# Safe Multi-Agent Reinforcement Learning in Polynomial Time
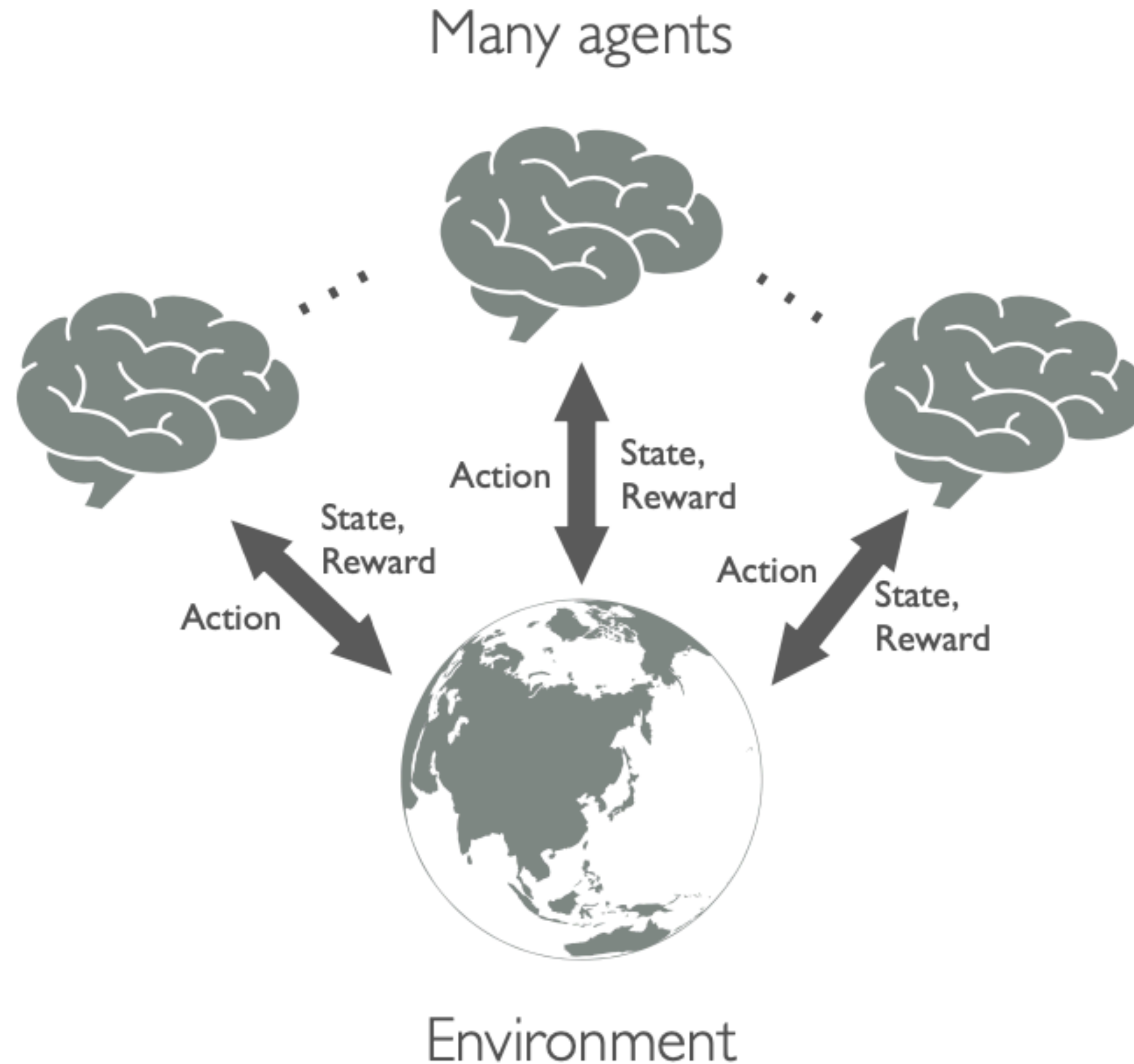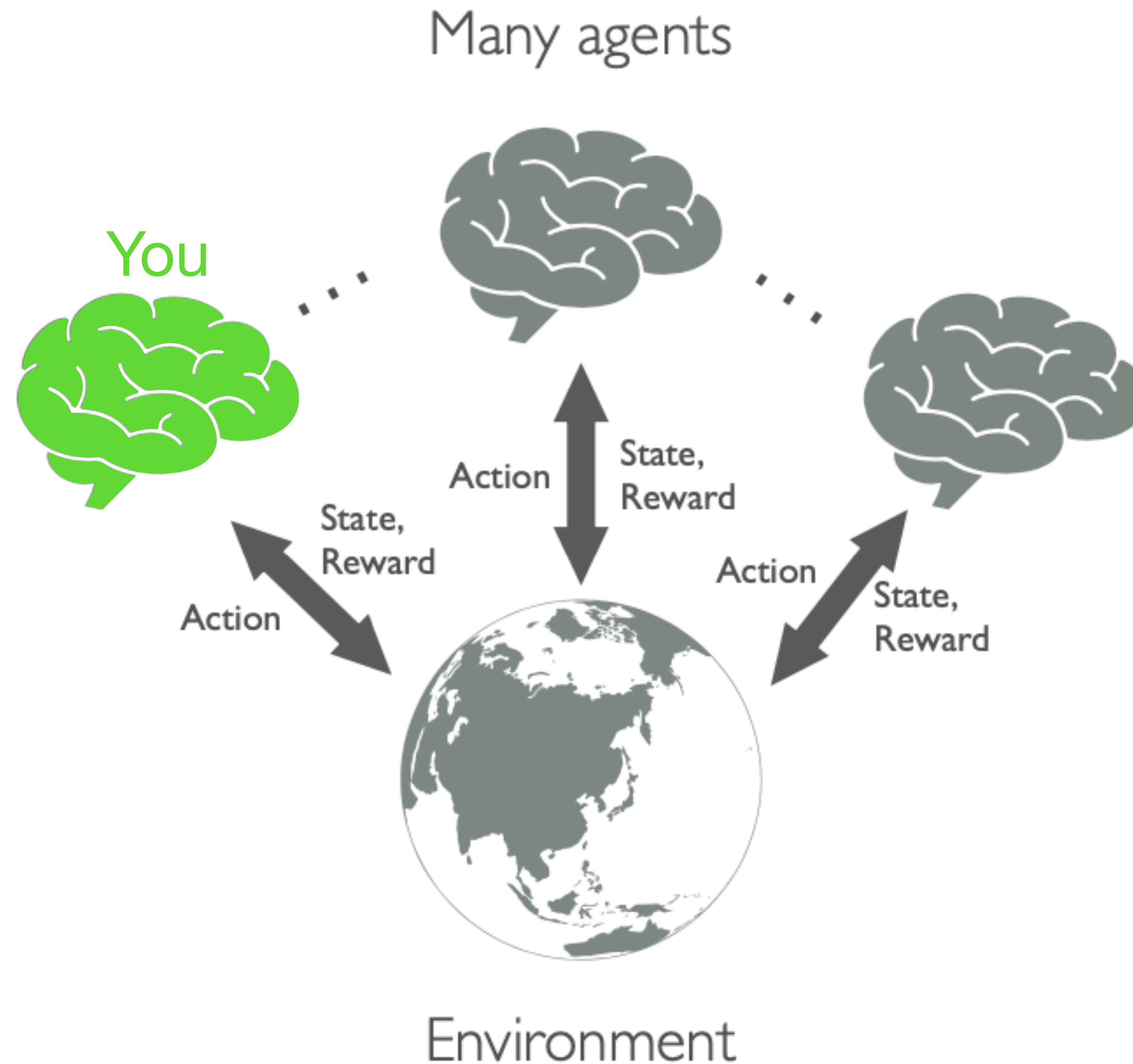
Jeremy McMahan

# Safety Concerns
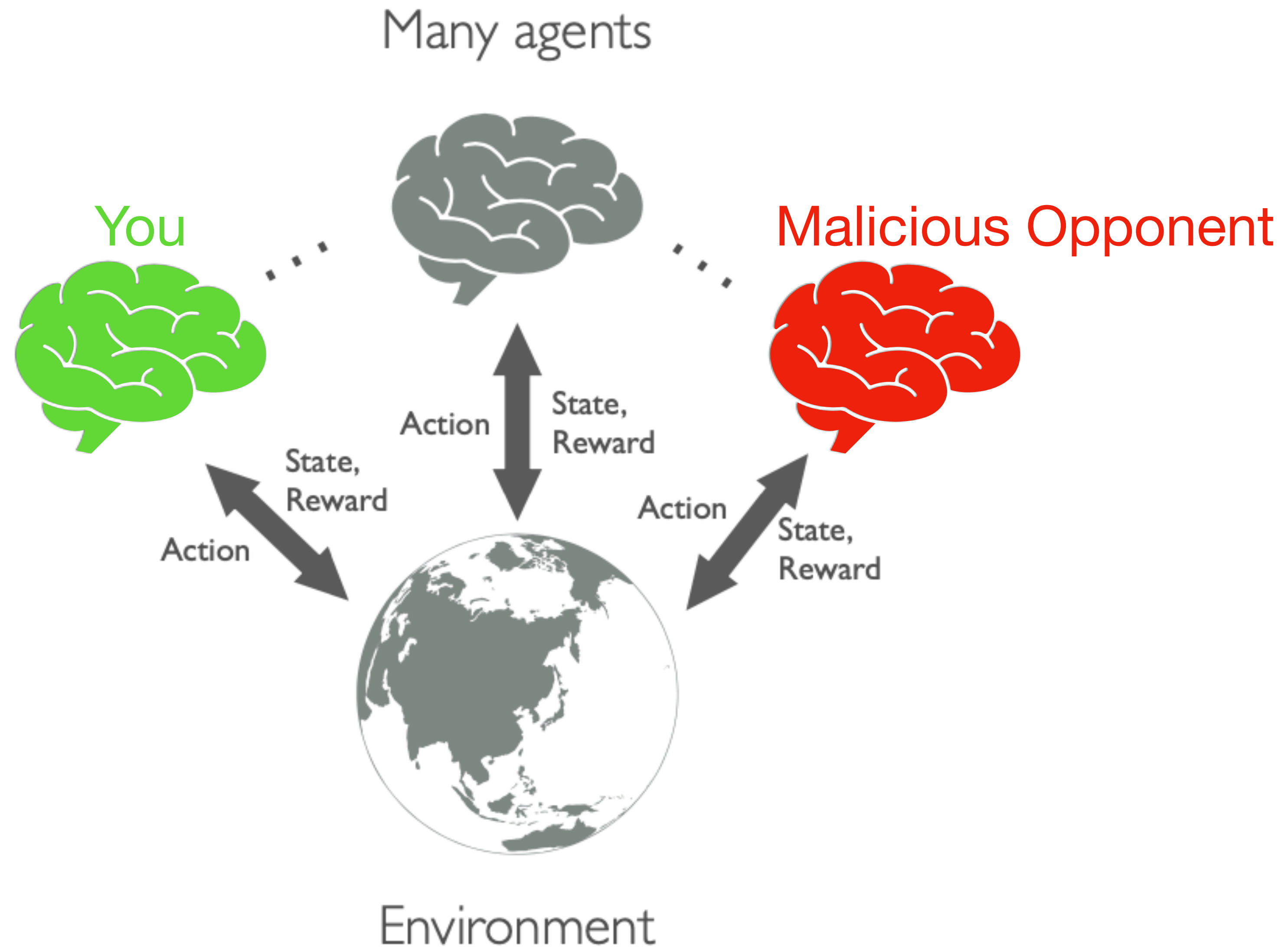
# Safety Concerns
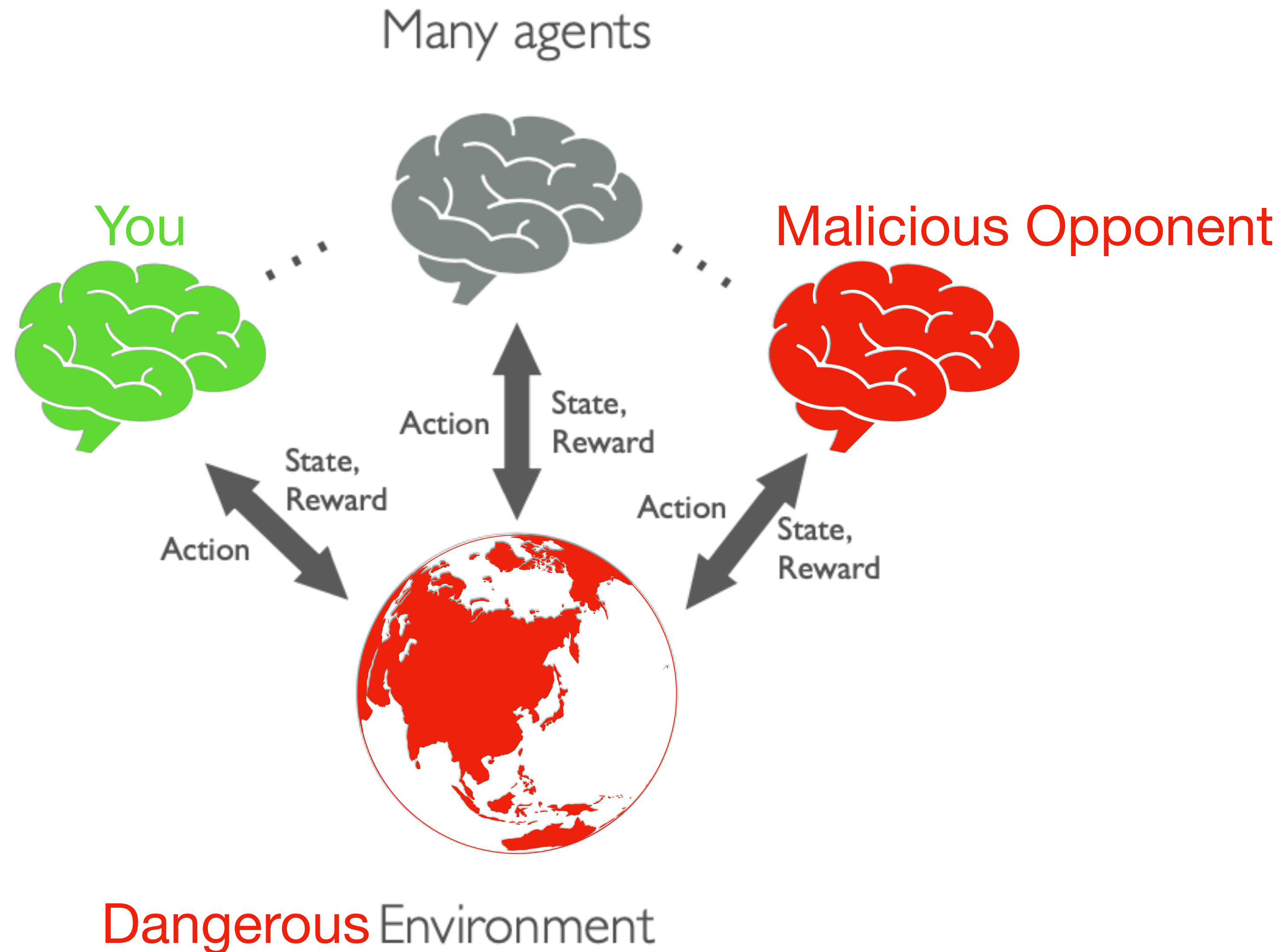
# Safety Concerns

# Safety Concerns



Many agents

You

Malicious Opponent

Action

State, Reward

Action

State, Reward

Action

State, Reward

Environment

# Safety Concerns



Many agents

You

Malicious Opponent

Action State, Reward

State, Reward Action

Action State, Reward

Dangerous Environment

# Safety Landscape

# Safety Landscape

Safety from **Agents**:

# Safety Landscape

Safety from **Agents**:

*Adversarial MARL*

# Safety Landscape

Safety from **Agents**:
*Adversarial MARL*

Safety from **Environment**:

# Safety Landscape

Safety from **Agents**:

*Adversarial MARL*

Safety from **Environment**:

*Constrained MARL*

# Safety Landscape

Safety from **Agents**:
*Adversarial MARL*

Safety from **Environment**:
*Constrained MARL*

# Safety Landscape

**Safety from Agents:**
*Adversarial MARL*

1. Manipulation Attacks

**Safety from Environment:**
*Constrained MARL*

# Safety Landscape

Safety from **Agents**:
*Adversarial MARL*

1. Manipulation Attacks

2. Misinformation Attacks

Safety from **Environment**:
*Constrained MARL*

# Safety Landscape

Safety from **Agents**:
*Adversarial MARL*

1. Manipulation Attacks

2. Misinformation Attacks

Safety from **Environment**:
*Constrained MARL*

1. Anytime Constraints

# Safety Landscape

Safety from **Agents**:
*Adversarial MARL*

1. Manipulation Attacks

2. Misinformation Attacks

Safety from **Environment**:
*Constrained MARL*

1. Anytime Constraints

2. Single-Constraint FPTAS

# Safety Landscape

Safety from **Agents**:
*Adversarial MARL*

1.   Manipulation Attacks

2. Misinformation Attacks

Safety from **Environment**:
*Constrained MARL*

1. Anytime Constraints

2. Single-Constraint FPTAS

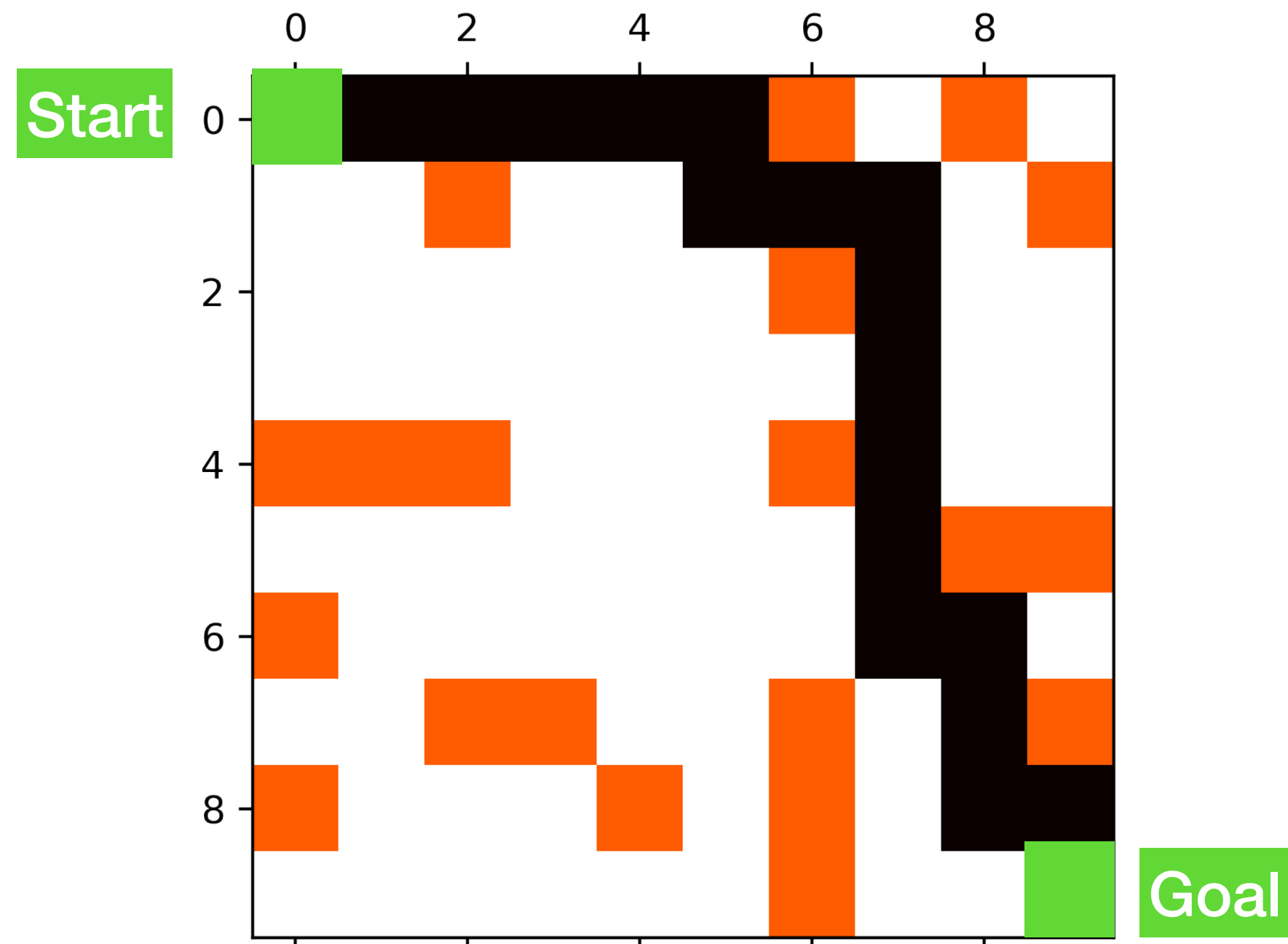3. Multi-Constraint Bicriteria

# Adversarial MARL
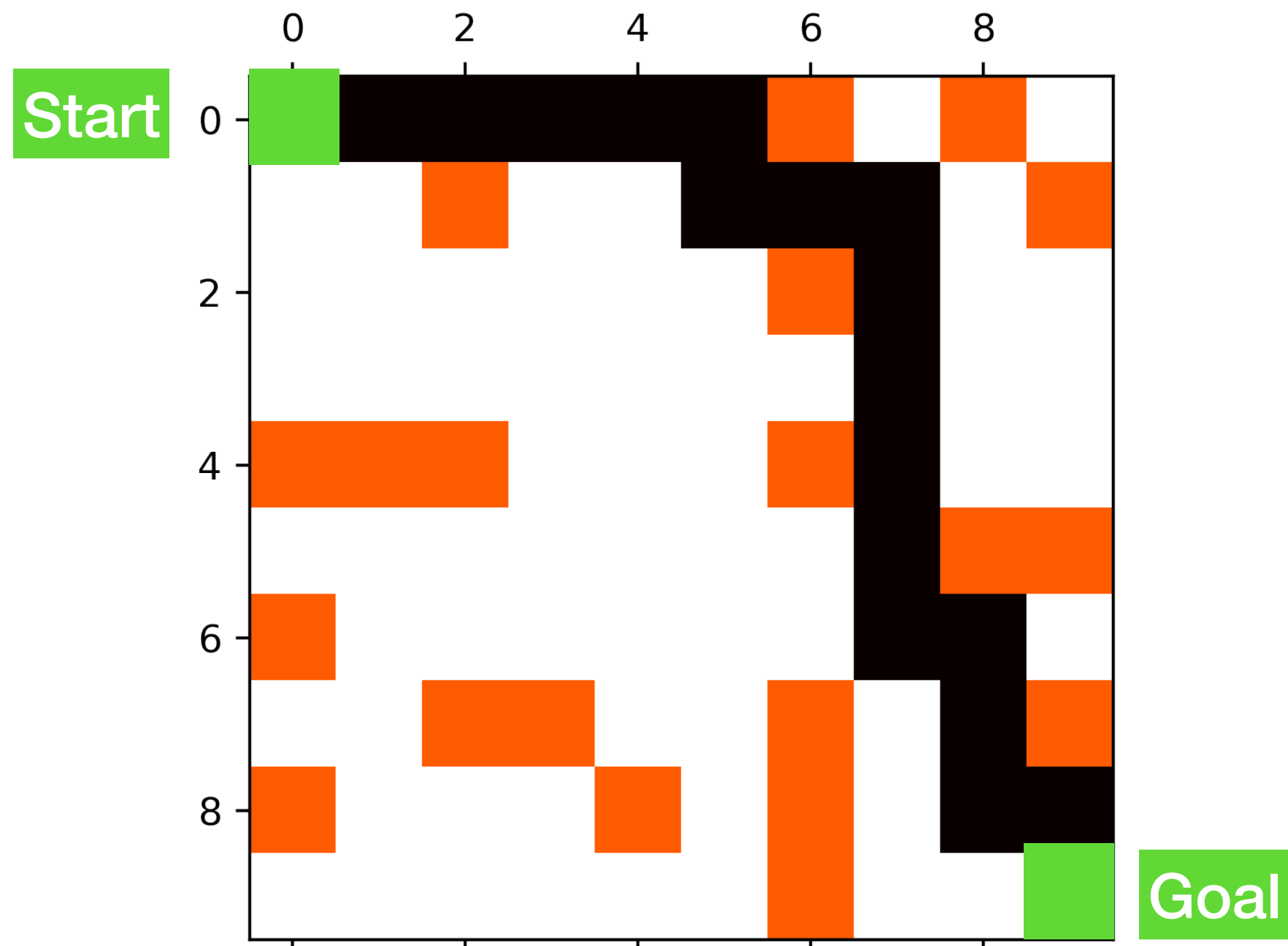
# Manipulation Attacks
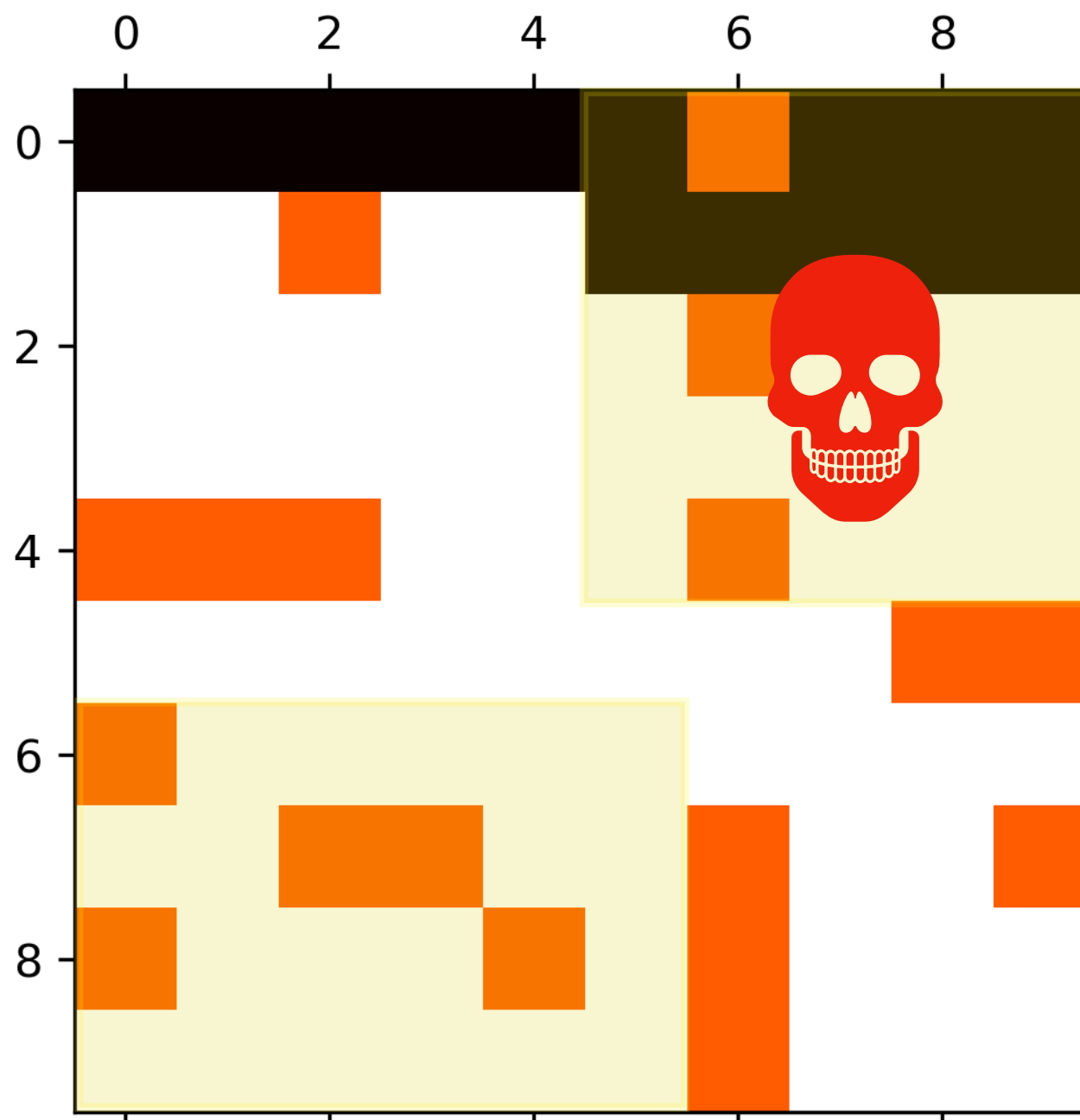
*AAAI 2024

# Motivation
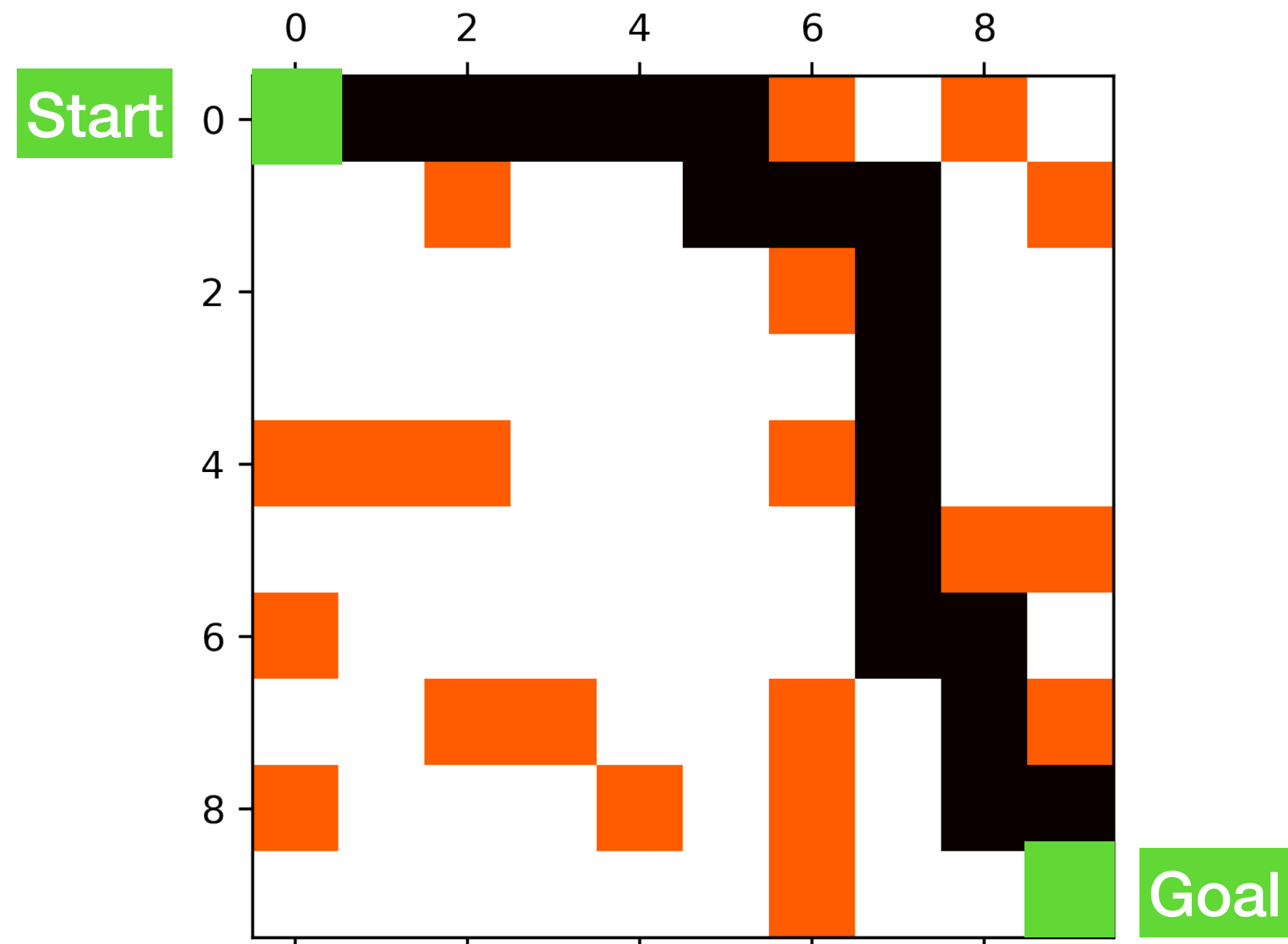
# Motivation

Optimal $\pi^*$

# Motivation



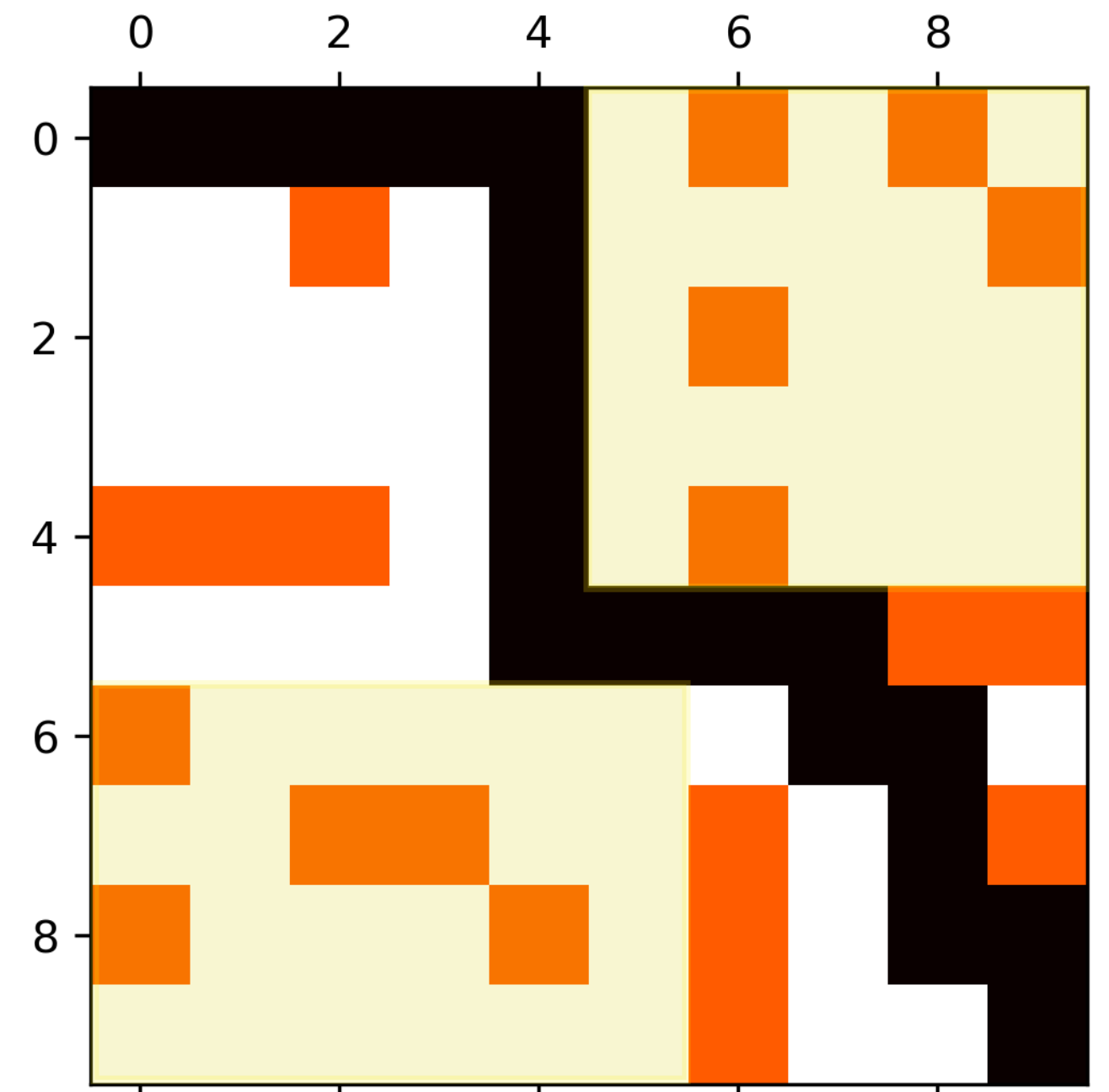Optimal $\pi^*$

Attacked $\pi^*$

# Motivation



Optimal $\pi^*$

Attacked $\pi^*$
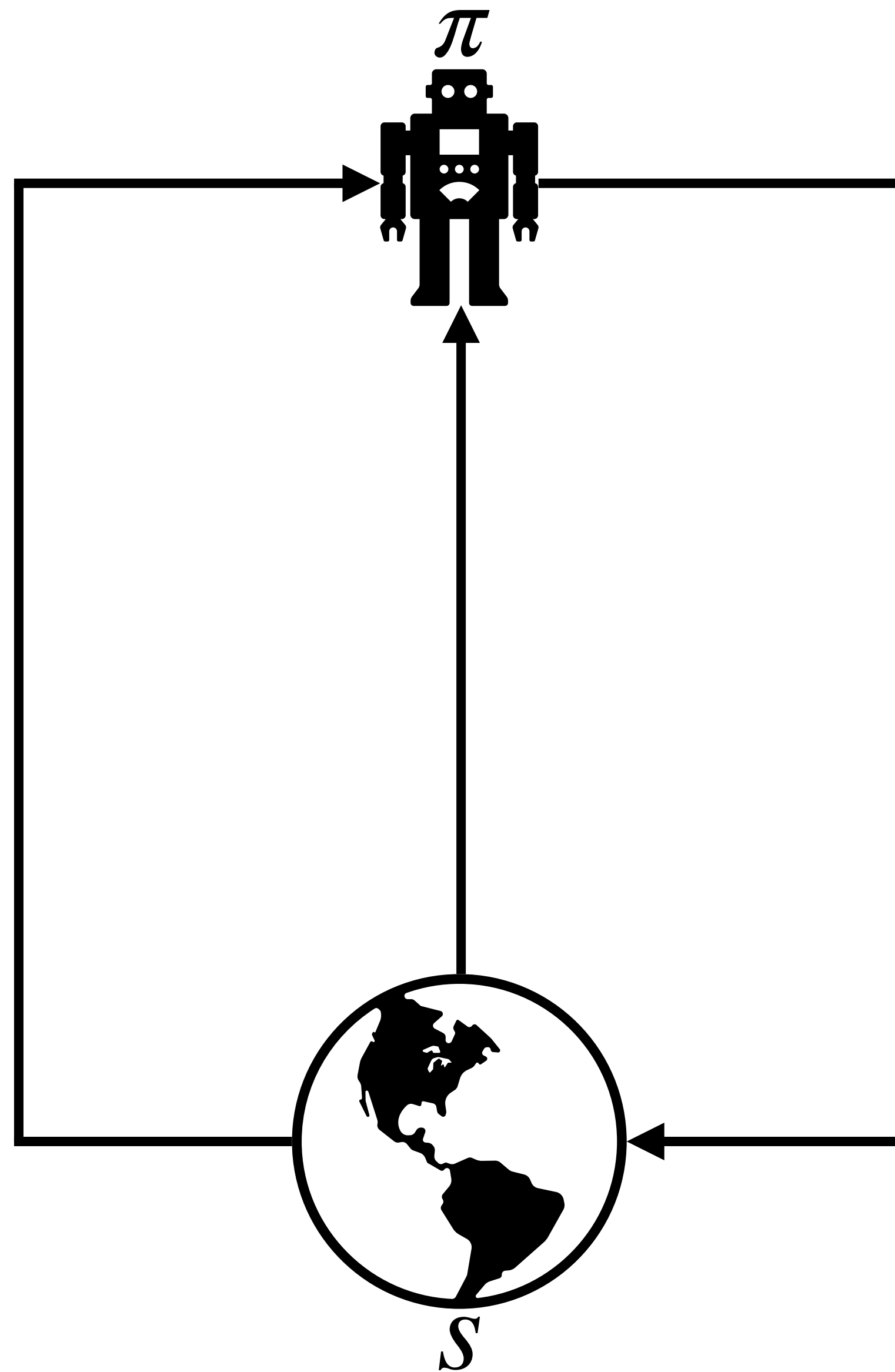
Robust $\hat{\pi}$
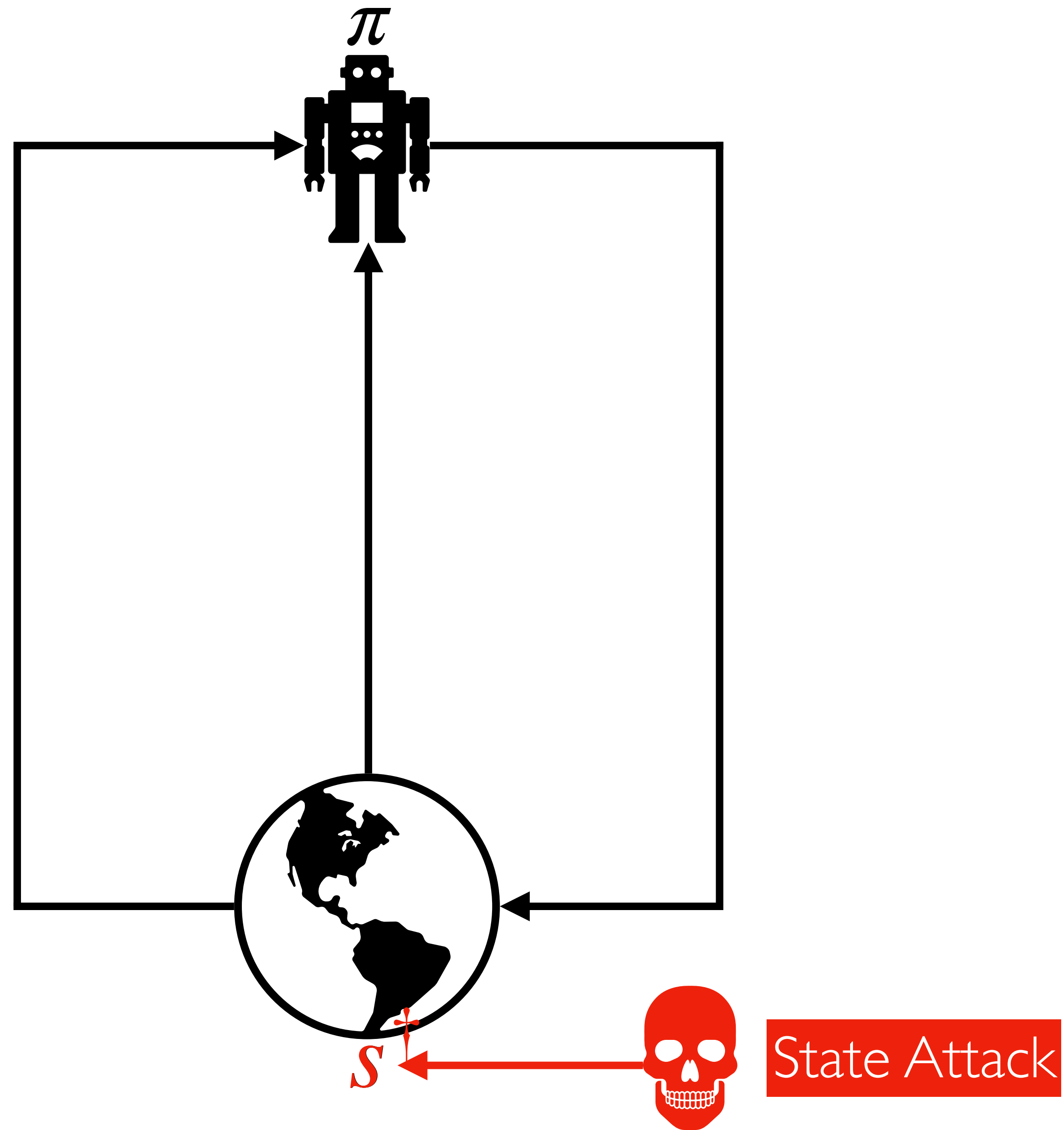
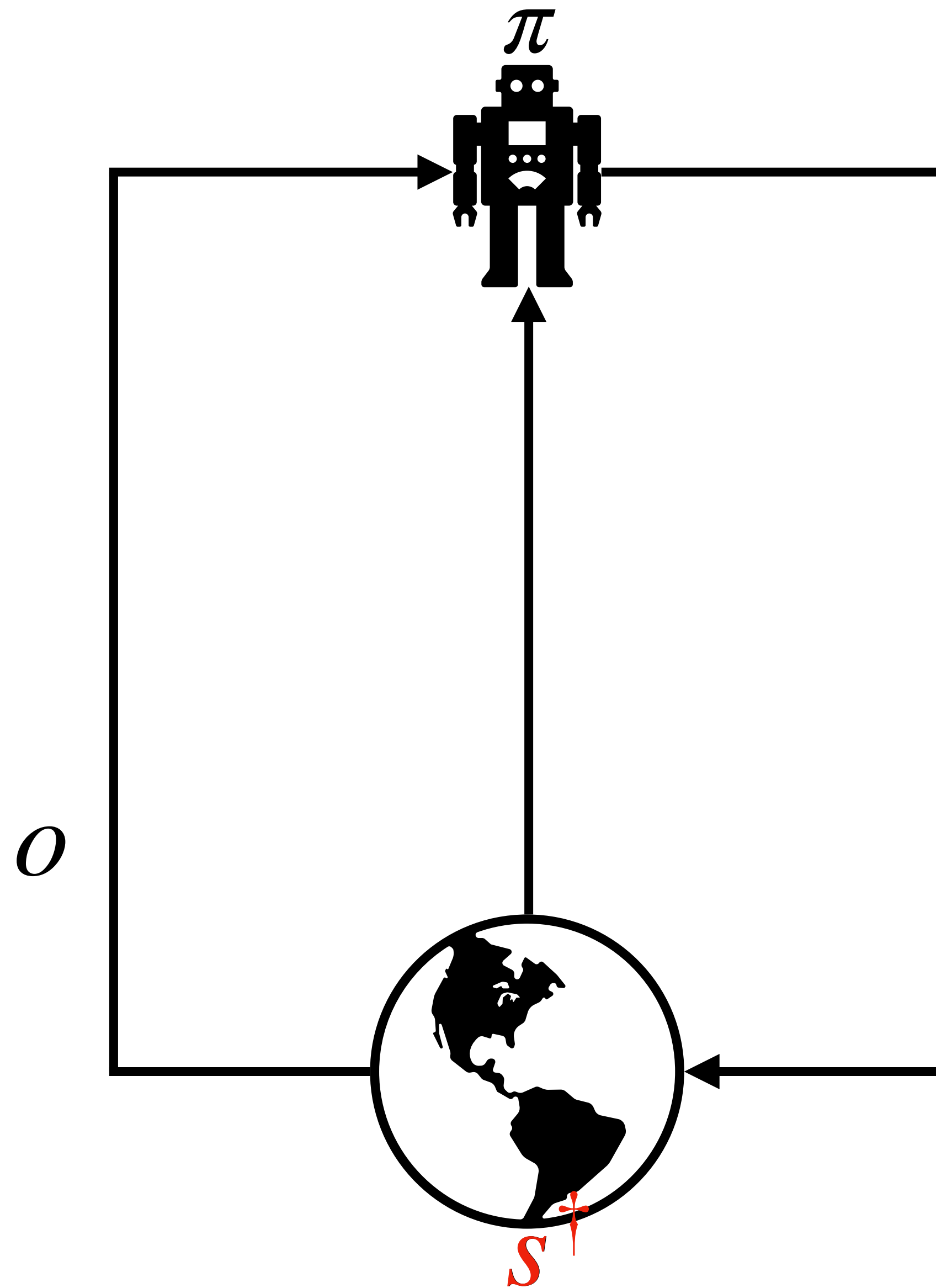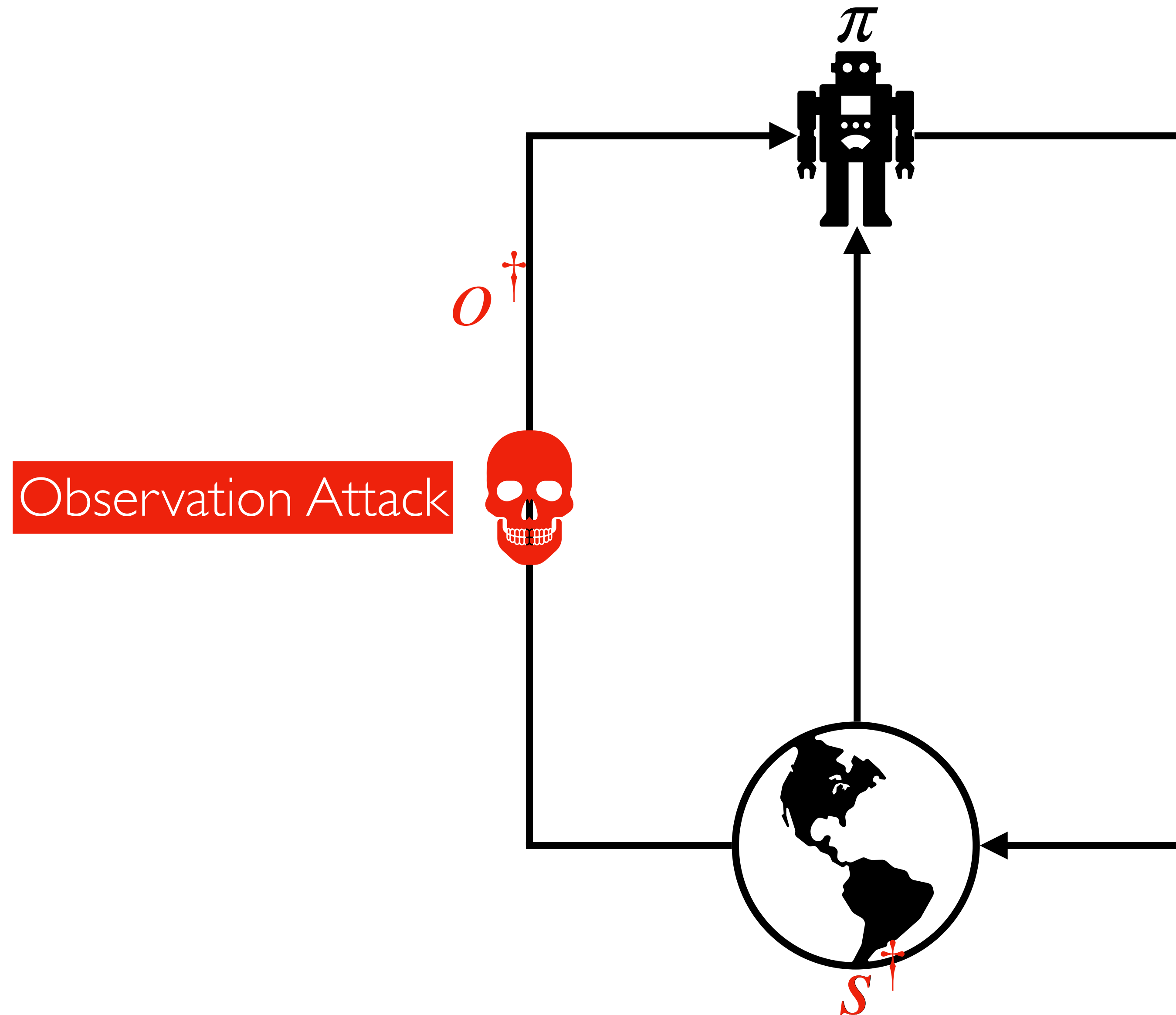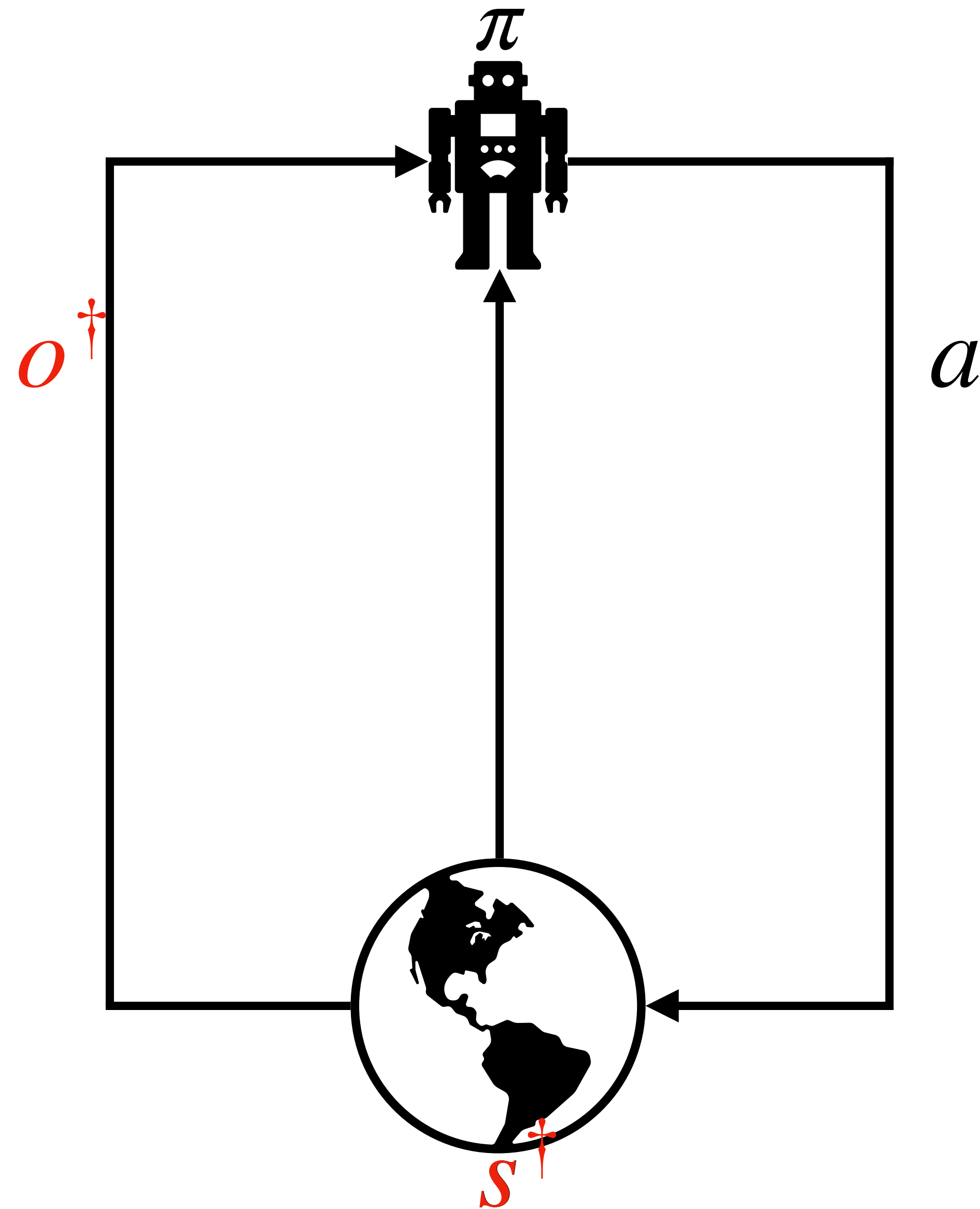# Attack Surfaces

# Attack Surfaces

$\pi$

$S$

# Attack Surfaces

# Attack Surfaces

# Attack Surfaces

# Attack Surfaces

# Attack Surfaces



$\pi$

$o^\dagger$

Action Attack

$a^\dagger$

$s^\dagger$

# Attack Surfaces

# Attack Surfaces

# Attacker's Perspective

# Attacker's Perspective

Attacker has its own reward $g(s_t, a_t, r_t)$ that depends on the victim's.

# Attacker's Perspective

Attacker has its own reward $g(s_t, a_t, r_t)$ that depends on the victim's.

**Definition 1** (Attack Problem)**.** For any $\pi$, the attacker's seeks a policy $\nu^* \in N$ that maximizes its expected reward from the victim-attacker-$M$ interaction:

$$\nu^* \in \arg\max_{\nu \in N} \mathbb{E}_M^{\pi,\nu} \left[ \sum_{t=0}^{\infty} \gamma^t g(s_t, a_t, r_t) \right].$$

# Adversarial Decomposition

# Adversarial Decomposition

# Adversarial Decomposition

# Adversarial Decomposition



Attacker MDP $\overline{M}$

$\pi$

$\pi(o^{\dagger})$

# Attack Results

# Attack Results

**Theorem**: *An optimal attack involving any combination of attack surfaces can be computed in time poly$(|M|, |\pi|)$.*

# Attack Results

**Theorem**: *An optimal attack involving any combination of attack surfaces can be computed in time poly$(|M|, |\pi|)$.*

*First results beyond observation attacks!*

# The Defense Problem

# The Defense Problem

Let $(V_1^{\pi,\nu}, V_2^{\pi,\nu})$ denote the victim's and attacker's value, respectively.

# The Defense Problem

Let $(V_1^{\pi,\nu}, V_2^{\pi,\nu})$ denote the victim's and attacker's value, respectively.

**Definition 2** (Defense Problem). The victim seeks a policy $\pi^*$ that maximizes its expected reward from the victim-attacker-$M$ interaction under the worst-case attack:

$$\pi^* \in \arg\max_{\pi \in \Pi} \ \min_{\nu \in BR(\pi)} \ V_1^{\pi,\nu}.$$

# The Defense Problem

Let $(V_1^{\pi,\nu}, V_2^{\pi,\nu})$ denote the victim's and attacker's value, respectively.

**Definition 2** (Defense Problem)**.** The victim seeks a policy $\pi^*$ that maximizes its expected reward from the victim-attacker-$M$ interaction under the worst-case attack:

$$\pi^* \in \arg\max_{\pi \in \Pi} \min_{\nu \in BR(\pi)} V_1^{\pi,\nu}.$$

$$BR(\pi) := \arg\max_{\nu \in N} V_2^{\pi,\nu}$$

# The Defense Problem

Let $(V_1^{\pi,\nu}, V_2^{\pi,\nu})$ denote the victim's and attacker's value, respectively.

**Definition 2** (Defense Problem)**.** The victim seeks a policy $\pi^*$ that maximizes its expected reward from the victim-attacker-$M$ interaction under the worst-case attack:

$$\pi^* \in \arg\max_{\pi \in \Pi} \ \min_{\nu \in BR(\pi)} \ V_1^{\pi,\nu}.$$

$$BR(\pi) := \arg\max_{\nu \in N} V_2^{\pi,\nu}$$

*Defense = WSE in a meta game.*

# Bottlenecks

# Bottlenecks

- WSE need not exist.

# Bottlenecks

- WSE need not exist.

- WSE are generally non-Markovian.

# Bottlenecks

- WSE need not exist.

- WSE are generally non-Markovian.

**Proposition**: *The defense problem is as hard as solving POMDPs. Thus, is NP-hard to even approximate.*

# Approach

# Approach

*Solution*: *ban observation attacks.*

# Approach

*Solution*: *ban observation attacks.*

$\overline{G}$

# Approach

*Solution: ban observation attacks.*

$\overline{G}$    Zero-sum:

# Approach

*Solution*: *ban observation attacks.*

$\overline{G}$

Zero-sum:

WSE → MPNE

# Approach

*Solution*: ban observation attacks.

$\overline{G}$

Zero-sum:

WSE → MPNE

General-sum:

# Approach

*Solution*: ban observation attacks.

$\overline{G}$

Zero-sum:

WSE → MPNE

General-sum:

WSE → Mutual Recursion

# Rollback Algorithm

Special Case: Action Attacks

# Rollback Algorithm

Special Case: Action Attacks

**1.** Victim determines Attacker's best response to any action $a$:

$$BR_h(s,a) = \arg\max_{a^\dagger \in \overline{\mathcal{A}}(s,a)} \left[ g_h(s,a,r_h(s,a)) + \mathbb{E}_{s' \sim P_h(s,a^\dagger)} \left[ V^*_{h+1,2}(s', \pi^*_{h+1}(s')) \right] \right]$$

# Rollback Algorithm

## Special Case: Action Attacks

**1. Victim determines Attacker's best response to any action $a$:**

$$BR_h(s,a) = \arg\max_{a^\dagger \in \overline{\mathcal{A}}(s,a)} \left[ g_h(s, a, r_h(s,a)) + \mathbb{E}_{s' \sim P_h(s,a^\dagger)} \left[ V^*_{h+1,2}(s', \pi^*_{h+1}(s')) \right] \right]$$

**2. Victim picks $a$ based on the worst-case best-response:**

$$V^*_{h,1}(s) = \max_{a \in \mathcal{A}} \min_{a^\dagger \in BR_h(s,a)} \left[ r_h(s, a^\dagger) + \mathbb{E}_{s' \sim P_h(s,a^\dagger)} \left[ V^*_{h+1,1}(s') \right] \right]$$

# Defense Results

# Defense Results

**Theorem**: *An optimal defense can be **computed** as the WSE of a meta game (POTBMG).*

# Defense Results

**Theorem**: *An optimal defense can be **computed** as the WSE of a meta game (POTBMG).*

*Moreover, the defense is computable in **polynomial time** if observation attacks are banned.*

# Defense Results

**Theorem**: *An optimal defense can be **computed** as the WSE of a meta game (POTBMG).*

*Moreover, the defense is computable in **polynomial time** if observation attacks are banned.*

*First results for the general defense problem!*

# Misinformation Attacks

*RLC 2024*

# Motivation

# Motivation

More **realistic** attacker: information advantage instead of environment control

# Motivation

More **realistic** attacker: information advantage instead of environment control

# Motivation

More **realistic** attacker: information advantage instead of environment control

# Motivation

More **realistic** attacker: information advantage instead of environment control
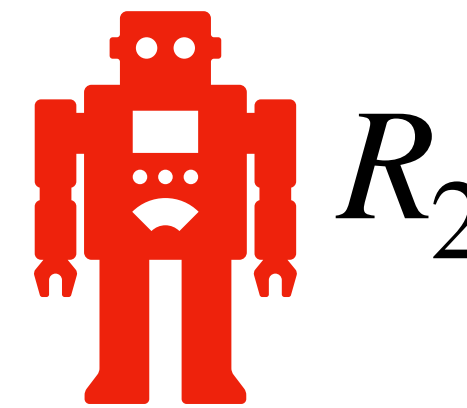
# Motivation

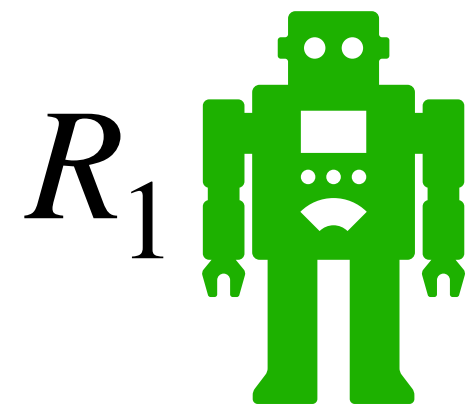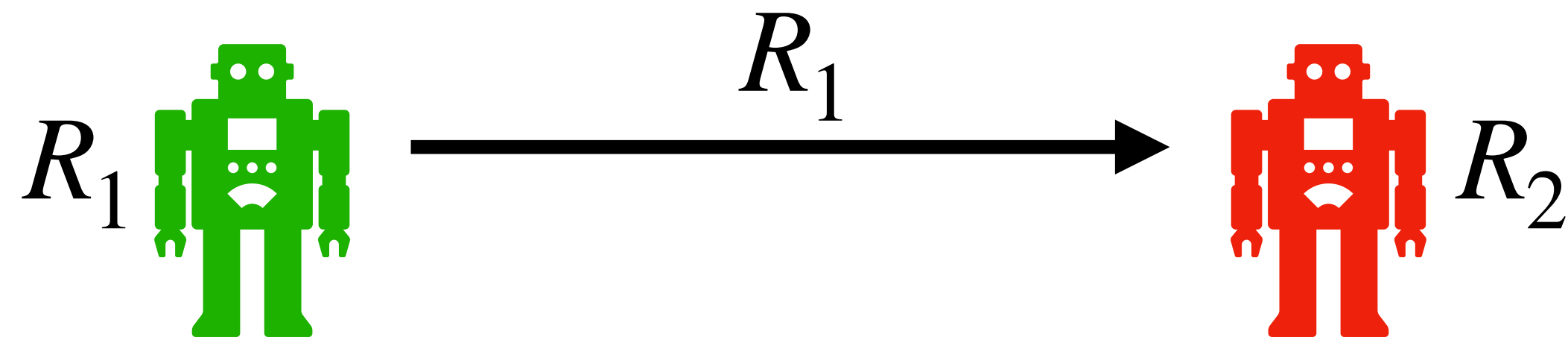More **realistic** attacker: information advantage instead of environment control

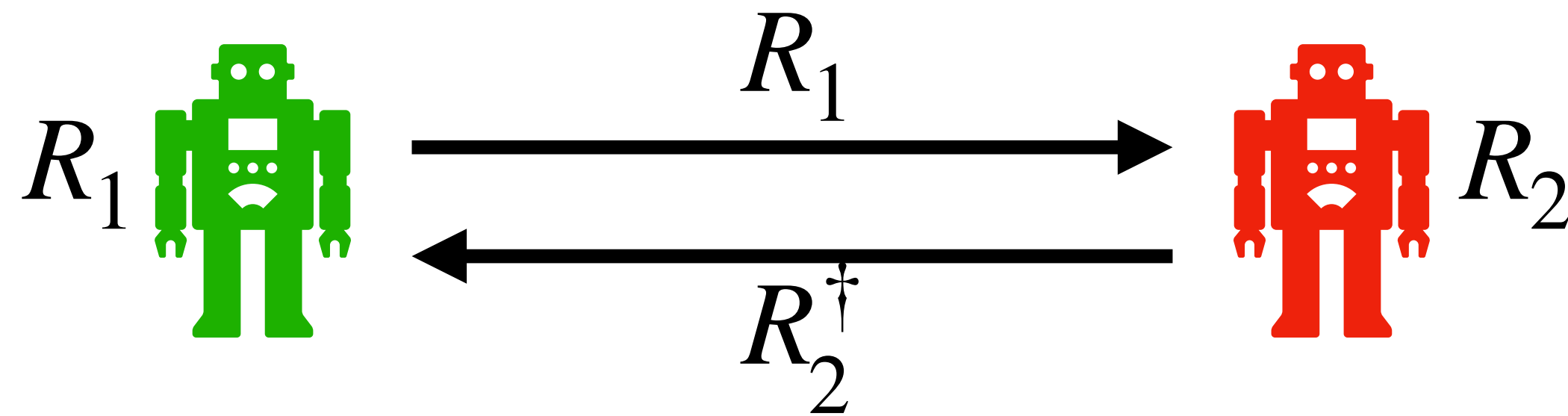# Motivation

More **realistic** attacker: information advantage instead of environment control
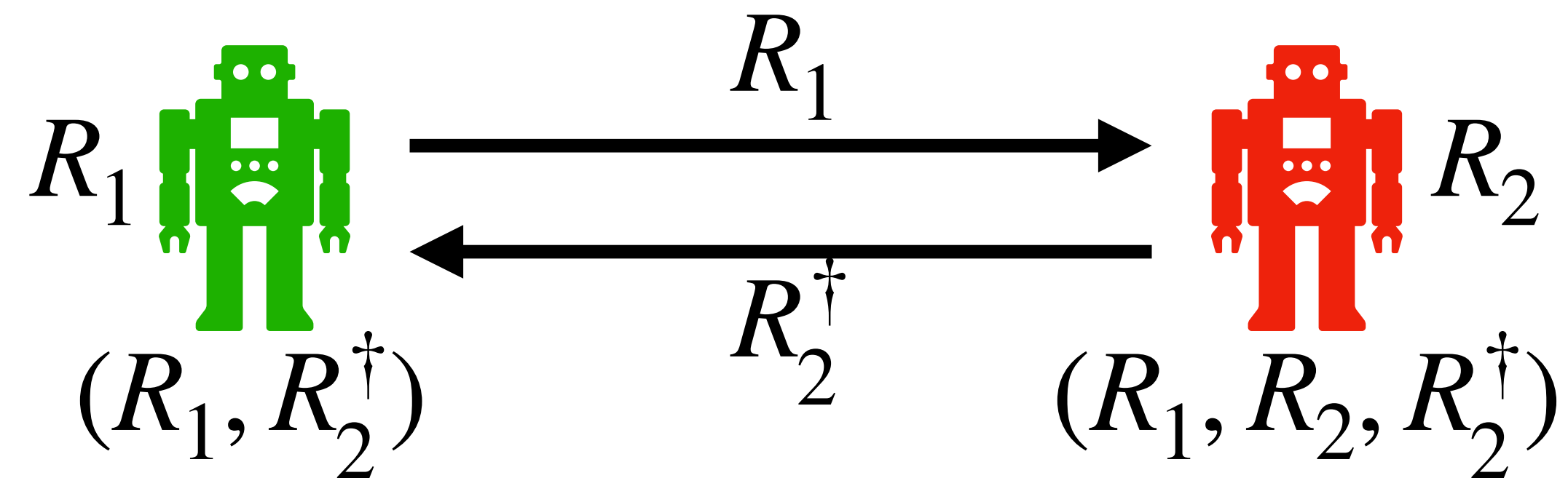
# Motivation

More **realistic** attacker: information advantage instead of environment control

# Motivation

More **realistic** attacker: information advantage instead of environment control

# Inception

# Inception

## Inception Problem

$$\max_{R_2^\dagger}$$ P2's best worst-case value given P1's beliefs about $R_2^\dagger$

# Inception

## Inception Problem

$$\max_{R_2^\dagger}$$ P2's best worst-case value given P1's beliefs about $R_2^\dagger$

Belief set: $\Pi_2^b(R_2^\dagger)$

# Inception

## Inception Problem

$$\max_{R_2^\dagger}$$ P2's best worst-case value given P1's beliefs about $R_2^\dagger$

Belief set: $\Pi_2^b(R_2^\dagger)$ —— P2 is "rational"

# Inception

## Inception Problem

$$\max_{R_2^\dagger}$$ P2's best worst-case value given P1's beliefs about $R_2^\dagger$

Belief set: $\Pi_2^b(R_2^\dagger)$ —— P2 is "rational"

$$\Pi_2^b(R_2^\dagger) = \left\{ \pi_2 \mid \exists R_2' \in \mathbb{B}_\epsilon(R_2^\dagger), (\cdot, \pi_2) \in SOL(R_1, R_2') \right\}$$

# Inception

## Inception Problem

$$\max_{R_2^\dagger} \max_{\pi_2^* \in \Pi_2} \min_{\pi_1^* \in \Pi_1^*} V_2^{\pi_1^*, \pi_2^*}$$

$$\text{s.t. } \Pi_1^* = \arg\max_{\pi_1 \in \Pi_1} \min_{\pi_2 \in \Pi_2^b(R_2^\dagger)} V_1^{\pi_1, \pi_2}$$

Belief set: $\Pi_2^b(R_2^\dagger)$ —— P2 is "rational"

$$\Pi_2^b(R_2^\dagger) = \left\{ \pi_2 \mid \exists R_2' \in \mathbb{B}_\epsilon(R_2^\dagger), (\,\cdot\,, \pi_2) \in SOL(R_1, R_2') \right\}$$

# Inception Approach

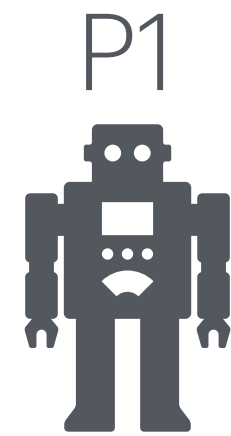# Inception Approach

P1

# Inception Approach

# Inception Approach

# Inception Approach

# Inception Approach

# Inception Approach

# Inception Approach

**Prediction**        Exploitation



Rational Belief        Linear Program

# Inception Approach

**Convincing**

P2

$$R_2^\dagger$$

Dom Strat

**Prediction**

P1

Best Response

Rational Belief

**Exploitation**

Best Response

Linear Program

# Inception Approach

# Inception Approach

**Convincing**

P2

$R_2^\dagger$

Dom Strat

iDSE

**Prediction**

P1

Best Response

Rational Belief

**Exploitation**

Best Response

Linear Program

Repeat to find the best **pure** strategy inception!

# Example: True Game

# Example: True Game

|     | L          | R          |
| --- | ---------- | ---------- |
| U   | 0, 5       | 1, 0       |
| D   | 1, $\epsilon$ | 0, 0    |
| S   | 1, 0       | 0, $\epsilon$ |

# Example: True Game

|   | L | R |
|---|---|---|
| U | 0, 5 | 1, 0 |
| D | 1, $\epsilon$ | 0, 0 |
| S | 1, 0 | 0, $\epsilon$ |

Unique NE

# Example: True Game

|   | L | R |
|---|---|---|
| U | 0, 5 | 1, 0 |
| D | 1, $\epsilon$ | 0, 0 |
| S | 1, 0 | 0, $\epsilon$ |

Unique NE

If P1 is rational, P2 gets 0!

# Example: True Game

|   | L | R |
|---|---|---|
| U | 0, 5 | 1, 0 |
| D | 1, $\epsilon$ | 0, 0 |
| S | 1, 0 | 0, $\epsilon$ |

P2 wants → (U, L: 0, 5)

Unique NE → (D, L: 1, $\epsilon$)

If P1 is rational, P2 gets 0!

# P2 fakes L

# P2 fakes L

|   | L | R |
|---|---|---|
| U | 0, 5 | 1, 0 |
| D | 1, $\epsilon$ | 0, 0 |
| S | 1, $2\epsilon$ | 0, $\epsilon$ |

Increased

# P2 fakes L

|   | L | R |
|---|---|---|
| U | 0, 5 | 1, 0 |
| D | 1, $\epsilon$ | 0, 0 |
| S | 1, $2\epsilon$ | 0, $\epsilon$ |

# P2 fakes L

|   | L | R |
|---|---|---|
| U | 0, 5 | 1, 0 |
| D | 1, $\epsilon$ | 0, 0 |
| S | 1, $2\epsilon$ | 0, $\epsilon$ |

P2

P1

# P2 fakes L

|     | L          | R          |
| --- | ---------- | ---------- |
| U   | 0, 5       | 1, 0       |
| D   | 1, $\epsilon$ | 0, 0    |
| S   | 1, $2\epsilon$ | 0, $\epsilon$ |

P2

P1

Tie!

# P2 fakes L



|   | L | R |
|---|---|---|
| U | 0, 5 | 1, 0 |
| D | 1, $\epsilon$ | 0, 0 |
| S | 1, 0 | 0, $\epsilon$ |

*Worst-Case Best Response*

# P2 fakes L



|   | L | R |
|---|---|---|
| U | 0, 5 | 1, 0 |
| D | 1, $\epsilon$ | 0, 0 |
| S | 1, 0 | 0, $\epsilon$ |

*Worst-Case Best Response*

P2 gets $\epsilon/2$ from (1/2,1/2) mix

# P2 fakes L



|   | L | R |
|---|---|---|
| U | 0, 5 | 1, 0 |
| D | 1, $\epsilon$ | 0, 0 |
| S | 1, 0 | 0, $\epsilon$ |

*Worst-Case Best Response*

P2 gets $\epsilon/2$ from (1/2,1/2) mix

Solved by Nash LP!

# P2 fakes R

# P2 fakes R

|   | L | R |
|---|---|---|
| U | 0, 5 | 1, 5+$\epsilon$ |
| D | 1, $\epsilon$ | 0, 2$\epsilon$ |
| S | 1, 0 | 0, $\epsilon$ |

Increased

# P2 fakes R

|   | L | R |
|---|---|---|
| U | 0, 5 | 1, 5+$\epsilon$ |
| D | 1, $\epsilon$ | 0, 2$\epsilon$ |
| S | 1, 0 | 0, $\epsilon$ |

Unique NE

# P2 fakes R

|   | L | R |
|---|---|---|
| U | 0, 5 | 1, 5+$\epsilon$ |
| D | 1, $\epsilon$ | 0, 2$\epsilon$ |
| S | 1, 0 | 0, $\epsilon$ |

Unique NE

P1 must play U!

# P2 fakes R

|   | L | R |
|---|---|---|
| U | 0, 5 | 1, 5+$\epsilon$ |
| D | 1, $\epsilon$ | 0, 2$\epsilon$ |
| S | 1, 0 | 0, $\epsilon$ |

P2 wins!

Unique NE

P1 must play U!

# P2 fakes R

## "Inception Attack"

# Exploitation

# Exploitation

**Assuming finite belief:** $\Pi_2^b(R_2^\dagger) = \{\pi_2^1, \ldots, \pi_2^K\}$

# Exploitation

**Assuming finite belief:** $\Pi_2^b(R_2^\dagger) = \{\pi_2^1, \ldots, \pi_2^K\}$

Complex

$$\max_{\pi_2^* \in \Pi_2} \min_{\pi_1^* \in \Pi_1^*} V_2^{\pi_1^*, \pi_2^*}$$

$$\text{s.t. } \Pi_1^* = \arg \max_{\pi_1 \in \Pi_1} \min_{\pi_2 \in \Pi_2^b(R_2^\dagger)} V_1^{\pi_1, \pi_2}$$

# Exploitation

**Assuming finite belief:** $\Pi_2^b(R_2^\dagger) = \{\pi_2^1, \ldots, \pi_2^K\}$

Complex

$$\max_{\pi_2^* \in \Pi_2} \min_{\pi_1^* \in \Pi_1^*} V_2^{\pi_1^*, \pi_2^*}$$

$$\text{s.t. } \Pi_1^* = \arg\max_{\pi_1 \in \Pi_1} \min_{\pi_2 \in \Pi_2^b(R_2^\dagger)} V_1^{\pi_1, \pi_2}$$

Duality

# Exploitation

**Assuming finite belief:** $\Pi_2^b(R_2^\dagger) = \{\pi_2^1, \ldots, \pi_2^K\}$

## Complex

$$\max_{\pi_2^* \in \Pi_2} \min_{\pi_1^* \in \Pi_1^*} V_2^{\pi_1^*, \pi_2^*}$$

$$\text{s.t. } \Pi_1^* = \arg\max_{\pi_1 \in \Pi_1} \min_{\pi_2 \in \Pi_2^b(R_2^\dagger)} V_1^{\pi_1, \pi_2}$$

Duality

## Linear

$$\max_{y \in \mathbb{R}^m, w \in \mathbb{R}^K, \alpha \in \mathbb{R}} z^* 1^\top w - \alpha$$

$$\text{s.t.} \quad \alpha + e_i^\top B y - e_i^\top A' w \geq 0 \quad \forall i \in [n]$$

$$1^\top y = 1, \quad y \geq 0 \quad w \geq 0.$$

# Exploitation

**Assuming finite belief:** $\Pi_2^b(R_2^\dagger) = \{\pi_2^1, \ldots, \pi_2^K\}$

### Complex

$$\max_{\pi_2^* \in \Pi_2} \min_{\pi_1^* \in \Pi_1^*} V_2^{\pi_1^*, \pi_2^*}$$

$$\text{s.t. } \Pi_1^* = \arg\max_{\pi_1 \in \Pi_1} \min_{\pi_2 \in \Pi_2^b(R_2^\dagger)} V_1^{\pi_1, \pi_2}$$

Duality

### Linear

$$\max_{y \in \mathbb{R}^m, w \in \mathbb{R}^K, \alpha \in \mathbb{R}} z^* 1^\top w - \alpha$$

$$\text{s.t.} \quad \alpha + e_i^\top B y - e_i^\top A' w \geq 0 \quad \forall i \in [n]$$

$$1^\top y = 1, \quad y \geq 0 \quad w \geq 0.$$

*Solve a sequence of LPs for MG case!*

# Results

# Results

**Theorem**: *rationality enables the **polynomial-time** computation of **misinformation attacks** that are optimal amongst the set of dominant-mixture reward functions.*

# Results

**Theorem**: *rationality enables the **polynomial-time** computation of **misinformation attacks** that are optimal amongst the set of dominant-mixture reward functions.*

*First efficient misinformation attacks on Markov games!*

# Constrained MARL

# Anytime Constraints

*AISTATS 2024*

# Motivation

# Motivation

# Motivation

# Motivation

# Motivation

# Motivation

# Motivation



$$\mathbb{P}_M^\pi \left[ \sum_{h=1}^{H} c_h \leq B \right] = 1$$

# Motivation



$$\mathbb{P}_M^{\pi}\left[\sum_{h=1}^{H} c_h \leq B\right] = 1$$

# Motivation



$$\mathbb{P}_M^{\pi}\left[\sum_{h=1}^{H} c_h \leq B\right] = 1$$

Cannot IOU a gas tank!

# Motivation



$$\mathbb{P}^{\pi}_{M}\left[\sum_{h=1}^{H} c_h \leq B\right] = 1$$

Cannot IOU a gas tank!

$$\mathbb{P}^{\pi}_{M}\left[\forall t \in [H], \sum_{h=1}^{t} c_h \leq B\right] = 1$$

# Motivation



❌    $\mathbb{P}_M^\pi \left[ \sum_{h=1}^{H} c_h \le B \right] = 1$    Cannot IOU a gas tank!

✅    $\mathbb{P}_M^\pi \left[ \forall t \in [H], \sum_{h=1}^{t} c_h \le B \right] = 1$

# Constrained Problem

# Constrained Problem

Agent's **goal** is to solve:

# Constrained Problem

Agent's **goal** is to solve:

$$\max_{\pi} \mathbb{E}_M^{\pi} \left[ \sum_{h=1}^{H} r_h(s_h, a_h) \right] \quad \text{s.t.} \quad \mathbb{P}_M^{\pi} \left[ \forall t \in [H], \ \sum_{h=1}^{t} c_h \leq B \right] = 1.$$

# Challenges

# Challenges

1. Feasible policies <span style="color:red">non-Markovian</span>

# Challenges

1. Feasible policies **non-Markovian**

2. Optimization is **NP-hard**

# Challenges

1. Feasible policies <span style="color:red">non-Markovian</span>

2. Optimization is <span style="color:red">NP-hard</span>

3. Determining feasibility of $\geq 2$ constraints is NP-hard
   $\implies$ <span style="color:red">Hardness of (value) Approximation</span>

# Reduction

# Reduction

*1. State-Cost
Augmentation*

# Reduction

*1. State-Cost Augmentation*

# Reduction



*1. State-Cost Augmentation*

# Reduction



*1. State-Cost Augmentation*

# Reduction



*1. State-Cost Augmentation*

# Reduction

*1. State-Cost Augmentation*

$$s$$

$$a$$

$$M$$

$$(s, \overline{c})$$

$$a'$$

$$\overline{M}$$

*2. BFS Generate Feasible Costs*

# Reduction

## 1. State-Cost Augmentation

$s$

$a$

$M$

$(s, \overline{c})$

$a'$

$\overline{M}$

## 2. BFS Generate Feasible Costs

$\overline{S}_1$

$(s_0, 0)$

# Reduction



1. State-Cost Augmentation

2. BFS Generate Feasible Costs

# Reduction

## 1. State-Cost Augmentation



## 2. BFS Generate Feasible Costs

# Reduction



*1. State-Cost Augmentation*

*2. BFS Generate Feasible Costs*

# Reduction

## 1. State-Cost Augmentation



## 2. BFS Generate Feasible Costs

# Reduction

# Exact Results

# Exact Results

$$\text{cost precision} \leq k \implies |\overline{S}| \leq SH2^{k+1}$$

# Exact Results

$$\text{cost precision} \leq k \implies |\overline{S}| \leq SH2^{k+1}$$

**Theorem (Fixed-Parameter Tractability):** *If the cost precision $k = O(\log(SAH))$, our algorithm outputs an optimal, anytime-constrained policy in polynomial time.*

# Approximate Feasibility

# Approximate Feasibility

**Definition 1** (Approximate Feasibility). For any $\epsilon > 0$, a policy $\pi$ is $\epsilon$-additive feasible if,

$$\mathbb{P}_M^\pi \left[ \forall k \in [H], \ \sum_{t=1}^k c_t \leq B + \epsilon \right] = 1,$$

and $\epsilon$-relative feasible if,

$$\mathbb{P}_M^\pi \left[ \forall k \in [H], \ \sum_{t=1}^k c_t \leq B(1 + \epsilon) \right] = 1.$$

# Approximation

# Approximation

*1. Truncate*

# Approximation

*1. Truncate* $\quad (\overline{c}, c)$

# Approximation

*1. Truncate* $\qquad (\overline{c}, c) \Rightarrow [B - Hc^{max}, B + 1]$

# Approximation

*1. Truncate*    $(\overline{c}, c)$ 🤖 ➡️ $[B - Hc^{max}, B + 1]$

*2. $\ell$-Discretize*

# Approximation

*1. Truncate*  $(\overline{c}, c)$ 🤖 ➡️ $[B - Hc^{max}, B + 1]$

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

*2. $\ell$-Discretize*  $(\overline{c}, c)$ 🤖

# Approximation

*1. Truncate*  $(\overline{c}, c)$  ➡️  $[B - Hc^{max}, B + 1]$

*2. $\ell$-Discretize*  $(\overline{c}, c)$  ➡️  $\left\lfloor \dfrac{\overline{c} + c}{\ell} \right\rfloor \ell$

# Approximation

*1. Truncate*   $(\bar{c}, c)$ 🤖 ➡️ $[B - Hc^{max}, B + 1]$

*2. $\ell$-Discretize*   $(\bar{c}, c)$ 🤖 ➡️ $\left\lfloor \dfrac{\bar{c} + c}{\ell} \right\rfloor \ell = \bar{c} + \left\lfloor \dfrac{c}{\ell} \right\rfloor \ell$

# Approximation Results

# Approximation Results

$$\ell = \frac{\epsilon}{H} \implies c \leq \hat{c} + \frac{\epsilon}{H} \implies \sum_h c_h \leq B + \epsilon$$

# Approximation Results

$$\ell = \frac{\epsilon}{H} \implies c \leq \hat{c} + \frac{\epsilon}{H} \implies \sum_h c_h \leq B + \epsilon$$

**Theorem (Approx):** *If $d$ is constant and $c^{max} \leq poly(|M|)$, our algorithm outputs an **optimal**-value, $\epsilon$-**feasible** policy in time $poly(|M|, \frac{1}{\epsilon})$*

*\*Guarantees are best-possible given hardness results.*

# Approximation Results

$$\ell = \frac{\epsilon}{H} \implies c \le \hat{c} + \frac{\epsilon}{H} \implies \sum_h c_h \le B + \epsilon$$

**Theorem (Approx):** *If $d$ is constant and $c^{max} \le poly(|M|)$, our algorithm outputs an **optimal**-value, $\epsilon$-**feasible** policy in time $poly(|M|, \frac{1}{\epsilon})$*

*First poly-time algorithm for anytime and almost sure constraints!*

*\*Guarantees are best-possible given hardness results.*

# Single-Constraint FPTAS

*NeurIPS 2024*

# Motivation

# Motivation

1. Previous approach cannot guarantee feasibility

# Motivation

1. Previous approach cannot guarantee feasibility

2. Only works for anytime constraints

# Motivation

1. Previous approach cannot guarantee feasibility

2. Only works for anytime constraints

# Packing Form

# Packing Form

$$\max_{\pi \in \Pi} \mathbb{E}_M^\pi \left[ \sum_{h=1}^{H} r_h(s_h, a_h) \right] \quad \text{s.t.} \quad \begin{cases} C_M^\pi \leq B \end{cases}$$

# Packing Form

$$\max_{\pi \in \Pi} \mathbb{E}^{\pi}_M \left[ \sum_{h=1}^{H} r_h(s_h, a_h) \right] \quad \text{s.t.} \quad \begin{cases} C^{\pi}_M \leq B \\ \pi \text{ deterministic} \end{cases}$$

# Packing Form

$$\max_{\pi \in \Pi} \mathbb{E}_M^\pi \left[ \sum_{h=1}^H r_h(s_h, a_h) \right] \quad \text{s.t.} \quad \begin{cases} C_M^\pi \leq B \\ \pi \text{ deterministic} \end{cases}$$

**Expectation:** $\quad C_M^\pi := \mathbb{E}_M^\pi \left[ \sum_{h=1}^H c_h \right]$

# Packing Form

$$\max_{\pi \in \Pi} \mathbb{E}_M^\pi \left[ \sum_{h=1}^H r_h(s_h, a_h) \right] \quad \text{s.t.} \quad \begin{cases} C_M^\pi \leq B \\ \pi \text{ deterministic} \end{cases}$$

**Expectation:** $\quad C_M^\pi := \mathbb{E}_M^\pi \left[ \sum_{h=1}^H c_h \right]$

**Anytime:** $\quad C_M^\pi := \max_t \max_{\tau : \mathbb{P}^\pi[\tau] > 0} \sum_{h=1}^t c_h$

# Duality

# Duality

Packing (Primal)

$$\max_{\pi \in \Pi^D} \quad V_M^\pi$$

$$\text{s.t.} \quad C_M^\pi \leq B$$

# Duality

Packing (Primal)

$$\max_{\pi \in \Pi^D} \quad V_M^\pi$$

$$\text{s.t.} \quad C_M^\pi \leq B$$

*Optimum value, but approximate cost*

# Duality

Packing (Primal)

$$V^* \text{---} \begin{array}{l} \max\limits_{\pi \in \Pi^D} \quad V_M^\pi \\[1em] \text{s.t.} \quad C_M^\pi \leq B \end{array}$$

*Optimum value, but approximate cost*

# Duality

Packing (Primal)

$$V^* — \begin{array}{ll} \max\limits_{\pi \in \Pi^D} & V_M^\pi \\[1em] \text{s.t.} & C_M^\pi \leq B \end{array}$$

Covering (Dual)

$$\begin{array}{ll} \min\limits_{\pi \in \Pi^D} & C_M^\pi \\[1em] \text{s.t.} & V_M^\pi \geq V^* \end{array}$$

*Optimum value, but approximate cost*

# Duality

Packing (Primal)

$$V^* \text{---} \begin{aligned} \max_{\pi \in \Pi^D} \quad & V_M^\pi \\ \text{s.t.} \quad & C_M^\pi \leq B \end{aligned}$$

Covering (Dual)

$$\begin{aligned} \min_{\pi \in \Pi^D} \quad & C_M^\pi \\ \text{s.t.} \quad & V_M^\pi \geq V^* \end{aligned}$$

*Optimum value, but approximate cost*

*Optimum cost, but approximate value*

# Duality

Packing (Primal)

$$V^* \text{---} \begin{array}{ll} \max\limits_{\pi \in \Pi^D} & V_M^\pi \\[2mm] \text{s.t.} & C_M^\pi \leq B \end{array}$$

Covering (Dual)

$$\begin{array}{ll} \min\limits_{\pi \in \Pi^D} & C_M^\pi \\[2mm] \text{s.t.} & V_M^\pi \geq V^* \end{array}$$

*Optimum value, but approximate cost*

*Optimum cost, but approximate value*

**Feasible!**

# Value-Demand Augmentation

# Value-Demand Augmentation

**Intuition**: Build $\overline{M}$ satisfying,

# Value-Demand Augmentation

**Intuition**: Build $\overline{M}$ satisfying,

$$\overline{C}_h^*(s, v) = \min_{\pi \in \Pi^D} \quad C_h^\pi(\tau_h)$$
$$\text{s.t.} \quad V_h^\pi(\tau_h) \geq v$$

# Value-Demand Augmentation

**Intuition**: Build $\overline{M}$ satisfying,

$$\overline{C}_h^*(s,v) = \min_{\pi \in \Pi^D} \quad C_h^\pi(\tau_h)$$
$$\text{s.t.} \quad V_h^\pi(\tau_h) \geq v$$

$\Rightarrow$ Dual $= \overline{C}_1^*(s_0, V^*)$

# Value-Demand Augmentation

**Intuition**: Build $\overline{M}$ satisfying,

$$\overline{C}_h^*(s, v) = \min_{\pi \in \Pi^D} \quad C_h^\pi(\tau_h)$$
$$\text{s.t.} \quad V_h^\pi(\tau_h) \geq v$$

$\Rightarrow$ Dual $= \overline{C}_1^*(s_0, V^*)$

# Value-Demand Augmentation

**Intuition**: Build $\overline{M}$ satisfying,

$$\overline{C}_h^*(s,v) = \min_{\pi \in \Pi^D} \quad C_h^\pi(\tau_h)$$
$$\text{s.t.} \quad V_h^\pi(\tau_h) \geq v$$

$\Rightarrow$ Dual $= \overline{C}_1^*(s_0, V^*)$

$\overline{S} = S \times \mathcal{V}$ — *all possible values*

# Value-Demand Augmentation

**Intuition**: Build $\overline{M}$ satisfying,

$$\overline{C}_h^*(s, v) = \min_{\pi \in \Pi^D} \quad C_h^\pi(\tau_h)$$
$$\text{s.t.} \quad V_h^\pi(\tau_h) \geq v$$

$\Longrightarrow$ Dual $= \overline{C}_1^*(s_0, V^*)$

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

$\overline{S} = S \times \mathcal{V}$ —— *all possible values*

**Invariant**: $v \leq \overline{V}_h^\pi(s, v)$

# Value-Demand Augmentation

**Intuition**: Build $\overline{M}$ satisfying,

$$\overline{C}_h^*(s, v) = \min_{\pi \in \Pi^D} \quad C_h^\pi(\tau_h)$$

$$\text{s.t.} \quad V_h^\pi(\tau_h) \geq v$$

$\Rightarrow$  Dual $= \overline{C}_1^*(s_0, V*)$

---

$\overline{S} = S \times \mathcal{V}$ — *all possible values*

**Invariant:** $v \leq \overline{V}_h^\pi(s, v) = r_h(s, a) + \sum_{s'} P_h(s' \mid s, a) \overline{V}_{h+1}^\pi(s', v_{s'})$  **PE**

# Value-Demand Augmentation

**Intuition**: Build $\overline{M}$ satisfying,

$$\overline{C}_h^*(s,v) = \min_{\pi \in \Pi^D} \quad C_h^\pi(\tau_h)$$
$$\text{s.t.} \quad V_h^\pi(\tau_h) \geq v$$

$\Rightarrow$ Dual $= \overline{C}_1^*(s_0, V^*)$

$\overline{S} = S \times \mathcal{V}$ — *all possible values*

**Invariant:** $v \leq \overline{V}_h^\pi(s,v) = r_h(s,a) + \sum_{s'} P_h(s' \mid s,a) \overline{V}_{h+1}^\pi(s', v_{s'})$   **PE**

$$\overline{\mathcal{A}}_h(s,v) := \left\{ (a, \mathbf{v}) \in \mathcal{A} \times \mathcal{V}^S \mid r_h(s,a) + \sum_{s'} P_h(s' \mid s,a) v_{s'} \geq v \right\}$$

# Outer Algorithm

# Outer Algorithm

*1. Solve:*

# Outer Algorithm

*1. Solve:*

$$\overline{C}_h^*(s,v) = \min_{a,\mathbf{v} \in \mathcal{A}_h(s,v)} c_h(s,a) + \overbrace{\max_{s'} \overline{C}_h^*(s',v_{s'})}^{\textit{Anytime Constraints}}$$

# Outer Algorithm

*1. Solve:* $\displaystyle \overline{C}_h^*(s,v) = \min_{a,\mathbf{v} \in \mathcal{A}_h(s,v)} c_h(s,a) + \sum_{s'} P_h(s' \mid s,a) \underbrace{\overline{C}_{h+1}^*(s',v_{s'})}_{\textit{Expectation Constraints}}$

# Outer Algorithm

*1. Solve:* $\quad \overline{C}_h^*(s, v) = \min_{a, \mathbf{v} \in \mathcal{A}_h(s,v)} c_h(s, a) + \sum_{s'} P_h(s' \mid s, a) \overline{C}_{h+1}^*(s', v_{s'})$

---

*2. Output:*

# Outer Algorithm

1. *Solve:* $\quad \overline{C}_h^*(s, v) = \min\limits_{a, \mathbf{v} \in \mathcal{A}_h(s,v)} c_h(s, a) + \sum\limits_{s'} P_h(s' \mid s, a) \overline{C}_{h+1}^*(s', v_{s'})$

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

2. *Output:* $\quad V_M^* = \max \left\{ v \in \mathcal{V} \mid \overline{C}_1^*(s_0, v) \leq B \right\}$

# Outer Algorithm

*1. Solve:*
$$\overline{C}_h^*(s, v) = \min_{a, \mathbf{v} \in \mathcal{A}_h(s,v)} c_h(s, a) + \sum_{s'} P_h(s' \mid s, a) \overline{C}_{h+1}^*(s', v_{s'})$$

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

*2. Output:*
$$V_M^* = \max \left\{ v \in \mathcal{V} \mid \overline{C}_1^*(s_0, v) \le B \right\}$$

**Feasible!**

# Outer Algorithm

*1. Solve:*
$$\overline{C}_h^*(s,v) = \min_{\boxed{a,\mathbf{v}\in\mathcal{A}_h(s,v)}} c_h(s,a) + \sum_{s'} P_h(s' \mid s,a)\overline{C}_{h+1}^*(s',v_{s'})$$

*Exponential!*

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

*2. Output:*
$$V_M^* = \max\left\{v \in \mathcal{V} \mid \overline{C}_1^*(s_0,v) \le B\right\}$$

**Feasible!**

# Solving $\overline{M}$ Fast

# Solving $\overline{M}$ Fast

**Optimality Equations**

$$\overline{C}_h^*(s, v) = \min_{(a, \mathbf{v})} c_h(s, a) + \sum_{s'} P_h(s' \mid s, a) \overline{C}_h^*(s, v_{s'})$$

$$\text{s.t. } r_h(s, a) + \sum_{s'} P_h(s' \mid s, a) v_{s'} \geq v$$

# Solving $\overline{M}$ Fast

**Optimality Equations**

$$\min_{\mathbf{v} \in \mathcal{V}^S} \sum_{s'} P_h(s' \mid s, a) \overline{C}_h^*(s, v_{s'})$$

$$\sum_{s'} P_h(s' \mid s, a) v_{s'} \qquad \geq v$$

# Solving $\overline{M}$ Fast

## Optimality Equations

$$\min_{\mathbf{v} \in \mathcal{V}^S} \sum_{s'} P_h(s' \mid s, a) \overline{C}_h^*(s, v_{s'})$$

$$\sum_{s'} \quad p_{s'} \qquad v_{s'} \qquad \geq v$$

# Solving $\overline{M}$ Fast

**Optimality Equations**

$$\min_{\mathbf{v} \in \mathcal{V}^S} \sum_{s'} w_{s'} \quad \overline{C}_h^*(s, v_{s'})$$

$$\sum_{s'} p_{s'} \quad v_{s'} \quad \geq v$$

# Solving $\overline{M}$ Fast

## Optimality Equations

$$\min_{\mathbf{v} \in \mathcal{V}^S} \sum_{s'} \quad w_{s'} \quad f(v_{s'})$$

$$\sum_{s'} \quad p_{s'} \quad v_{s'} \quad \geq v$$

# Solving $\overline{M}$ Fast

**Optimality Equations**

$$\min_{\mathbf{v} \in \mathcal{V}^S} \sum_{s'} w_{s'} \, f(v_{s'})$$

$$\sum_{s'} p_{s'} \quad v_{s'} \qquad \geq v$$

$\longleftrightarrow$

**Knapsack Problem**

$$\min_{x \in X^n} \sum_i w_i x_i$$

$$\text{s.t.} \sum_i p_i x_i \geq P$$

# Solving $\overline{M}$ Fast

**Optimality Equations**

$$\min_{\mathbf{v} \in \mathcal{V}^S} \sum_{s'} w_{s'} \quad f(v_{s'})$$

$$\sum_{s'} p_{s'} \quad v_{s'} \quad \geq v$$

$\longleftrightarrow$

**Knapsack Problem**

$$\min_{x \in X^n} \sum_i w_i x_i$$

$$\text{s.t. } \sum_i p_i x_i \geq P$$

**Knapsack Approx!**

$$MC(s', p)$$

# Solving $\overline{M}$ Fast

**Optimality Equations**

$$\min_{\mathbf{v} \in \mathcal{V}^S} \sum_{s'} w_{s'} \quad f(v_{s'})$$

$$\sum_{s'} p_{s'} \quad v_{s'} \quad \geq v$$

**Knapsack Problem**

$$\min_{x \in X^n} \sum_i w_i x_i$$

$$\text{s.t.} \sum_i p_i x_i \geq P$$

**Knapsack Approx!**

$$\sum_{i=1}^{s'-1} p_i v_i$$

$$MC(s', p)$$

# Solving $\overline{M}$ Fast

**Optimality Equations**

$$\min_{\mathbf{v} \in \mathcal{V}^S} \sum_{s'} w_{s'} \quad f(v_{s'})$$

$$\sum_{s'} p_{s'} \quad v_{s'} \quad \geq v$$

$\longleftrightarrow$

**Knapsack Problem**

$$\min_{x \in X^n} \sum_i w_i x_i$$

$$\text{s.t.} \sum_i p_i x_i \geq P$$

**Knapsack Approx!**

$$\sum_{i=1}^{s'-1} p_i v_i$$

$$MC(s', p) = \min_{v_{s'} \in \mathcal{V}} \quad w_{s'} \quad f(v_{s'})$$

# Solving $\overline{M}$ Fast

**Optimality Equations**

$$\min_{\mathbf{v} \in \mathcal{V}^S} \sum_{s'} w_{s'} \quad f(v_{s'})$$

$$\sum_{s'} p_{s'} \quad v_{s'} \qquad \geq v$$

**Knapsack Problem**

$$\min_{x \in X^n} \sum_i w_i x_i$$

$$\text{s.t.} \sum_i p_i x_i \geq P$$

**Knapsack Approx!**

$$\sum_{i=1}^{s'-1} p_i v_i$$

$$MC(s', p) = \min_{v_{s'} \in \mathcal{V}} w_{s'} \quad f(v_{s'}) \quad + MC(s'+1, p + \quad p_{s'} \quad v_{s'})$$

# Solving $\overline{M}$ Fast

**Optimality Equations**

$$\min_{\mathbf{v} \in \mathcal{V}^S} \sum_{s'} w_{s'} \quad f(v_{s'})$$

$$\sum_{s'} p_{s'} \quad v_{s'} \quad \geq v$$

**Knapsack Problem**

$$\min_{x \in X^n} \sum_i w_i x_i$$

$$\text{s.t.} \sum_i p_i x_i \geq P$$

**Knapsack Approx!**

$$\sum_{i=1}^{s'-1} p_i v_i$$

$$\sum_{i=1}^{s'} p_i v_i$$

$$MC(s', p) = \min_{v_{s'} \in \mathcal{V}} \quad w_{s'} \quad f(v_{s'}) \quad + MC(s'+1, p + \quad p_{s'} \quad v_{s'})$$

# Solving $\overline{M}$ Fast

**Optimality Equations**

$$\min_{\mathbf{v} \in \mathcal{V}^S} \sum_{s'} w_{s'} \quad f(v_{s'})$$

$$\sum_{s'} p_{s'} \quad v_{s'} \quad \geq v$$

**Knapsack Problem**

$$\min_{x \in X^n} \sum_i w_i x_i$$

$$\text{s.t.} \sum_i p_i x_i \geq P$$

**Knapsack Approx!**

$$MC(s', p) = \min_{v_{s'} \in \mathcal{V}} \quad w_{s'} \quad f(v_{s'}) \quad + MC(s' + 1, p + \quad p_{s'} \quad v_{s'})$$
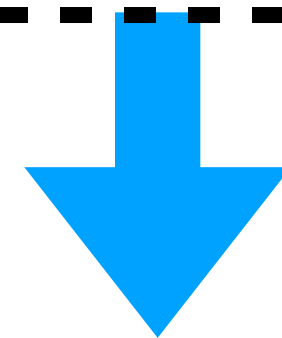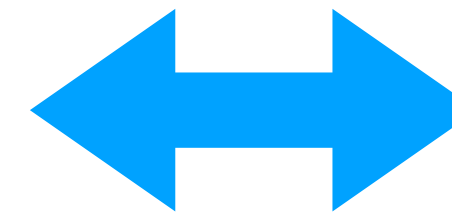
Round for approx

# Solving $\overline{M}$ Fast

**Optimality Equations**

$$\min_{\mathbf{v} \in \mathcal{V}^S} \sum_{s'} \quad w_{s'} \quad f(v_{s'})$$

$$\sum_{s'} \quad p_{s'} \quad v_{s'} \quad \geq v$$

$\longleftrightarrow$

**Knapsack Problem**

$$\min_{x \in X^n} \sum_i w_i x_i$$

$$\text{s.t.} \sum_i p_i x_i \geq P$$

**Knapsack Approx!**

$$MC(s', p) = \min_{v_{s'} \in \mathcal{V}} P_h(s' \mid s, a) \overline{C}^*_{h+1}(s', v_{s'}) + MC(s' + 1, p + P_h(s' \mid s, a) v_{s'})$$

Round for approx

# Time-Space Rounding

# Time-Space Rounding

Round $v's$ down $\implies$ cost goes down!

# Time-Space Rounding

Round $v's$ down $\implies$ cost goes down!

**Feasible!**

# Time-Space Rounding

Round $v's$ down $\implies$ cost goes down!

**Feasible!**

Rounding $v's$ causes error over time

# Time-Space Rounding

Round $v's$ down $\implies$ cost goes down!

**Feasible!**

Rounding $v's$ causes error over **time**

$\hat{V}^{\pi} \geq V^* - \ell H$

# Time-Space Rounding

Round $v's$ down $\implies$ cost goes down!

**Feasible!**

Rounding $v's$ causes error over <span style="color:orange">time</span>

$\hat{V}^\pi \geq V^* - \ell H$

Rounding $p's$ causes error over <span style="color:red">space</span>

# Time-Space Rounding

Round $v's$ down $\implies$ cost goes down! **Feasible!**

Rounding $v's$ causes error over <span style="color:orange">time</span> $\implies$ $\hat{V}^\pi \geq V^* - \ell H$

Rounding $p's$ causes error over <span style="color:red">space</span> $\implies$ $\hat{V}^\pi \geq V^* - \ell SH$

# Time-Space Rounding

Round $v's$ down $\Longrightarrow$ cost goes down!

**Feasible!**

Rounding $v's$ causes error over **time** $\Longrightarrow$ $\hat{V}^\pi \geq V^* - \ell H$

Rounding $p's$ causes error over **space** $\Longrightarrow$ $\hat{V}^\pi \geq V^* - \ell SH$

$$\ell = \frac{\epsilon}{SH} \implies \hat{V}^\pi \geq V^* - \epsilon$$

# Results

# Results

**Theorem (FPTAS)**: *If the rewards are poly-bounded, our algorithm outputs a **feasible** policy with value $V^* - \epsilon$ in time $poly(|M|, \frac{1}{\epsilon})$*

*Guarantees are best-possible given hardness results.*

# Results

**Theorem (FPTAS)**: *If the rewards are poly-bounded, our algorithm outputs a **feasible** policy with value $V* - \epsilon$ in time $poly(|M|, \frac{1}{\epsilon})$*

*First ever poly-time algorithm for **deterministic**, expectation-constrained policies!*

*\*Guarantees are best-possible given hardness results.*

# Multi-Constraint Bicriteria

*ICML 2025*

# Motivation

# Motivation

**Full Problem**

$$\max_{\pi \in \Pi^D} \quad V^\pi$$

$$\text{s.t.} \quad C_1^\pi \leq B_1$$

$$C_2^\pi \leq B_2$$

$$\vdots$$

$$C_m^\pi \leq B_m$$

# Motivation

## Full Problem

$$\max_{\pi \in \Pi^D} \quad V^\pi$$

$$\text{s.t.} \quad C_1^\pi \leq B_1$$

$$C_2^\pi \leq B_2$$

$$\vdots$$

$$C_m^\pi \leq B_m$$

**Expectation:** $\mathbb{E}_M^\pi \left[ \sum_{h=1}^H c_h \right] \leq B$

**Chance:** $\mathbb{P}_M^\pi \left[ \sum_{h=1}^H c_h > B \right] \leq \delta$

**Almost Sure:** $\mathbb{P}_M^\pi \left[ \sum_{h=1}^H c_h \leq B \right] = 1$

**Anytime:** $\mathbb{P}_M^\pi \left[ \forall t, \ \sum_{h=1}^t c_h \leq B \right] = 1$

# Motivation

### Full Problem

$$\max_{\pi \in \Pi^D} \quad V^\pi$$

$$\text{s.t.} \quad C_1^\pi \leq B_1$$

$$C_2^\pi \leq B_2$$

$$\vdots$$

$$C_m^\pi \leq B_m$$

**Expectation:** $\quad \mathbb{E}_M^\pi \left[ \sum_{h=1}^{H} c_h \right] \leq B$

**Chance:** $\quad \mathbb{P}_M^\pi \left[ \sum_{h=1}^{H} c_h > B \right] \leq \delta$

**Almost Sure:** $\quad \mathbb{P}_M^\pi \left[ \sum_{h=1}^{H} c_h \leq B \right] = 1$

**Anytime:** $\quad \mathbb{P}_M^\pi \left[ \forall t, \ \sum_{h=1}^{t} c_h \leq B \right] = 1$

*Can we create a framework that works for **any combination** of constraints?*

# Budget Augmentation

# Budget Augmentation

**Full Form**

$$\max_{\pi \in \Pi^D} \quad V^\pi$$

$$\text{s.t.} \quad C_1^\pi \leq B_1$$

$$C_2^\pi \leq B_2$$

$$\vdots$$

$$C_m^\pi \leq B_m$$

# Budget Augmentation

**Full Form**

$$\max_{\pi \in \Pi^D} \quad V^\pi$$

$$\text{s.t.} \quad C_1^\pi \leq B_1$$

$$C_2^\pi \leq B_2$$

$$\vdots$$

$$C_m^\pi \leq B_m$$

**Intuition**: Build $\overline{M}$ satisfying,

$$\overline{V}_h^*(s, b) = \max_{\pi \in \Pi^D} \quad V_h^\pi(\tau_h)$$

$$\text{s.t.} \quad C_h^\pi(\tau_h) \leq b$$

# Budget Augmentation

**Full Form**

$$\max_{\pi \in \Pi^D} \quad V^\pi$$
$$\text{s.t.} \quad C_1^\pi \leq B_1$$
$$C_2^\pi \leq B_2$$
$$\vdots$$
$$C_m^\pi \leq B_m$$

**Intuition**: Build $\overline{M}$ satisfying,

$$\overline{V}_h^*(s, b) = \max_{\pi \in \Pi^D} \quad V_h^\pi(\tau_h)$$
$$\text{s.t.} \quad C_h^\pi(\tau_h) \leq b$$

Primal $= \overline{V}_1^*(s_0, B)$

# Budget Augmentation

**Full Form**

$$\max_{\pi \in \Pi^D} \quad V^\pi$$

$$\text{s.t.} \quad C_1^\pi \leq B_1$$

$$C_2^\pi \leq B_2$$

$$\vdots$$

$$C_m^\pi \leq B_m$$

**Intuition**: Build $\overline{M}$ satisfying,

$$\overline{V}_h^*(s, b) = \max_{\pi \in \Pi^D} \quad V_h^\pi(\tau_h)$$

$$\text{s.t.} \quad C_h^\pi(\tau_h) \leq b$$

Primal $= \overline{V}_1^*(s_0, B)$

*Use previous approach but with rounding up!*

# Constraint Assumptions

# Constraint Assumptions

*1. Recursion:*

# Constraint Assumptions

*1. Recursion:*  $C_h^\pi(\tau_h) = c_h(s,a) + f_{s'} g(P_h(s' \mid s,a)) C_{h+1}^\pi(s')$

# Constraint Assumptions

*1. Recursion:*    $C_h^\pi(\tau_h) = c_h(s, a) + f_{s'} g(P_h(s' \mid s, a)) C_{h+1}^\pi(s')$

Required for inner DP

# Constraint Assumptions

*1. Recursion:*  $C_h^\pi(\tau_h) = c_h(s, a) + f_{s'} g(P_h(s' \mid s, a)) C_{h+1}^\pi(s')$

| | Exp | AS |
|---|---|---|
| $f$ | $\sum_{s'}$ | $\max_{s'}$ |
| $g$ | $id$ | $[x > 0]$ |

Required for inner DP

# Constraint Assumptions

*1. Recursion:*
$$C_h^\pi(\tau_h) = c_h(s,a) + f_{s'} g(P_h(s' \mid s,a)) C_{h+1}^\pi(s')$$

Required for inner DP

|       | Exp           | AS          |
| ----- | ------------- | ----------- |
| $f$   | $\sum_{s'}$   | $\max_{s'}$ |
| $g$   | $id$          | $[x > 0]$   |

# Constraint Assumptions

*1. Recursion:*

$$C_h^\pi(\tau_h) = c_h(s,a) + f_{s'} g(P_h(s' \mid s,a)) C_{h+1}^\pi(s')$$

Required for inner DP

| | Exp | AS |
|---|---|---|
| $f$ | $\sum\limits_{s'}$ | $\max\limits_{s'}$ |
| $g$ | $id$ | $[x > 0]$ |

*2. 1-Lipschitz:*

# Constraint Assumptions

| | Exp | AS |
|---|---|---|
| $f$ | $\displaystyle\sum_{s'}$ | $\displaystyle\max_{s'}$ |
| $g$ | $id$ | $[x > 0]$ |

*1. Recursion:*  $\qquad C_h^\pi(\tau_h) = c_h(s,a) + f_{s'}g(P_h(s' \mid s,a))C_{h+1}^\pi(s')$

Required for inner DP

*2. 1-Lipschitz:* $\qquad f(x, \mathrm{round}(y)) \leq f(x, y + \ell) \leq f(x, y) + \ell$

# Constraint Assumptions

*1. Recursion:*
$$C_h^\pi(\tau_h) = c_h(s,a) + f_{s'}g(P_h(s' \mid s,a))C_{h+1}^\pi(s')$$

Required for inner DP

|     | Exp | AS |
| --- | --- | --- |
| $f$ | $\sum\limits_{s'}$ | $\max\limits_{s'}$ |
| $g$ | $id$ | $[x > 0]$ |

*2. 1-Lipschitz:*
$$f(x, \mathrm{round}(y)) \le f(x, y + \ell) \le f(x, y) + \ell$$

Required for rounding error analysis

# Results

# Results

**Theorem (Bicriteria):** *Our algorithm computes an **optimal**-value, $\epsilon$-**feasible** policy in **polynomial time**, so long as the costs are poly-bounded and satisfy the SR condition.*

*Guarantees are best-possible given hardness results.*

# Results

**Theorem (Bicriteria):** *Our algorithm computes an **optimal**-value, $\epsilon$-**feasible** policy in **polynomial time**, so long as the costs are poly-bounded and satisfy the SR condition.*

*Includes **all** classical constraints!*

*\*Guarantees are best-possible given hardness results.*

# Results

**Theorem (Bicriteria):** *Our algorithm computes an **optimal**-value, $\epsilon$-**feasible** policy in **polynomial time**, so long as the costs are poly-bounded and satisfy the SR condition.*

*Includes **all** classical constraints!*

*First ever poly-time algorithm for **chance** constraints and **non-homogenous** constraints!*

*\*Guarantees are best-possible given hardness results.*

# Future Directions

1. Beyond Worst-case Analysis for all works
   (especially POMDPs for defense and anytime constraints)

2. Submodular Constrained Reinforcement Learning

3. Optimal learning under constraints.

# Thank you!

# Backup

# Motivating Example

# Motivating Example

# Motivating Example

1. Robust to visual noise (ash)

# Motivating Example

1. Robust to visual noise (ash)

2. Robust to other rescue vehicles

# Motivating Example

1. Robust to visual noise (ash)

2. Robust to other rescue vehicles

3. Coordinate well with teammates

# Motivating Example

1. Robust to visual noise (ash)

2. Robust to other rescue vehicles

3. Coordinate well with teammates



1. Effective fuel management

# Motivating Example

1. Robust to visual noise (ash)

2. Robust to other rescue vehicles

3. Coordinate well with teammates



1. Effective fuel management

2. Avoids dangerous terrain (lava)

# Motivating Example



1. Robust to visual noise (ash)

2. Robust to other rescue vehicles

3. Coordinate well with teammates

1. Effective fuel management

2. Avoids dangerous terrain (lava)

3. Balances risks of difficult terrain

# Framework Extensions

# Framework Extensions

1. Multiple agents

# Framework Extensions

1. Multiple agents

2. Infinite discounting

# Framework Extensions

1. Multiple agents

2. Infinite discounting

3. Stochastic costs

# Framework Extensions

1. Multiple agents

2. Infinite discounting

3. Stochastic costs

   1. Discrete

# Framework Extensions

1. Multiple agents

2. Infinite discounting

3. Stochastic costs

   1. Discrete

   2. Bounded Continuous

# Framework Extensions

1. Multiple agents

2. Infinite discounting

3. Stochastic costs

    1. Discrete

    2. Bounded Continuous

4. Continuous States

# Chance Constraints

# Chance Constraints

1. Use Discretized $\hat{M}$ from anytime constraints section

# Chance Constraints

1. Use Discretized $\hat{M}$ from anytime constraints section

2. Define $C_h^\pi(s, \bar{c}) = \mathbb{P}^\pi \left[ \exists k, \bar{c} + \sum_{t=h}^{k} c_t > B \right]$

# Chance Constraints

1. Use Discretized $\hat{M}$ from anytime constraints section

2. Define $C_h^\pi(s, \bar{c}) = \mathbb{P}^\pi \left[ \exists k, \bar{c} + \sum_{t=h}^{k} c_t > B \right]$ satisfies,

$$C_h^\pi(s, \bar{c}) = [\bar{c} + c_h(s, a) > B] + \sum_{s'} P_h(s' \mid s, a) C_{h+1}^\pi(s, \bar{c} + c_h(s, a))$$

# Chance Constraints

1. Use Discretized $\hat{M}$ from anytime constraints section

2. Define $\quad C_h^\pi(s, \bar{c}) = \mathbb{P}^\pi \left[ \exists k, \bar{c} + \sum_{t=h}^k c_t > B \right] \quad$ satisfies,

$$C_h^\pi(s, \bar{c}) = \underbrace{[\bar{c} + c_h(s, a) > B]}_{\text{New } c_h'((s, \bar{c}), a)} + \sum_{s'} P_h(s' \mid s, a) C_{h+1}^\pi(s, \bar{c} + c_h(s, a))$$

# Action Space

# Action Space

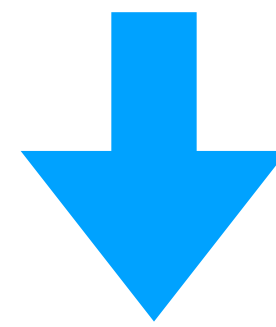**Policy Evaluation Equation:** $C_h^\pi(s) = c_h(s, a) + \sum_{s'} P_h(s' \mid s, a) C_{h+1}^\pi(s')$

# Action Space

**Policy Evaluation Equation:** $C_h^\pi(s) = c_h(s, a) + \sum_{s'} P_h(s' \mid s, a) C_{h+1}^\pi(s')$

*Same form as before!*

$$\overline{\mathcal{A}}_h(s, b) := \left\{ (a, \mathbf{b}) \in \mathcal{A} \times \mathbb{R}^S \,\middle|\, c_h(s, a) + \sum_{s'} P_h(s' \mid s, a) b_{s'} \leq b \right\}$$

# Budget Augmentation

# Budget Augmentation

**Intuition**: Build $\overline{M}$ satisfying,

# Budget Augmentation

**Intuition**: Build $\overline{M}$ satisfying,

$$\overline{V}^*_h(s, b) = \max_{\pi \in \Pi^D} \quad V^\pi_h(\tau_h)$$

$$\text{s.t.} \quad C^\pi_h(\tau_h) \leq b$$

# Budget Augmentation

**Intuition**: Build $\overline{M}$ satisfying,

$$\overline{V}_h^*(s, b) = \max_{\pi \in \Pi^D} \quad V_h^\pi(\tau_h)$$
$$\text{s.t.} \quad C_h^\pi(\tau_h) \leq b$$

$\Rightarrow$ Primal $= \overline{V}_1^*(s_0, B)$

# Budget Augmentation

**Intuition**: Build $\overline{M}$ satisfying,

$$\overline{V}_h^*(s,b) = \max_{\pi \in \Pi^D} \quad V_h^\pi(\tau_h)$$
$$\text{s.t.} \quad C_h^\pi(\tau_h) \leq b$$

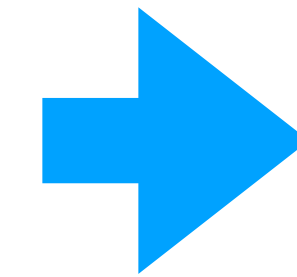$\Longrightarrow$ Primal $= \overline{V}_1^*(s_0, B)$

# Budget Augmentation

**Intuition**: Build $\overline{M}$ satisfying,

$$\overline{V}_h^*(s, b) = \max_{\pi \in \Pi^D} \quad V_h^\pi(\tau_h)$$

$$\text{s.t.} \quad C_h^\pi(\tau_h) \le b$$

$\Rightarrow \quad \text{Primal} = \overline{V}_1^*(s_0, B)$

Current budget

*Not Unique!*

$\pi$

$b$

$s, h$

$a$

$s_1', h+1 \quad b_{s_1'}$

$s_n', h+1 \quad b_{s_n'}$

# Budget Augmentation

**Intuition**: Build $\overline{M}$ satisfying,

$$\overline{V}_h^*(s, b) = \max_{\pi \in \Pi^D} \quad V_h^\pi(\tau_h)$$
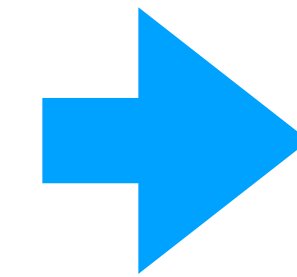$$\text{s.t.} \quad C_h^\pi(\tau_h) \leq b$$
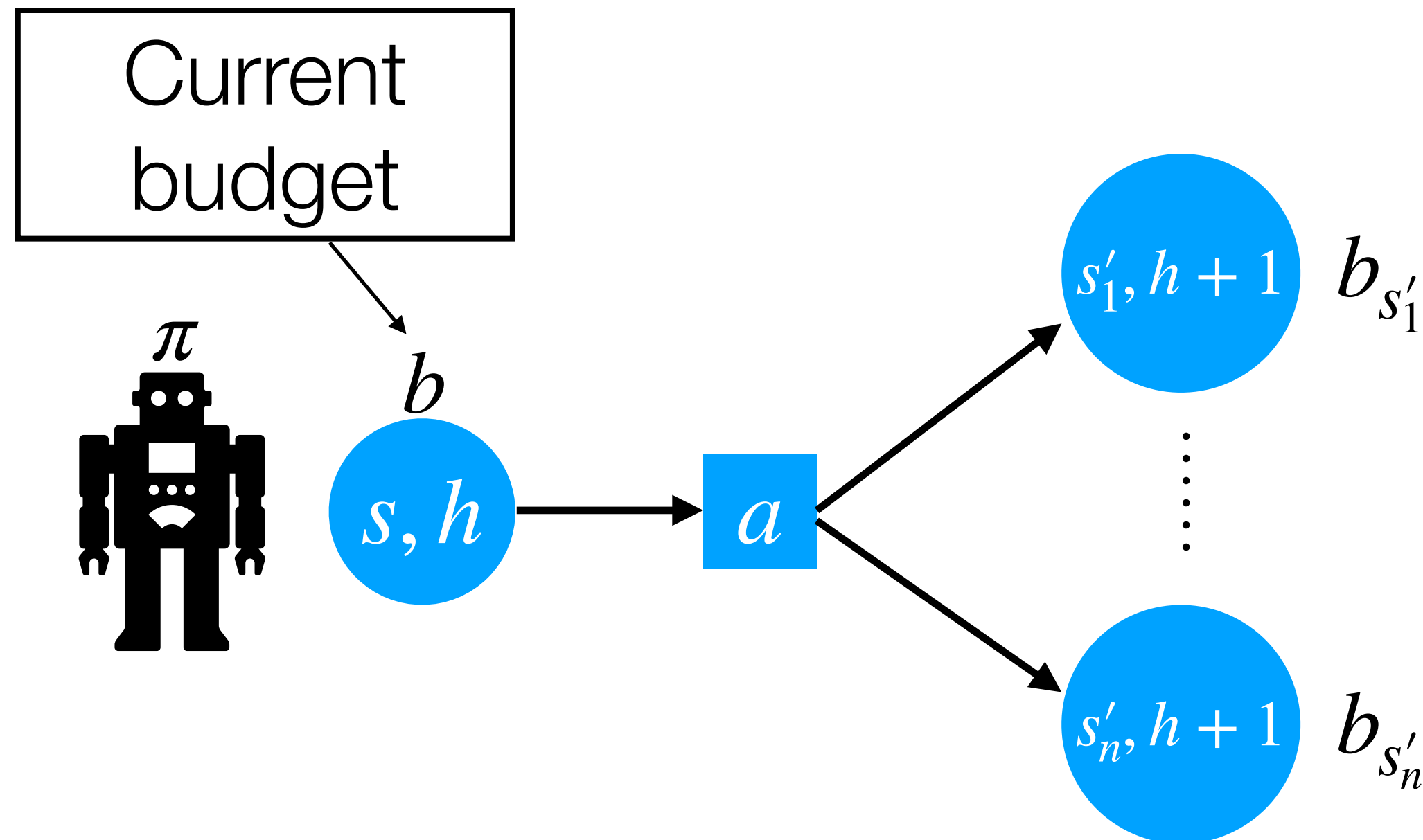
$\blacktriangleright$ Primal $= \overline{V}_1^*(s_0, B)$

# Budget Augmentation

**Intuition**: Build $\overline{M}$ satisfying,

$$\overline{V}_h^*(s,b) = \max_{\pi \in \Pi^D} \quad V_h^\pi(\tau_h)$$
$$\text{s.t.} \quad C_h^\pi(\tau_h) \leq b$$
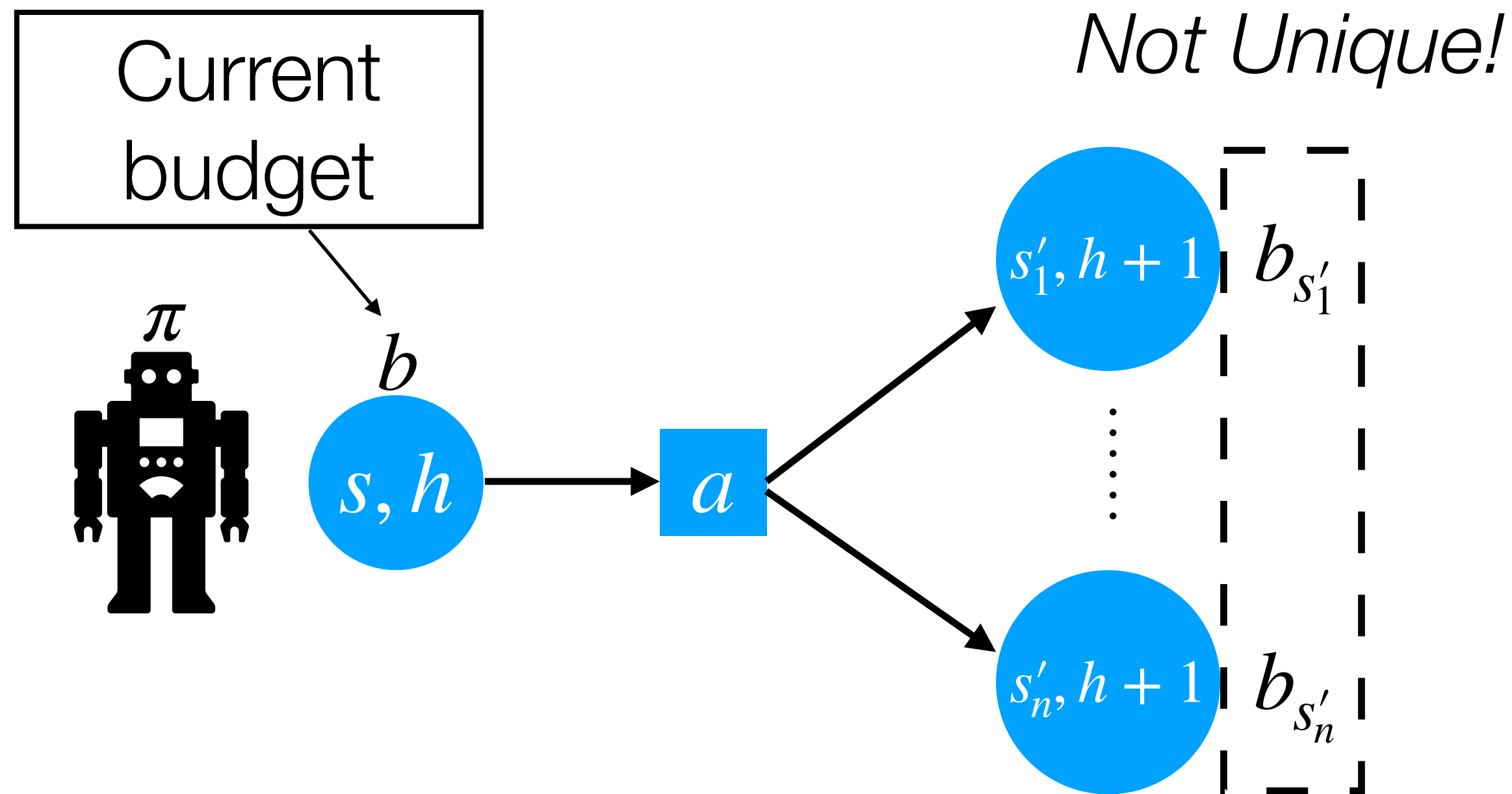
➡️ Primal $= \overline{V}_1^*(s_0, B)$

---

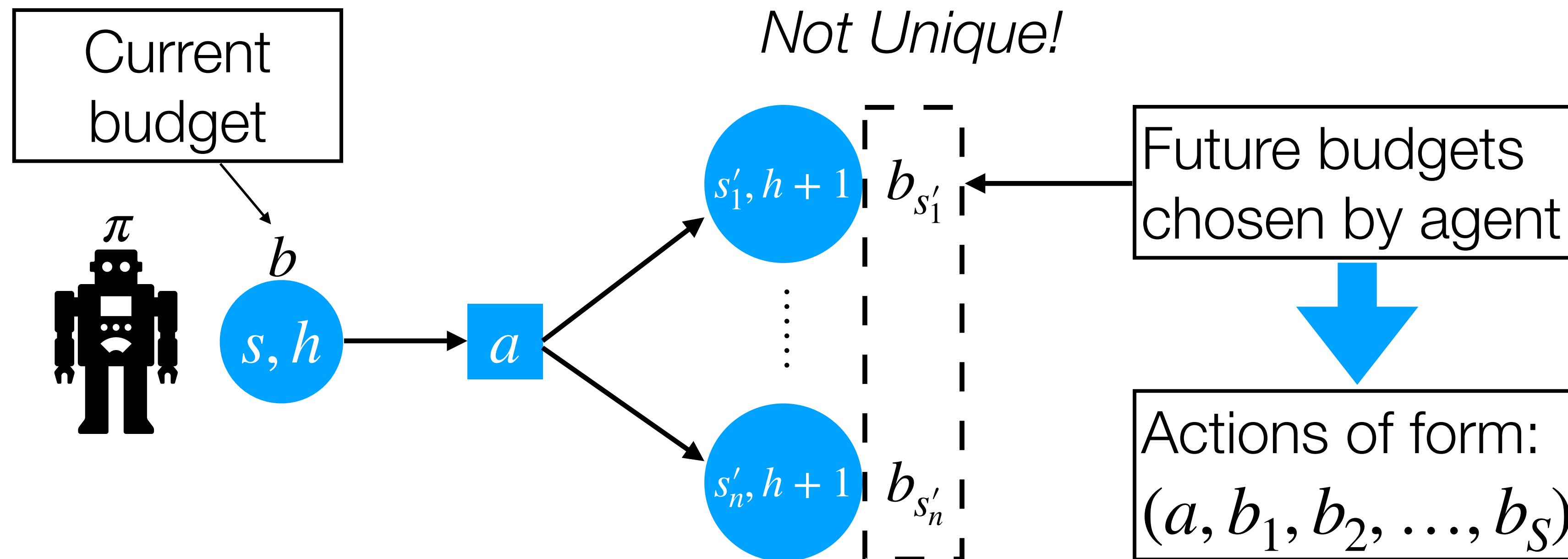Current budget

*Not Unique!*



$\pi$

$b$

$s, h$ → $a$ → $s_1', h+1$ $b_{s_1'}$

$s_n', h+1$ $b_{s_n'}$

Future budgets chosen by agent

⬇️

Actions of form:
$(a, b_1, b_2, \ldots, b_S)$

**Definition 1** (TSR). We call a cost criterion $C$ *time-recursive* (TR) if for any cMDP $M$ and policy $\pi \in \Pi^D$, $\pi$'s cost decomposes recursively into $C_M^\pi = C_1^\pi(s_0)$. Here, $C_{H+1}^\pi(\cdot) = \mathbf{0}$ and for any $h \in [H]$ and $\tau_h \in \mathcal{H}_h$,

$$C_h^\pi(\tau_h) = c_h(s, a) + f\left(\left(P_h(s' \mid s, a), C_{h+1}^\pi(\tau_h, a, s')\right)_{s' \in P_h(s,a)}\right), \tag{TR}$$

where $s = s_h(\tau_h)$, $a = \pi_h(\tau_h)$, and $f$ is a non-decreasing function[1] computable in $O(S)$ time. For technical reasons, we also require that $f(x) = \infty$ whenever $\infty \in x$.

We further say $C$ is *time-space-recursive* (TSR) if the $f$ term above is equal to $g_h^{\tau_h, a}(1)$. Here, $g_h^{\tau_h, a}(S + 1) = 0$ and for any $t \leq S$,

$$g_h^{\tau_h, a}(t) = \alpha\left(\beta\left(P_h(t \mid s, a), C_{h+1}^\pi(\tau_h, a, t)\right), g_h^{\tau_h, a}(t + 1)\right), \tag{SR}$$

where $\alpha$ is a non-decreasing function, and both $\alpha, \beta$ are computable in $O(1)$ time. We also assume that $\alpha(\cdot, \infty) = \infty$, and $\beta$ satisfies $\alpha(\beta(0, \cdot), x) = x$ to match $f$'s condition.

# Generalization

# Generalization

**Recursive cost optimization** suffices for our algorithm

# Generalization

**Recursive cost optimization** suffices for our algorithm

**Assumption [time-space recursive]:** *the optimal cost is computable recursively over both **time** and state **space***

# Generalization

**Recursive cost optimization** suffices for our algorithm

**Assumption [time-space recursive]:** *the optimal cost is computable recursively over both **time** and state **space***

*holds for expectation, almost sure, and anytime constraints*

# Action Space

# Action Space

**Policy Evaluation Equation**: $V_h^\pi(s) = r_h(s, a) + \sum_{s'} P_h(s' \mid s, a) V_{h+1}^\pi(s')$

# Action Space

**Policy Evaluation Equation:** $V_h^\pi(s) = r_h(s, a) + \sum_{s'} P_h(s' \mid s, a) V_{h+1}^\pi(s')$

Guarantee demand by:

# Action Space

**Policy Evaluation Equation:** $V_h^\pi(s) = r_h(s, a) + \sum_{s'} P_h(s' \mid s, a) V_{h+1}^\pi(s')$

Guarantee demand by:

1. $\forall i, \ V_{h+1}^\pi(s_i') \geq v_{s_i'}$

# Action Space

**Policy Evaluation Equation:** $V_h^\pi(s) = r_h(s,a) + \sum_{s'} P_h(s' \mid s,a)V_{h+1}^\pi(s')$

Guarantee demand by:

1. $\forall i,\ V_{h+1}^\pi(s_i') \geq v_{s_i'}$

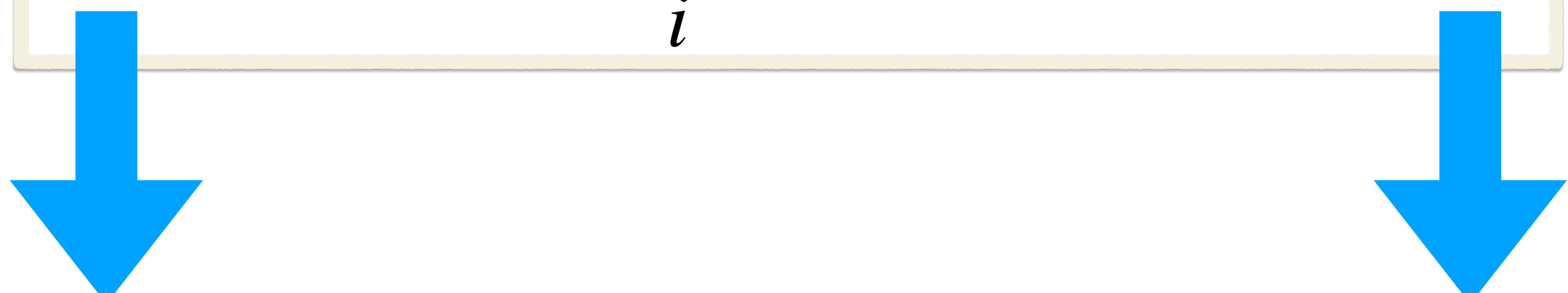2. $r_h(s,a) + \sum_i P_h(s_i' \mid s,a)v_{s_i'} \geq v$

# Action Space

**Policy Evaluation Equation:** $V_h^\pi(s) = r_h(s, a) + \sum_{s'} P_h(s' \mid s, a) V_{h+1}^\pi(s')$

Guarantee demand by:

1. $\forall i, \ V_{h+1}^\pi(s_i') \geq v_{s_i'}$

2. $r_h(s, a) + \sum_i P_h(s_i' \mid s, a) v_{s_i'} \geq v$

$$\overline{\mathcal{A}}_h(s, v) := \left\{ (a, \mathbf{v}) \in \mathcal{A} \times \mathcal{V}^S \mid r_h(s, a) + \sum_{s'} P_h(s' \mid s, a) v_{s'} \geq v \right\}$$

# Action-Space Dynamic Programming

# Action-Space Dynamic Programming

$$\min_{\mathbf{v} \in \mathcal{V}^S} \quad P_h(1 \mid s, a)\overline{C}^*_{h+1}(1, v_1) + \cdots + P_h(S \mid s, a)\overline{C}^*_{h+1}(S, v_S)$$

$$\text{s.t.} \quad P_h(1 \mid s, a)v_1 + \cdots + P_h(S \mid s, a)v_S \geq v - r_h(s, a)$$

# Action-Space Dynamic Programming

$$\min_{\mathbf{v} \in \mathcal{V}^S} \quad P_h(1 \mid s, a)\overline{C}^*_{h+1}(1, \boxed{v_1}) + \cdots + P_h(S \mid s, a)\overline{C}^*_{h+1}(S, v_S)$$

$$\text{s.t.} \quad P_h(1 \mid s, a)\boxed{v_1} + \cdots + P_h(S \mid s, a)v_S \geq v - r_h(s, a)$$

# Action-Space Dynamic Programming

$$\min_{\mathbf{v} \in \mathcal{V}^S} \quad P_h(1 \mid s, a)\overline{C}^*_{h+1}(1, \boxed{v_1}) + \cdots + P_h(S \mid s, a)\overline{C}^*_{h+1}(S, v_S)$$

$$\text{s.t.} \quad P_h(1 \mid s, a)\boxed{v_1} + \cdots + P_h(S \mid s, a)v_S \geq v - r_h(s, a)$$

*Can choose each $v_i$ independently if track the partial demand*

# Action-Space Dynamic Programming

$$\min_{\mathbf{v} \in \mathcal{V}^S} \quad P_h(1 \mid s,a)\overline{C}^*_{h+1}(1, v_1) + \cdots + P_h(S \mid s,a)\overline{C}^*_{h+1}(S, v_S)$$
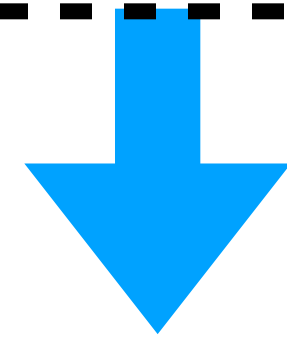
$$\text{s.t.} \quad P_h(1 \mid s,a)v_1 + \cdots + P_h(S \mid s,a)v_S \geq v - r_h(s,a)$$

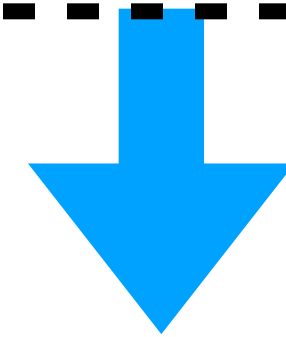*Can choose each $v_i$ independently if track the partial demand*

***Space Recursion!***
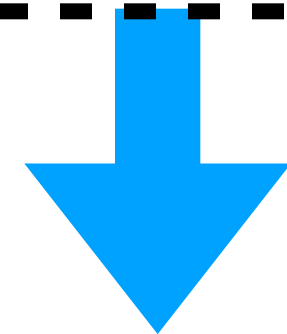
# Action-Space Dynamic Programming

$$\min_{\mathbf{v} \in \mathcal{V}^S} \quad P_h(1 \mid s, a)\overline{C}^*_{h+1}(1, \boxed{v_1}) + \cdots + P_h(S \mid s, a)\overline{C}^*_{h+1}(S, v_S)$$

$$\text{s.t.} \quad P_h(1 \mid s, a)\boxed{v_1} + \cdots + P_h(S \mid s, a)v_S \geq v - r_h(s, a)$$

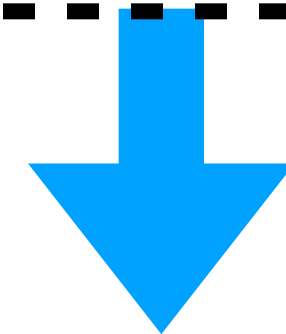*Can choose each $v_i$ independently if track the partial demand*

**Space Recursion!**
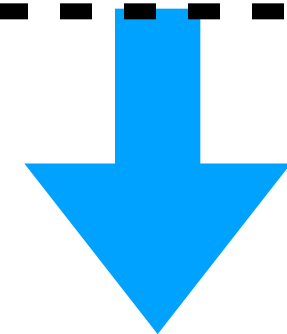
$$g(t, u) = \min_{v_t \in \mathcal{V}} P_h(t \mid s, a)C^*_{h+1}(t, v_t) + g(t + 1, u + P_h(t \mid s, a)v_t)$$
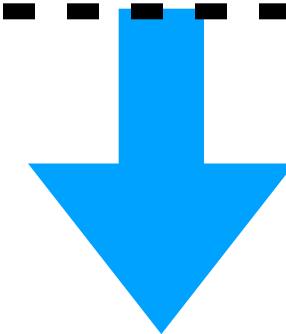
# Action-Space Dynamic Programming

$$\min_{\mathbf{v}\in\mathcal{V}^S} \quad P_h(1\mid s,a)\overline{C}^*_{h+1}(1,\boxed{v_1}) + \cdots + P_h(S\mid s,a)\overline{C}^*_{h+1}(S,v_S)$$

$$\text{s.t.} \quad P_h(1\mid s,a)\boxed{v_1} + \cdots + P_h(S\mid s,a)v_S \geq v - r_h(s,a)$$

*Can choose each $v_i$ independently if track the partial demand*

**Space Recursion!**

Partial demand

$$g(t,\boxed{u}) = \min_{v_t\in\mathcal{V}} P_h(t\mid s,a)C^*_{h+1}(t,v_t) + g(t+1, u + P_h(t\mid s,a)v_t)$$

# Action-Space Dynamic Programming

$$\min_{\mathbf{v} \in \mathcal{V}^S} \quad P_h(1 \mid s, a)\overline{C}^*_{h+1}(1, \boxed{v_1}) + \cdots + P_h(S \mid s, a)\overline{C}^*_{h+1}(S, v_S)$$

$$\text{s.t.} \quad P_h(1 \mid s, a)\boxed{v_1} + \cdots + P_h(S \mid s, a)v_S \geq v - r_h(s, a)$$

*Can choose each $v_i$ independently if track the partial demand*

## Space Recursion!

**Partial demand**

$$g(t, \boxed{u}) = \min_{v_t \in \mathcal{V}} P_h(t \mid s, a)C^*_{h+1}(t, v_t) + g(t+1, u + P_h(t \mid s, a)v_t)$$

*Value check at end:* $\quad g(S+1, u) := \chi_{\{u \geq v\}}$

# Why Deterministic Policies?

# Why Deterministic Policies?

- Cheap [1]

# Why Deterministic Policies?

- Cheap [1]

- Multi-agent coordination [2]

# Why Deterministic Policies?

- Cheap [1]

- Multi-agent coordination [2]

- Trust-worthy [3]

# Why Deterministic Policies?

- Cheap [1]

- Multi-agent coordination [2]

- Trust-worthy [3]

# Why Deterministic Policies?

- Cheap [1]

- Multi-agent coordination [2]

- Trust-worthy [3]

  - Predictable

# Why Deterministic Policies?

- Cheap [1]

- Multi-agent coordination [2]

- Trust-worthy [3]

  - Predictable

# Why Deterministic Policies?

- Cheap [1]

- Multi-agent coordination [2]

- Trust-worthy [3]

  - Predictable

- Optimal for modern constraints [4]

# Value-Demand Augmentation

# Value-Demand Augmentation

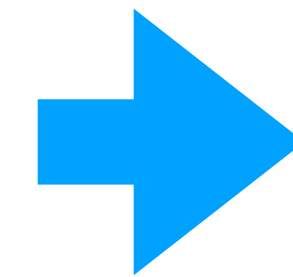**Intuition**: Build $\overline{M}$ satisfying,

# Value-Demand Augmentation

Intuition: Build $\overline{M}$ satisfying,

$$\overline{C}_h^*(s,v) = \min_{\pi \in \Pi^D} \quad C_h^\pi(\tau_h)$$
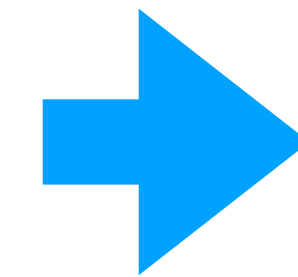
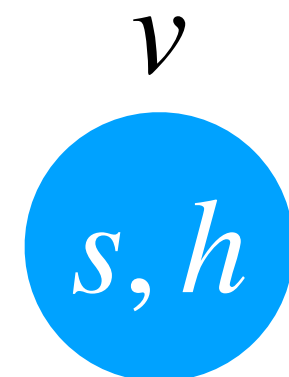$$\text{s.t.} \quad V_h^\pi(\tau_h) \geq v$$

# Value-Demand Augmentation

**Intuition**: Build $\overline{M}$ satisfying,

$$\overline{C}_h^*(s, v) = \min_{\pi \in \Pi^D} \quad C_h^\pi(\tau_h)$$
$$\text{s.t.} \quad V_h^\pi(\tau_h) \geq v$$

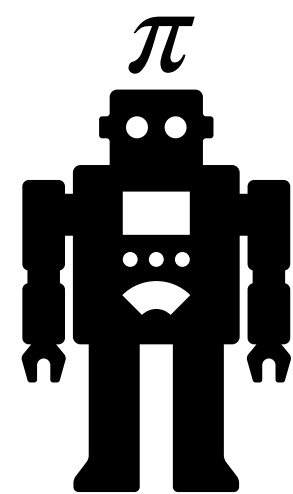$\Rightarrow$ Dual $= \overline{C}_1^*(s_0, V^*)$

# Value-Demand Augmentation

**Intuition**: Build $\overline{M}$ satisfying,

$$\overline{C}_h^*(s,v) = \min_{\pi \in \Pi^D} \quad C_h^\pi(\tau_h)$$
$$\text{s.t.} \quad V_h^\pi(\tau_h) \geq v$$

$\Rightarrow$ Dual $= \overline{C}_1^*(s_0, V^*)$

Future value demand

$\pi$

$v$

$s, h$

# Value-Demand Augmentation

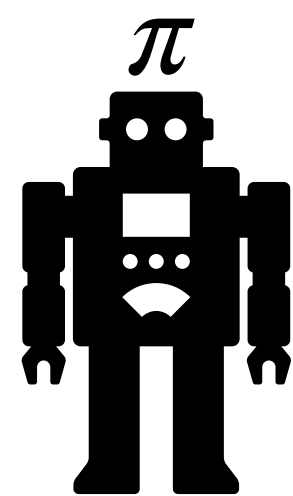**Intuition**: Build $\overline{M}$ satisfying,

$$\overline{C}_h^*(s, v) = \min_{\pi \in \Pi^D} \quad C_h^\pi(\tau_h)$$

$$\text{s.t.} \quad V_h^\pi(\tau_h) \geq v$$

$\Rightarrow$ Dual $= \overline{C}_1^*(s_0, V^*)$
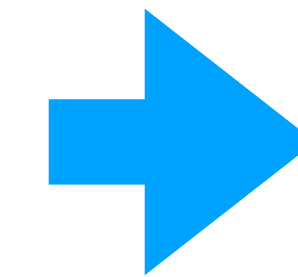
Future value demand

$\pi$

$v$

$s, h \longrightarrow a$

# Value-Demand Augmentation

Intuition: Build $\overline{M}$ satisfying,

$$\overline{C}_h^*(s,v) = \min_{\pi \in \Pi^D} \quad C_h^\pi(\tau_h)$$

$$\text{s.t.} \quad V_h^\pi(\tau_h) \geq v$$

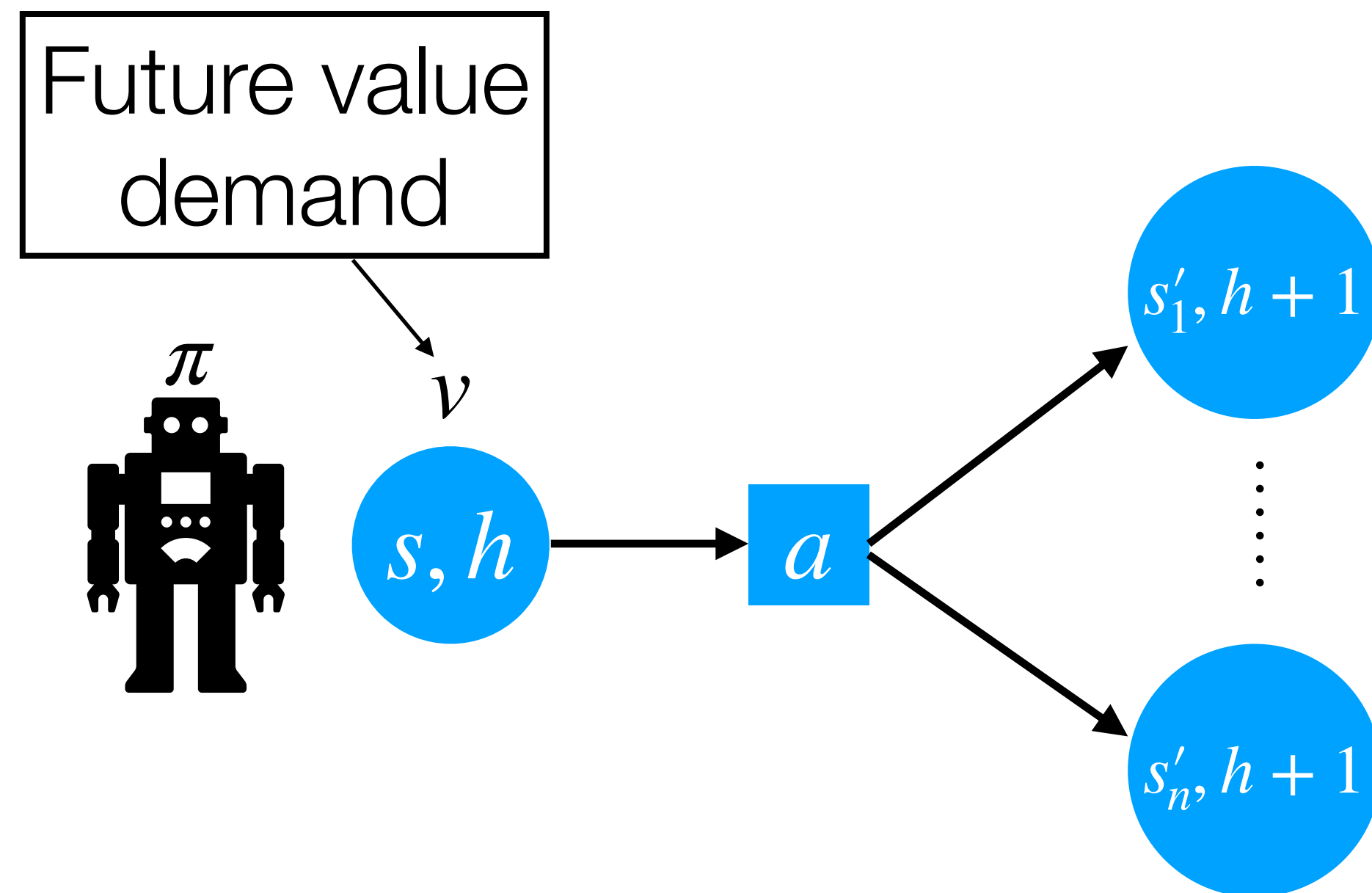$\Longrightarrow$ Dual $= \overline{C}_1^*(s_0, V^*)$

Future value demand

# Value-Demand Augmentation

**Intuition**: Build $\overline{M}$ satisfying,

$$\overline{C}_h^*(s,v) = \min_{\pi \in \Pi^D} \quad C_h^\pi(\tau_h)$$
$$\text{s.t.} \quad V_h^\pi(\tau_h) \geq v$$

$\Rightarrow$ Dual $= \overline{C}_1^*(s_0, V^*)$



Future value demand

$\pi$

$v$

$s, h$
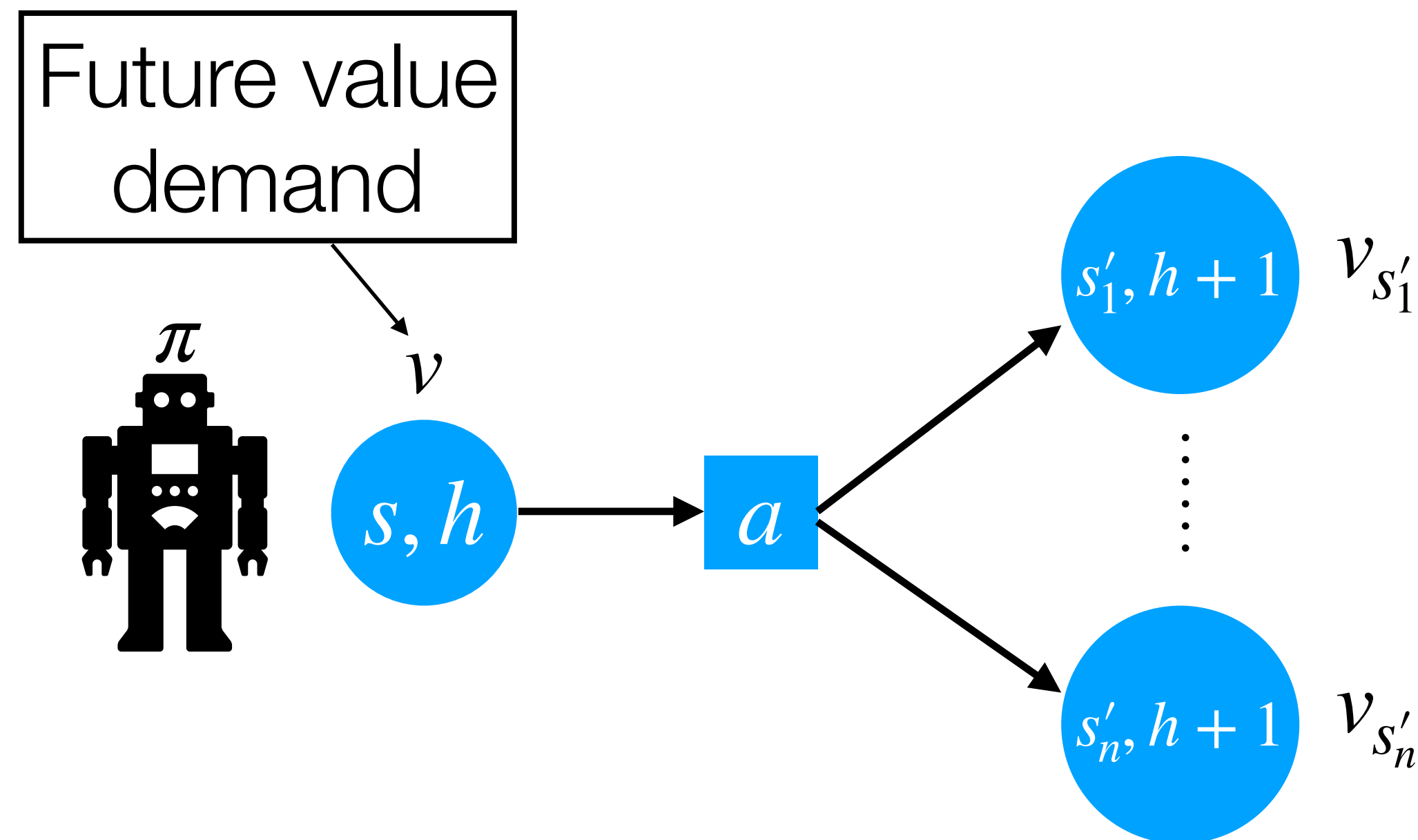
$a$

$s_1', h+1 \quad v_{s_1'}$

$s_n', h+1 \quad v_{s_n'}$

# Value-Demand Augmentation

**Intuition**: Build $\overline{M}$ satisfying,

$$\overline{C}_h^*(s,v) = \min_{\pi \in \Pi^D} \quad C_h^\pi(\tau_h)$$
$$\text{s.t.} \quad V_h^\pi(\tau_h) \geq v$$

$\Rightarrow$  Dual $= \overline{C}_1^*(s_0, V^*)$

Future value demand

*Not Unique!*

$\pi$

$v$

$s, h$
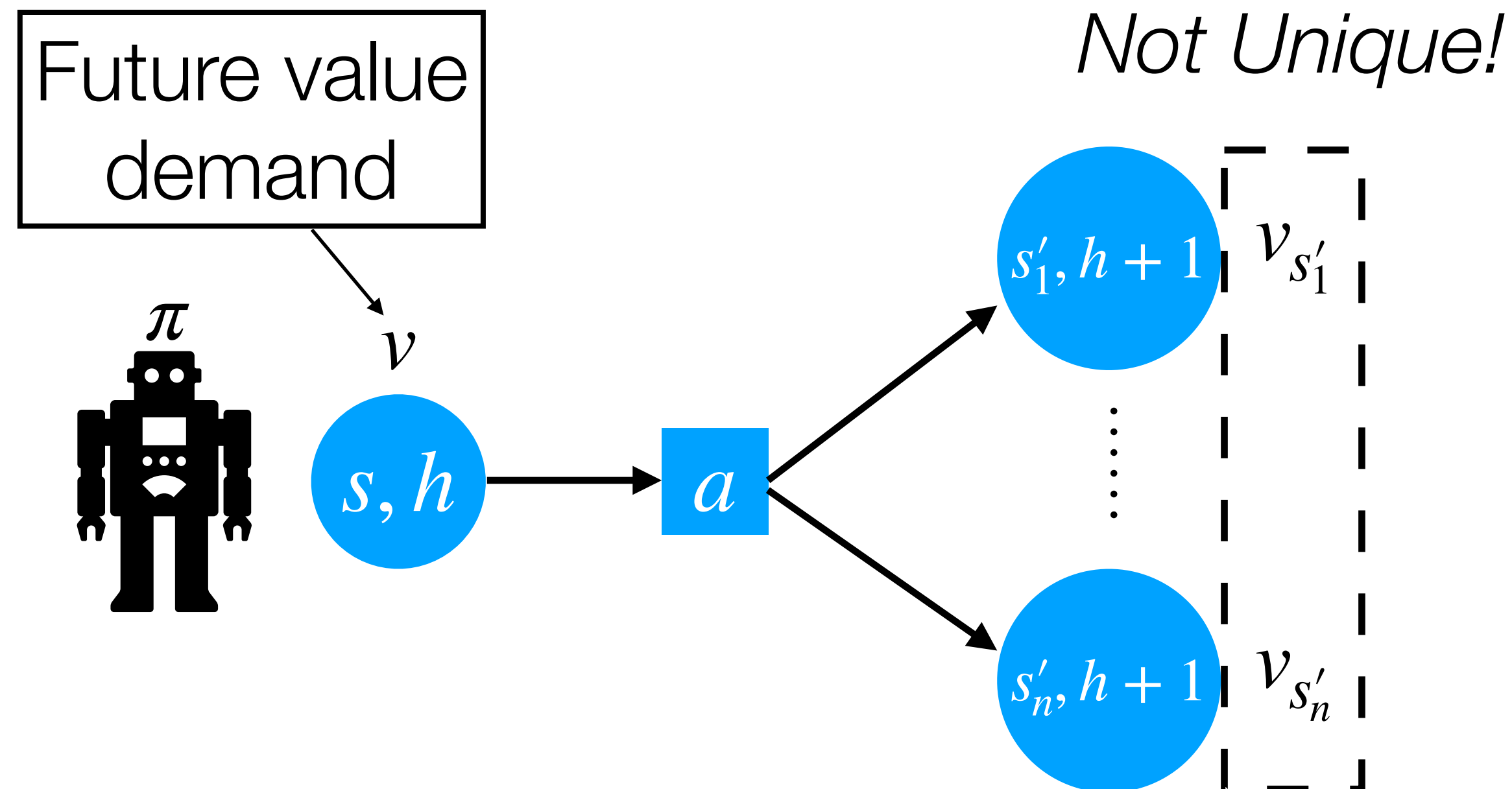
$a$

$s_1', h+1 \quad v_{s_1'}$

$s_n', h+1 \quad v_{s_n'}$

# Value-Demand Augmentation

**Intuition**: Build $\overline{M}$ satisfying,

$$\overline{C}_h^*(s,v) = \min_{\pi \in \Pi^D} \quad C_h^\pi(\tau_h)$$
$$\text{s.t.} \quad V_h^\pi(\tau_h) \geq v$$

$\blacktriangleright$ Dual $= \overline{C}_1^*(s_0, V^*)$

# Value-Demand Augmentation

**Intuition**: Build $\overline{M}$ satisfying,

$$\overline{C}_h^*(s, v) = \min_{\pi \in \Pi^D} \quad C_h^\pi(\tau_h)$$

$$\text{s.t.} \quad V_h^\pi(\tau_h) \geq v$$

$\Rightarrow$ Dual $= \overline{C}_1^*(s_0, V*)$
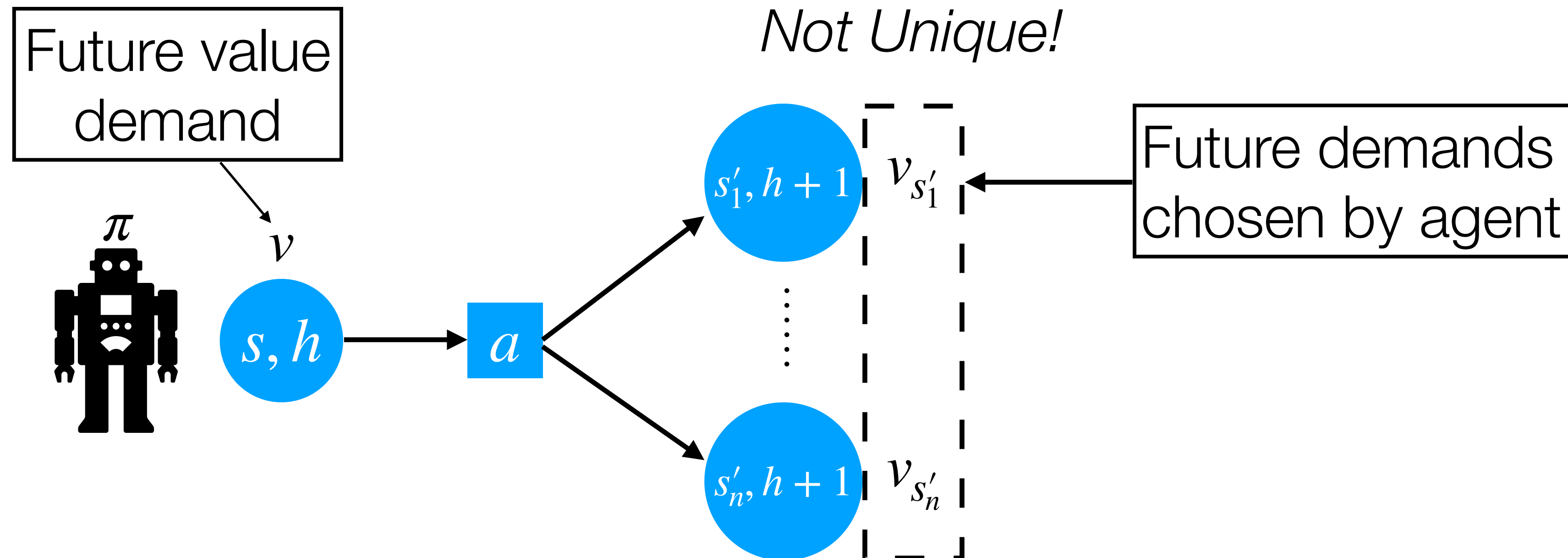


Future value demand

$\pi$

$v$

$s, h$

$a$

*Not Unique!*

$s_1', h+1$   $v_{s_1'}$

$s_n', h+1$   $v_{s_n'}$

Future demands chosen by agent

Actions of form:
$(a, v_1, v_2, \ldots, v_S)$

# Constraint Landscape

Put the formulas in here

# Constraint Landscape

Put the formulas in here

**Constraints**

- Expectation

- Chance

- Almost Sure

- Anytime

# Constraint Landscape

Put the formulas in here

**Constraints**

Flexibility ——

- Expectation

- Chance

- Almost Sure

- Anytime

# Constraint Landscape

Put the formulas in here

## Constraints

Flexibility ———

Precision ———

- Expectation

- Chance

- Almost Sure

- Anytime

# Constraint Landscape

Put the formulas in here

**Constraints**

Flexibility ——— • Expectation

Precision ——— • Chance

Safety { • Almost Sure

• Anytime

# Constraint Landscape

Put the formulas in here

**Constraints**

**Policies**

Flexibility ——— • Expectation

Precision ——— • Chance

Safety { • Almost Sure

• Anytime

# Constraint Landscape

Put the formulas in here

**Constraints**

**Policies**

Flexibility —— • Expectation

Precision —— • Chance

Safety { • Almost Sure

• Anytime

• Stochastic

# Constraint Landscape

Put the formulas in here

**Constraints**

**Policies**

Flexibility —— • Expectation

Precision —— • Chance

• Almost Sure

Safety {

• Anytime

• Stochastic

• Deterministic

# Constraint Landscape

Put the formulas in here

**Constraints**

- Expectation
- Chance
- Almost Sure
- Anytime

Flexibility
Precision
Safety {

**Policies**

- Stochastic
- Deterministic

} Predictable

# Constraint Landscape

Put the formulas in here

**Constraints**

Flexibility —— • Expectation

Precision —— • Chance

Safety { • Almost Sure

• Anytime

**Policies**

• Stochastic

• Deterministic

} Predictable

Cheap

# Constraint Landscape

Put the formulas in here

**Constraints**

Flexibility —— • Expectation

Precision —— • Chance

Safety { 
• Almost Sure
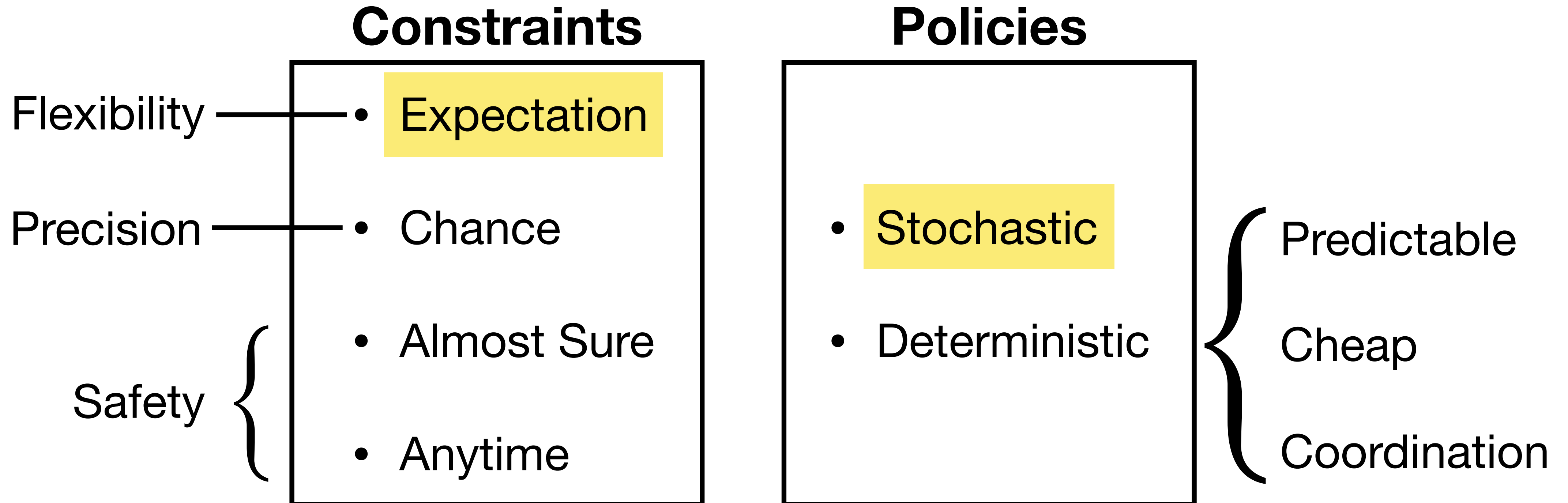
• Anytime

**Policies**

• Stochastic

• Deterministic

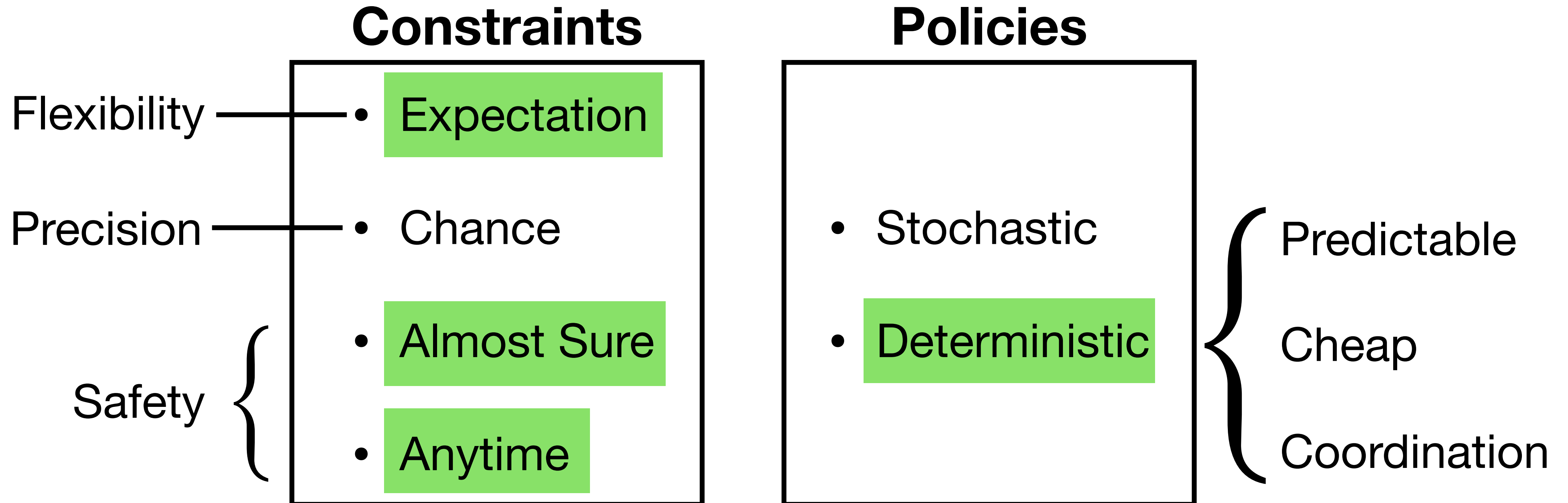} Predictable

Cheap

Coordination

# Constraint Landscape

Put the formulas in here

**Constraints**

Flexibility —— • <mark>Expectation</mark>

Precision —— • Chance

Safety { • Almost Sure

• Anytime

**Policies**

• <mark>Stochastic</mark>

• Deterministic } Predictable

Cheap

Coordination

# Constraint Landscape

# General Formulation

# General Formulation

# General Formulation



$$\pi \xrightarrow{a} s \rightarrow \begin{cases} r_h \sim R_h(s, a) \\ c_h \sim C_h(s, a) \end{cases}$$

# General Formulation



$\pi$

$a$

$s$

$$\begin{cases} r_h \sim R_h(s, a) \\ c_h \sim C_h(s, a) \end{cases}$$

**Agent's goal:**

$$\max_{\pi} V^{\pi}$$

$$\text{s.t. constraints on } \sum_{h=1}^{H} c_h$$