# Privacy-Preserving Policy Synthesis for MDPs

**Matthew Hale**

**Department of Mechanical and Aerospace Engineering**
**University of Florida**

**Joint work with Parham Gohari, Bo Wu, and Ufuk Topcu at UT Austin**

**AFOSR Center of Excellence on Assured Autonomy in Contested Environments**
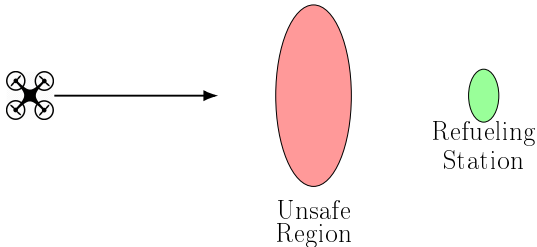**April 15th, 2020**

▶ Adversaries can observe us, and actions can reveal intent/knowledge

- Adversaries can observe us, and actions can reveal intent/knowledge
  - Direction of travel can reveal a destination
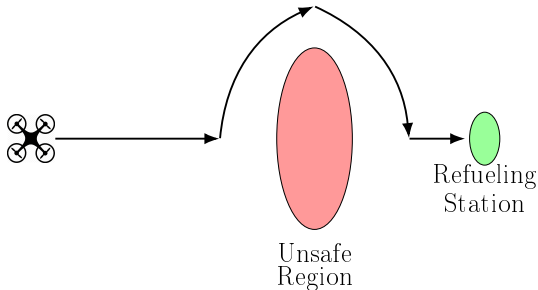


Unsafe
Region

Refueling
Station

- ▶ Adversaries can observe us, and actions can reveal intent/knowledge
  - ▶ Direction of travel can reveal a destination
  - ▶ Avoiding an area can reveal knowledge of hazards
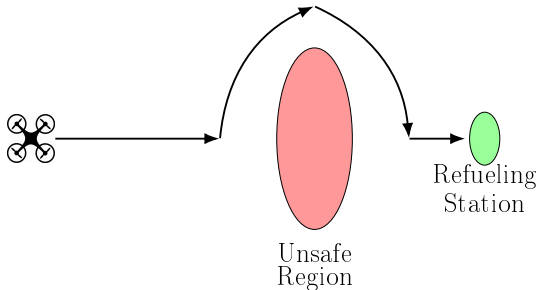


Refueling
Station

Unsafe
Region

▶ Adversaries can observe us, and actions can reveal intent/knowledge
  ▶ Direction of travel can reveal a destination
  ▶ Avoiding an area can reveal knowledge of hazards



Refueling
Station

Unsafe
Region

**Fundamental Problem**

A task must be completed without revealing the information driving it.

▶ There are 3 goals in this work:

► There are 3 goals in this work:



Goal #1

**1** **Provably** protect the information driving a decision

▶ There are 3 goals in this work:



1. **Provably** protect the information driving a decision
2. Synthesize an altered, privacy-preserving decision policy
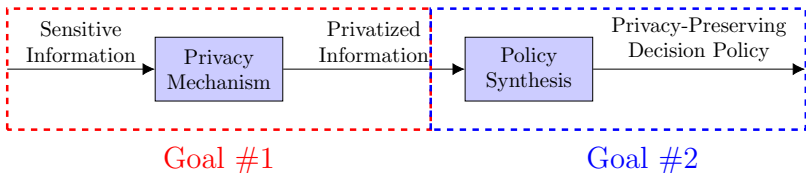
► There are 3 goals in this work:



Goal #1          Goal #2

1  **Provably** protect the information driving a decision
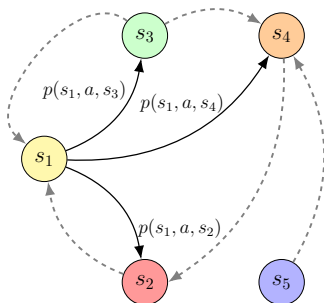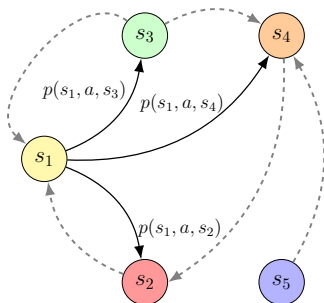2  Synthesize an altered, privacy-preserving decision policy
3  Quantify the "cost of privacy," formalize tradeoffs between privacy and performance

▶ We consider MDP models:

▶ We consider MDP models:



▶ We want to take actions that don't reveal transition probabilities

- We consider MDP models:



- We want to take actions that don't reveal transition probabilities
- In state $s$, taking action $a$ transitions to state $s'$ with prob. $p(s, a, s')$

▶ We consider MDP models:



▶ We want to take actions that don't reveal transition probabilities
▶ In state $s$, taking action $a$ transitions to state $s'$ with prob. $p(s, a, s')$
▶ For all $s$, we have $p(s, a, s') \geq 0$ and $\sum_{s'} p(s, a, s') = 1$

## Differential Privacy (DP)

DP is a privacy framework with several key features:

- ▶ It offers a formal definition of "privacy"

## Differential Privacy (DP)

DP is a privacy framework with several key features:

- ▶ It offers a formal definition of "privacy"
- ▶ It is immune to post-processing
  - ▶ $x$ private $\Rightarrow f(x)$ private for all $f$

## Differential Privacy (DP)

DP is a privacy framework with several key features:

- ▶ It offers a formal definition of "privacy"
- ▶ It is immune to post-processing
  - ▶ $x$ private $\Rightarrow f(x)$ private for all $f$
- ▶ It is robust to side information

## Differential Privacy (DP)

DP is a privacy framework with several key features:

- ▶ It offers a formal definition of "privacy"
- ▶ It is immune to post-processing
    - ▶ $x$ private $\Rightarrow f(x)$ private for all $f$
- ▶ It is robust to side information

▶ Used by:

Apple

Google

Uber

## Differential Privacy (DP)

DP is a privacy framework with several key features:

- It offers a formal definition of "privacy"
- It is immune to post-processing
    - $x$ private $\Rightarrow f(x)$ private for all $f$
- It is robust to side information

| Apple | Google | Uber |
|-------|--------|------|

- Used by:



## DP Idea

**Make probability vectors look "similar"**

**Fundamental Inequality of Differential Privacy**

For probability vectors $p$ and $q$, we generate private forms $\tilde{p}$, $\tilde{q}$ to satisfy
$$\mathbb{P}(\tilde{p}) \leq e^{\epsilon}\mathbb{P}(\tilde{q}) + \delta$$

**Fundamental Inequality of Differential Privacy**

For probability vectors $p$ and $q$, we generate private forms $\tilde{p}$, $\tilde{q}$ to satisfy
$$\mathbb{P}(\tilde{p}) \leq e^{\epsilon}\mathbb{P}(\tilde{q}) + \delta$$

**Fundamental Inequality of Differential Privacy**

For probability vectors $p$ and $q$, we generate private forms $\tilde{p}$, $\tilde{q}$ to satisfy
$$\mathbb{P}(\tilde{p}) \leq e^{\epsilon}\mathbb{P}(\tilde{q}) + \delta$$

▶ For us this will take the form

$$\xrightarrow[p]{\text{Vector}} \boxed{\text{Dirichlet}(kp)} \xrightarrow[\tilde{p}]{\text{Private Vector}}$$

where
$$\text{Dirichlet}(kp) = \frac{\Gamma\left(\sum_{i=1}^{n} kp_i\right)}{\prod_{i=1}^{n} \Gamma(kp_i)} \prod_{i=1}^{n} x_i^{kp_i}$$

**Fundamental Inequality of Differential Privacy**

For probability vectors $p$ and $q$, we generate private forms $\tilde{p}$, $\tilde{q}$ to satisfy
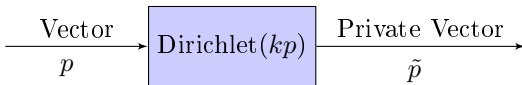$$\mathbb{P}(\tilde{p}) \leq e^\epsilon \mathbb{P}(\tilde{q}) + \delta$$

▶ For us this will take the form

$$\xrightarrow{\text{Vector } p} \boxed{\text{Dirichlet}(kp)} \xrightarrow{\text{Private Vector } \tilde{p}}$$

where
$$\text{Dirichlet}(kp) = \frac{\Gamma\left(\sum_{i=1}^n kp_i\right)}{\prod_{i=1}^n \Gamma(kp_i)} \prod_{i=1}^n x_i^{kp_i}$$
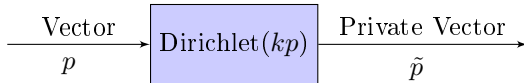
**Privacy Theorem (ACC 2020 Paper; Gohari, Wu, Hale, and Topcu)**

The Dirichlet mechanism provides $\big(\epsilon(k), \delta(k)\big)$-differential privacy.

▶ Example: $k = 24$ gives $(1.18, 0.05)$-DP

▶ Objective is to maximize the accumulated reward
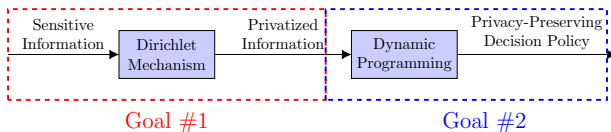
$$\sum_{t=1}^{T} \gamma^t R(s_t, a_t)$$

▶ Objective is to maximize the accumulated reward

$$\sum_{t=1}^{T} \gamma^t R(s_t, a_t)$$

▶ We privatize transition probabilities, then synthesize a decision policy

- Objective is to maximize the accumulated reward

$$\sum_{t=1}^{T} \gamma^t R(s_t, a_t)$$

- We privatize transition probabilities, then synthesize a decision policy



Goal #1          Goal #2

- Synthesis is just post-processing, so its output protects $p$ as well

► Objective is to maximize the accumulated reward

$$\sum_{t=1}^{T} \gamma^t R(s_t, a_t)$$

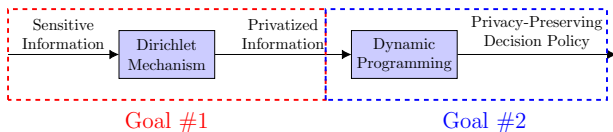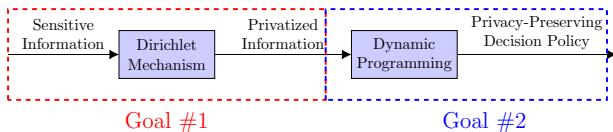► We privatize transition probabilities, then synthesize a decision policy



Goal #1                          Goal #2

► Synthesis is just post-processing, so its output protects $p$ as well
► Our actions protect transition probabilities!

▶ Set **Cost of privacy** $=$ (Reward without DP) - (Reward with DP)

- Set **Cost of privacy** = (Reward without DP) - (Reward with DP)

### Theorem: Cost of Privacy

- Privatize all transition probabilities with the Dirichlet mechanism. Then:

$$\text{Cost of privacy} \leq w_0 - v_0,$$

where, for all $t \in \{0, \ldots, T\}$,

$$v_t = \sum_{a \in \mathcal{A}_s} \pi(a \mid s) \left( R(s, a) + \gamma \min_{p \in \hat{\mathcal{P}}} p(s, a, s') v_{t+1}(s') \right)$$

$$w_t = \sum_{a \in \mathcal{A}_s} \pi(a \mid s) \left( R(s, a) + \gamma \max_{p \in \hat{\mathcal{P}}} p(s, a, s') w_{t+1}(s') \right)$$

UF UNIVERSITY of FLORIDA    Duke    TEXAS The University of Texas at Austin    UC SANTA CRUZ

▶ Set **Cost of privacy** = (Reward without DP) - (Reward with DP)

**Theorem: Cost of Privacy**

▶ Privatize all transition probabilities with the Dirichlet mechanism. Then:
$$\text{Cost of privacy} \leq w_0 - v_0,$$

where, for all $t \in \{0, \ldots, T\}$,

$$v_t = \sum_{a \in \mathcal{A}_s} \pi(a \mid s) \left( R(s,a) + \gamma \min_{p \in \hat{\mathcal{P}}} p(s,a,s') v_{t+1}(s') \right)$$

$$w_t = \sum_{a \in \mathcal{A}_s} \pi(a \mid s) \left( R(s,a) + \gamma \max_{p \in \hat{\mathcal{P}}} p(s,a,s') w_{t+1}(s') \right)$$
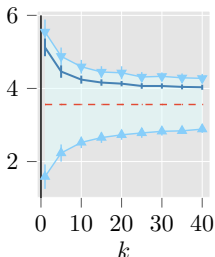
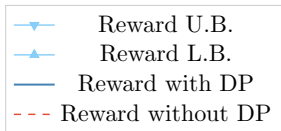▶ This is computable in $O\left(T|S|^{4.5}|A_s|\right)$ time

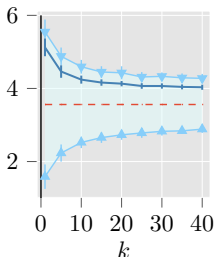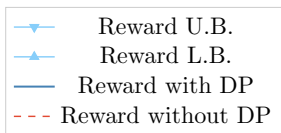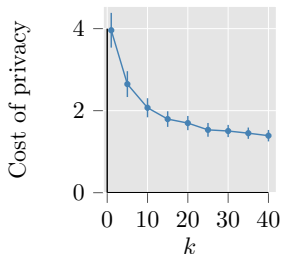- Implement privacy and synthesize a policy for a $30$-state MDP

- Implement privacy and synthesize a policy for a $30$-state MDP
- Total time required is $4.88$s on a desktop computer

- Implement privacy and synthesize a policy for a $30$-state MDP
- Total time required is $4.88$s on a desktop computer

- Implement privacy and synthesize a policy for a $30$-state MDP
- Total time required is $4.88$s on a desktop computer

- Incorporating temporal logic specifications:
  - What are the tradeoffs in privacy, safety, and performance?
  - What is the complexity of computing a safe, private policy and bounds on the cost of safety & privacy?

- ▶ Incorporating temporal logic specifications:
  - ▶ What are the tradeoffs in privacy, safety, and performance?
  - ▶ What is the complexity of computing a safe, private policy and bounds on the cost of safety & privacy?
- ▶ What are the effects of privatizing other characteristics of MDPs?

# Thank you