of the non-private terms and use the arguments from LSVI-UCB++ to find the upper confidence bonuses using a Bernstein self-normalized concentration inequality, uniform covering arguments, elliptical potentials, and utility of the privatized terms
3. Use the bonuses to prove optimism and pessimism of the privatized Q-value function

DP-LSVI-UCB++ Regret Proof Sketch

First, let us define the following events

\mathcal{E} : Accurate value predictions despite privatized regression

$\tilde{\mathcal{E}}$: Sharper bounds using Bernstein-style control of noise and variance

\mathcal{E}_1 : Estimated values of learned policy \approx actual value of true policy (uniform stability)

\mathcal{E}_2 : Optimistic and pessimistic values sandwich the truth (stability under noise)

\mathcal{E}_3 : Total estimation variance is bounded (learnability of the environment)

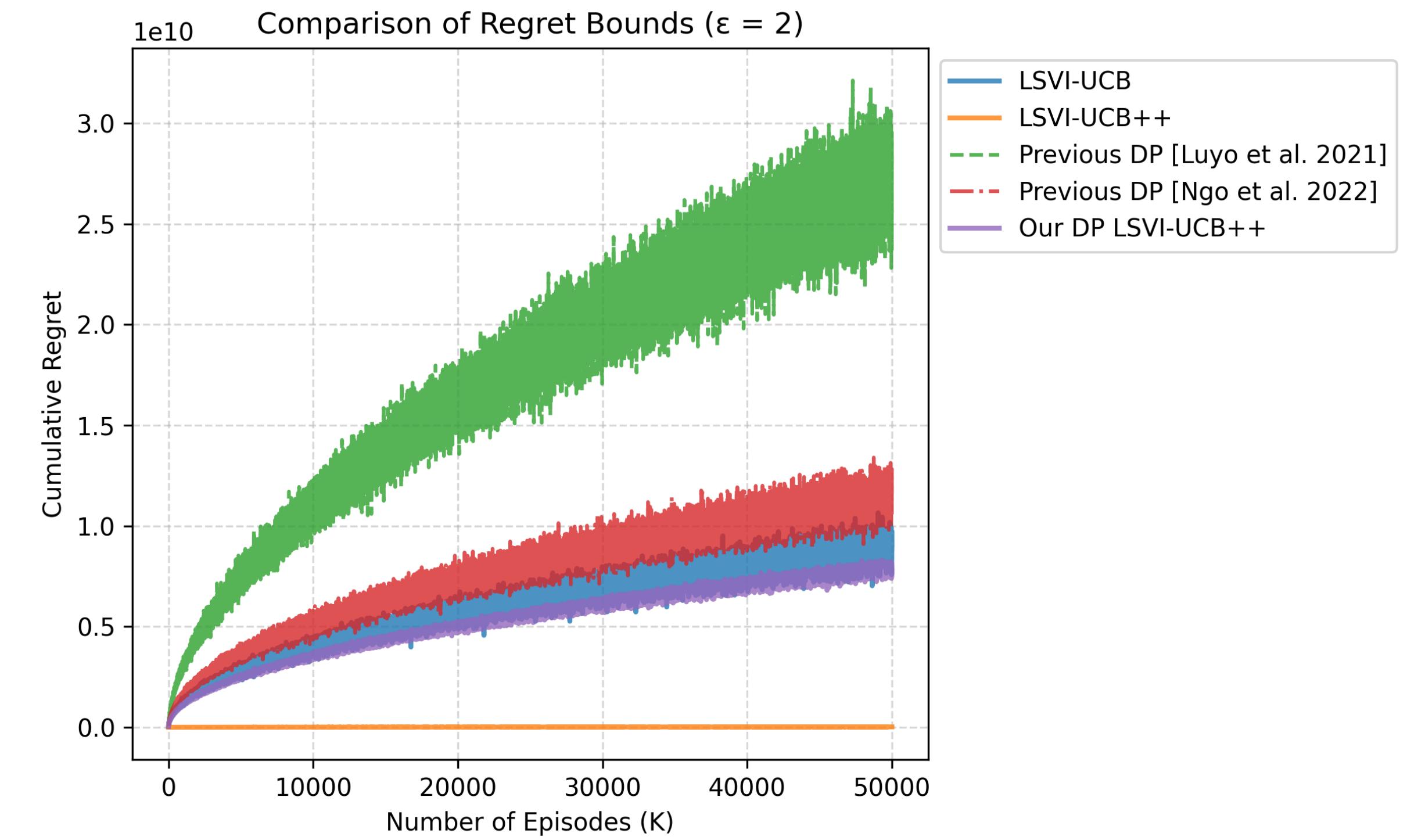
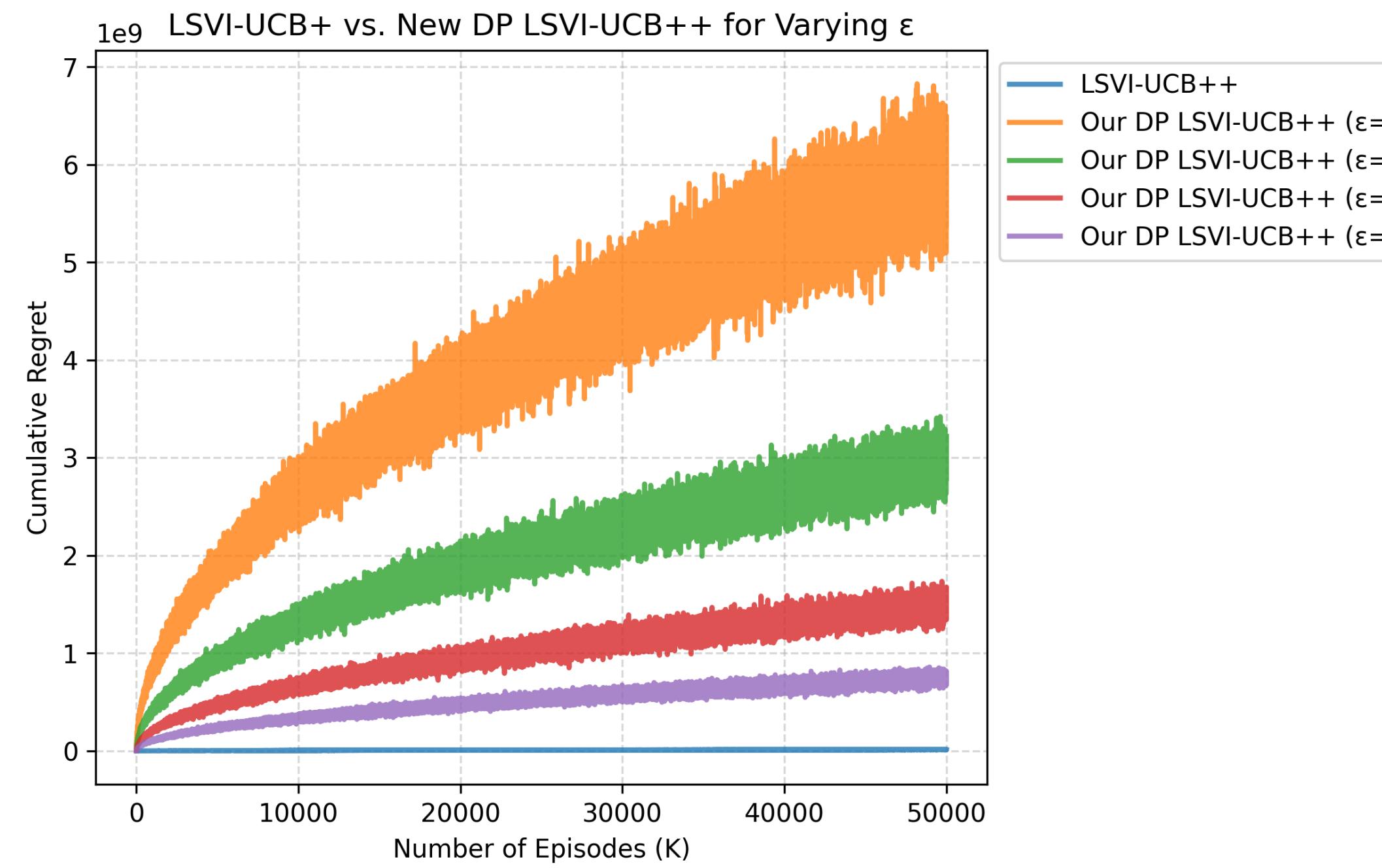
DP-LSVI-UCB++ Regret Proof Sketch

Conditioned on the event $\mathcal{E} \cap \widetilde{\mathcal{E}} \cap \mathcal{E}_1 \cap \mathcal{E}_2 \cap \mathcal{E}_3$

$$\begin{aligned}\mathcal{R}(K) &= \sum_{k=1}^K \left(V_1^*(s_1^k) - \widetilde{\widehat{V}}_{k,1}^{\pi_k}(s_1^k) \right) \\ &\leq \sum_{k=1}^K \left(\widetilde{\widehat{V}}_{k,1}(s_1^k) - \widetilde{\widehat{V}}_{k,1}^{\hat{\pi}^k}(s_1^k) \right) \\ &\leq 16d^4H^8\iota + 40\beta d^7H^5\iota + 8\beta \sqrt{2dH\iota \sum_{h=1}^H \sum_{k=1}^K (\tilde{\sigma}_{k,h}^2 + H)} + 4\sqrt{H^3K \log(H/\delta)} \\ &\leq \widetilde{O} \left(d\sqrt{H^3K} + \frac{H^{15/4}d^{7/6}K^{1/2} \log(10dKH/\delta)}{\epsilon} \right)\end{aligned}$$

DP-LSVI-UCB++ enjoys a privacy guarantee at (almost) no drop in utility

Environment Setup: We use a 6-state chain environment with two actions: left and right. The agent starts on the left and aims to reach the rightmost state for higher rewards. we set the planning horizon $H = 20$ and run $K = 50,000$ episodes.



Conclusion

- DP-LSVI-UCB++ is a (ϵ, δ) -**JDP** algorithm for linear MDP that achieves the new **state-of-the-art regret bound** by using **Bernstein concentration**, **Gaussian mechanisms**, and **GOE perturbations** for tight utility-privacy tradeoff
- Our results show that theoretically and empirically, we match or outperform non-private baselines and do better than previous work in this area
- Future directions
 - Extending these results to the low-rank MDP setting
 - Exploring alternative mechanism that adapt noise dynamically based on the observed data's sensitivity could lead to improved regret bounds

Thanks For Listening! Questions?