# Optimal Attack and Defense for Reinforcement Learning

*Jeremy McMahan*, Young Wu, Xiaojin Zhu, Qiaomin Xie

# RL Basics
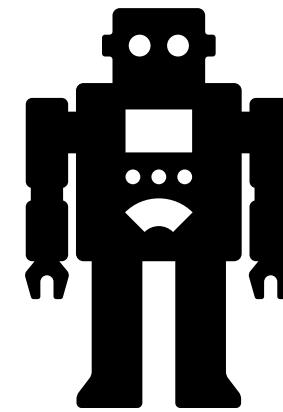
# RL Interaction Protocol

# RL Interaction Protocol

Models sequential decision making in uncertain environments

# RL Interaction Protocol

Models sequential decision making in uncertain environments
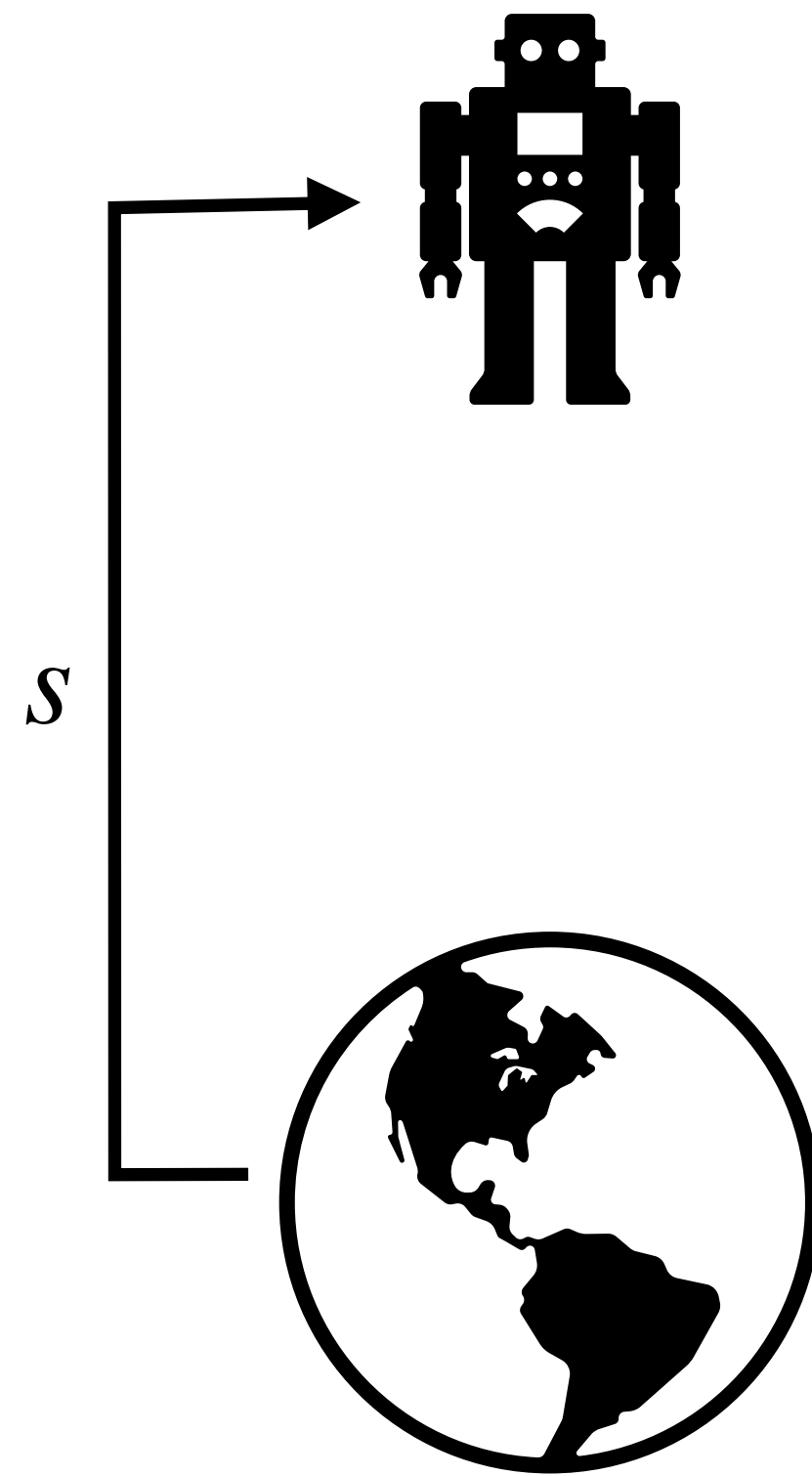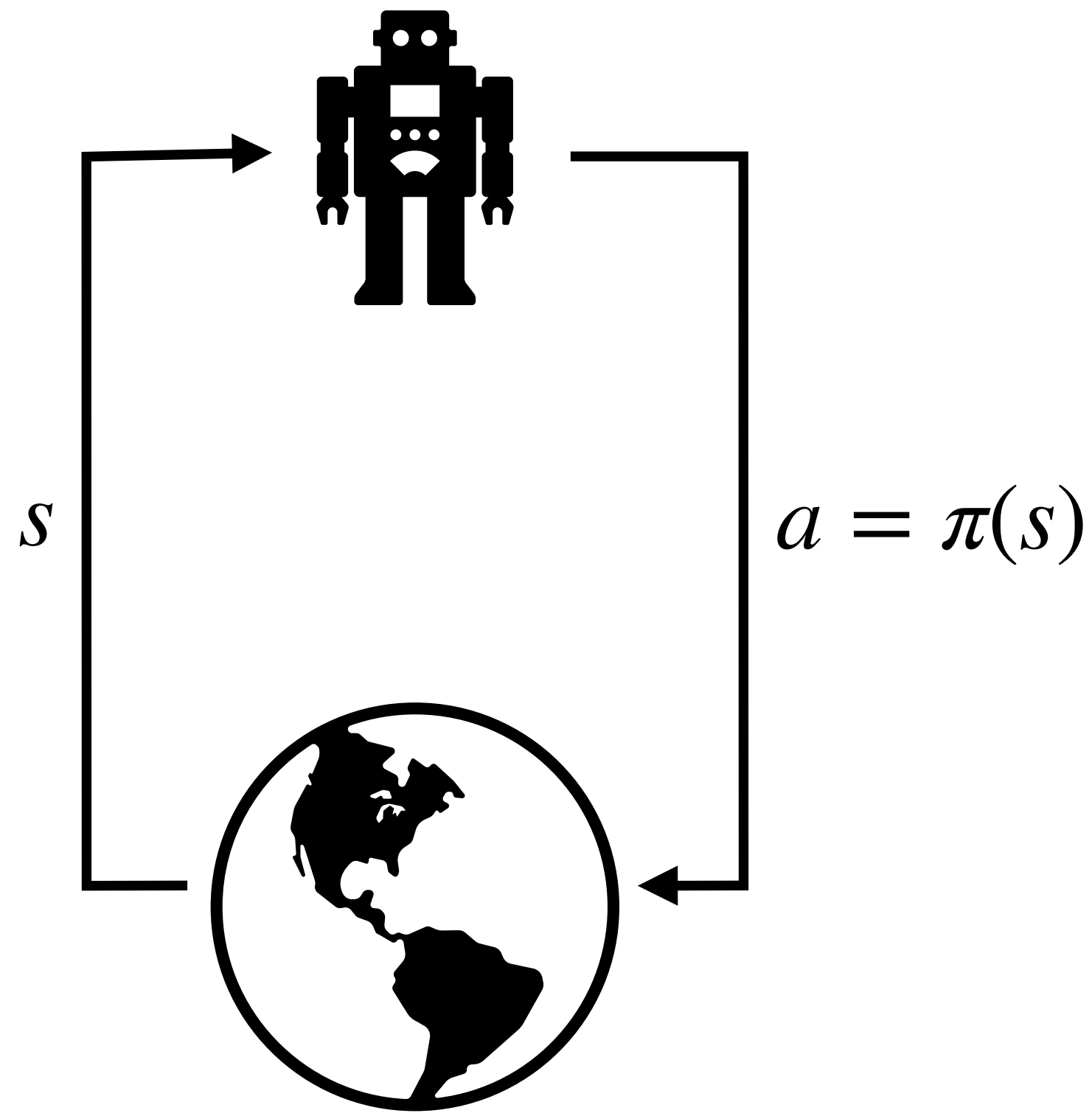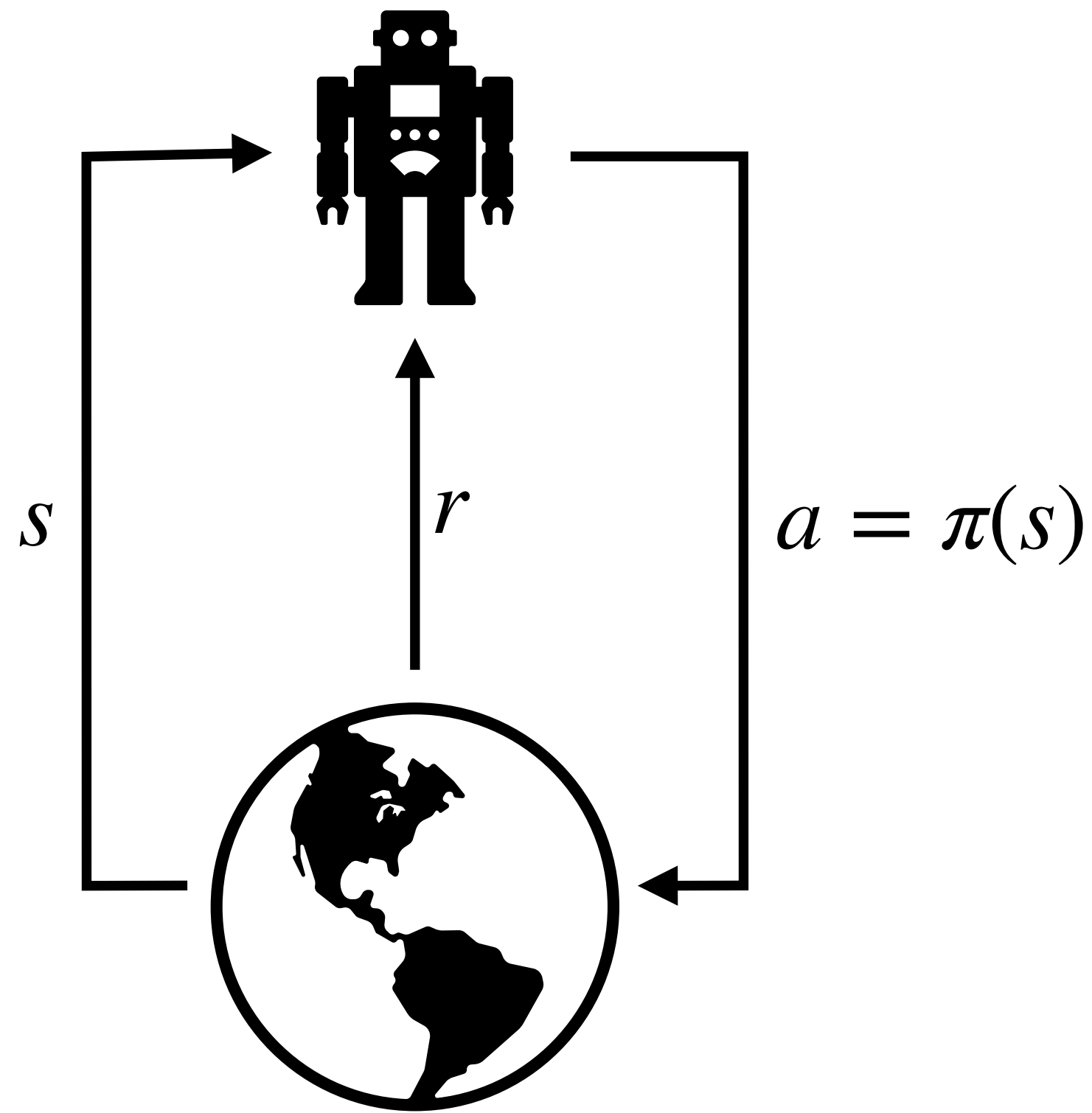
# RL Interaction Protocol

Models sequential decision making in uncertain environments

# RL Interaction Protocol

Models sequential decision making in uncertain environments

$s$

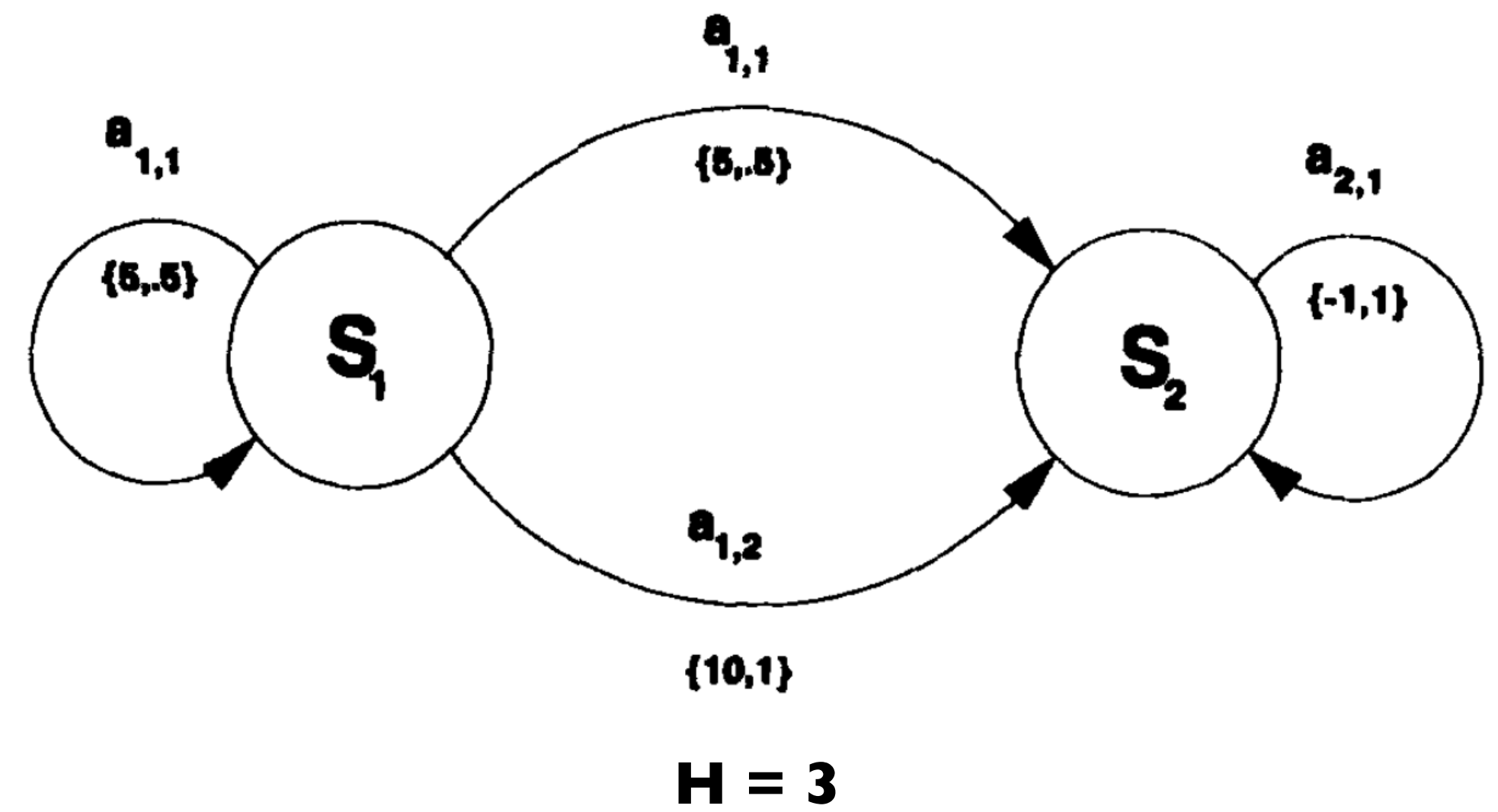$a = \pi(s)$

# RL Interaction Protocol

Models sequential decision making in uncertain environments

# Model: MDPs
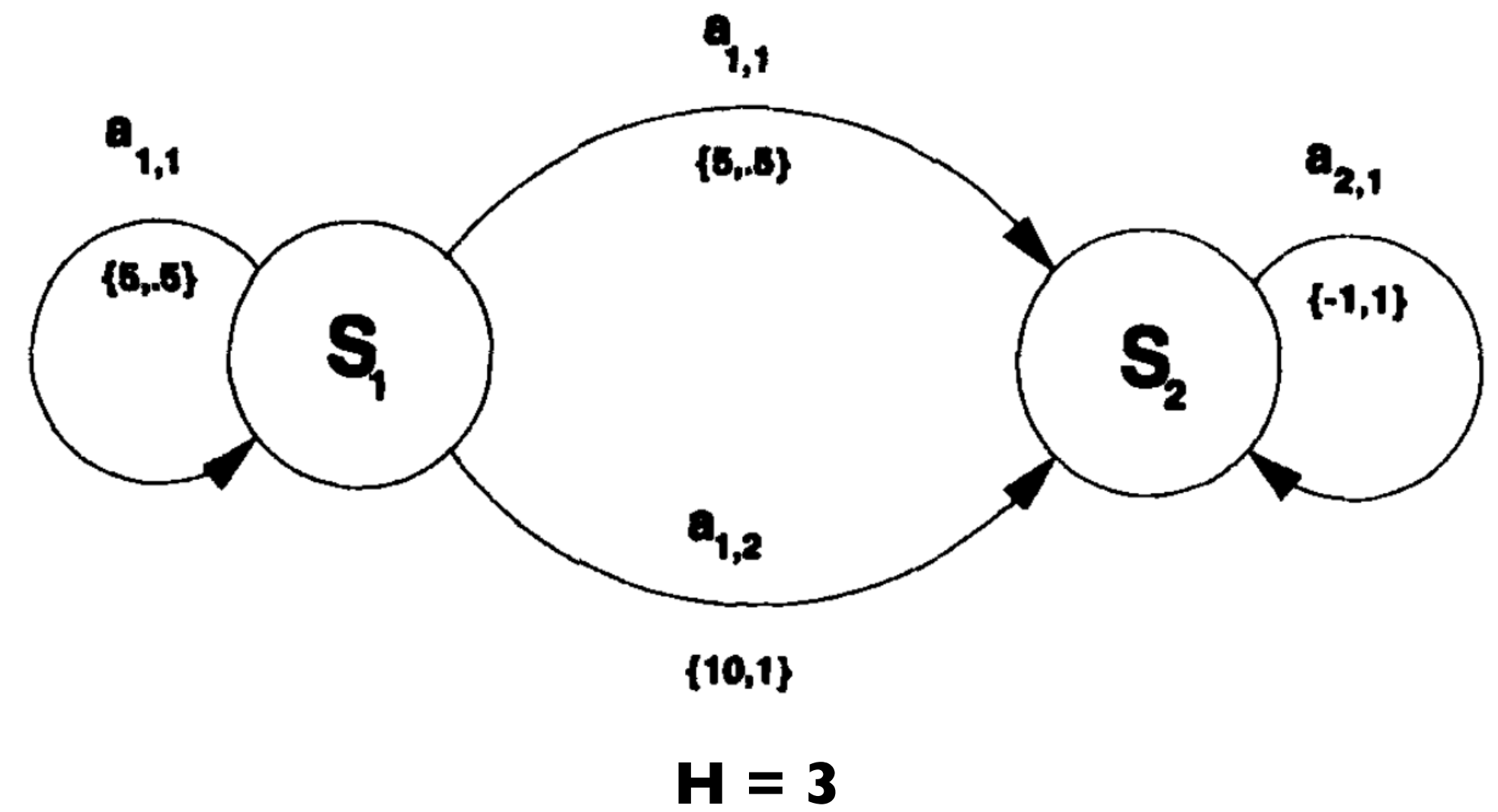
**H = 3**

# Model: MDPs
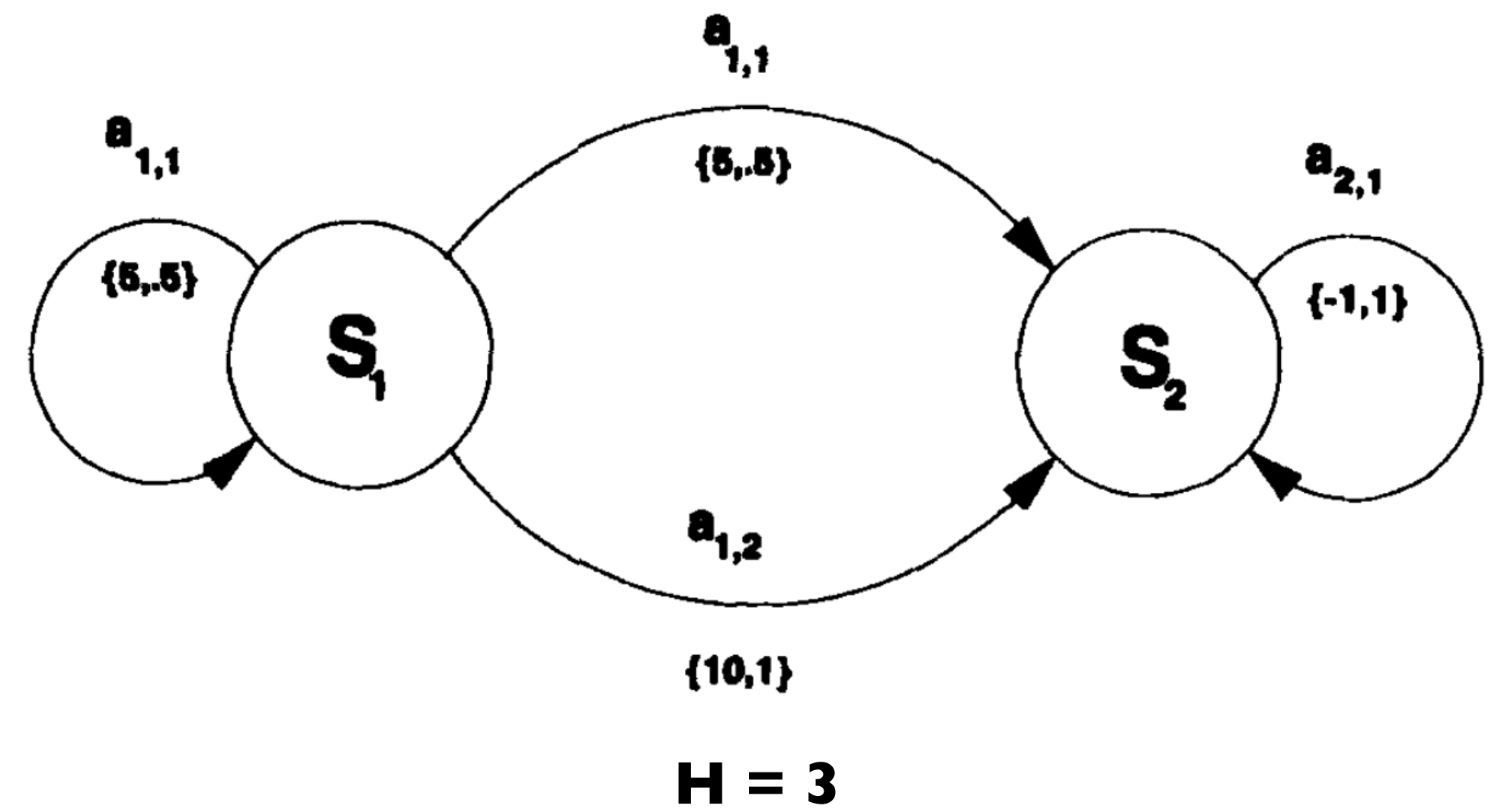
# Model: MDPs

- States, $S$

# Model: MDPs

- States, $S$

- Actions, $A$

# Model: MDPs

- States, $S$

- Actions, $A$

- Rewards, $r(s, a)$

# Model: MDPs

- States, $S$

- Actions, $A$

- Rewards, $r(s, a)$

- Transition Probabilities, $P(s' \mid s, a)$

# Model: MDPs

- States, $S$

- Actions, $A$

- Rewards, $r(s, a)$

- Transition Probabilities, $P(s' \mid s, a)$

- Time Horizon, $H$

# Policies

# Policies

A *policy* is a plan of what action to take in each state.

# Policies

A *policy* is a plan of what action to take in each state.

# Policies

A *policy* is a plan of what action to take in each state.

$$\pi(s_1) = a_{1,2} \quad \pi(s_2) = a_{2,1}$$

# Value

The *value* of M under $\pi$ is: $V^\pi(s) = E_\pi \left[ \sum_{h=1}^{H} r_h(s, a) \mid s_0 = s \right]$.

# Value



$$\pi(s_1) = a_{1,2} \qquad \text{Reward} = 10$$

# Value



$\pi(s_2) = a_{2,1}$          Reward = -1

# Value



$$\pi(s_2) = a_{2,1} \qquad \text{Reward} = -1$$

# Value

$$V^\pi(s_1) = 10 - 1 - 1 = 8$$

# Optimal Policies

# Optimal Policies

$$\pi^* = \sup_{\pi} V^{\pi}(s_0)$$

# Example MDP

# Example MDP

# Example MDP

Disaster Relief with Autonomous Vehicles

- State Space is $\mathbb{R}^2$

- Action Space is $[-1,1]^2$

- New location is $s + a$

- Reward for finding people in need.

# Performance of Optimal Policies

# Performance of Optimal Policies

# Performance of Optimal Policies

Unique optimal policy $\pi^*$ is:

$$\pi^*$$

| t/S | s_0 | s_1 |
|-----|-----|-----|
| h   | a_0 | a_1 |
| H   | a_1 | a_0 |

$$M$$

# Performance of Optimal Policies

$$M$$

Unique optimal policy $\pi^*$ is:

$$\pi^*$$

| t/S | s_0 | s_1 |
|-----|-----|-----|
| h   | a_0 | a_1 |
| H   | a_1 | a_0 |

The optimal policy achieves value:

$$V_M^{\pi^*} = 2(H-1)$$

# Optimal Policies are NOT Robust

# Optimal Policies are NOT Robust

Optimal Policies may behave poorly under measurement noise or adversarial manipulations!

# Optimal Policies are NOT Robust

Optimal Policies may behave poorly under measurement noise or adversarial manipulations!

# Optimal Policies are NOT Robust

Optimal Policies may behave poorly under measurement noise or adversarial manipulations!

- If first state is actually $s_1$ or M receives $a_1$ instead, $\pi^*$ at best gets $1/2$ of its value.

# Optimal Policies are NOT Robust

Optimal Policies may behave poorly under measurement noise or adversarial manipulations!

- If first state is actually $s_1$ or M receives $a_1$ instead, $\pi^*$ at best gets $1/2$ of its value.

- If states are swapped consistently, $\pi^*$ gets no value!

# Optimal Policies are NOT Robust

Optimal Policies may behave poorly under measurement noise or adversarial manipulations!

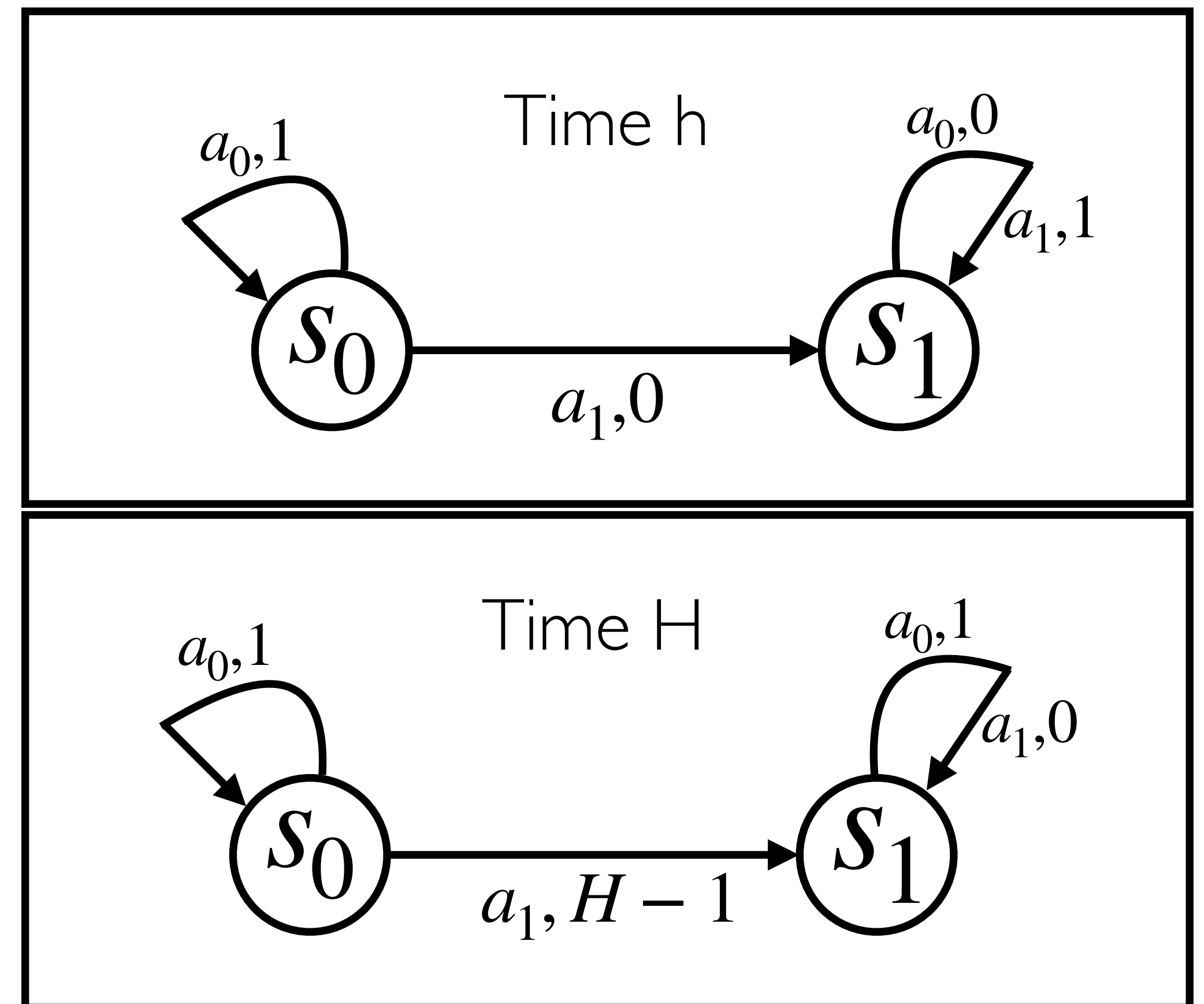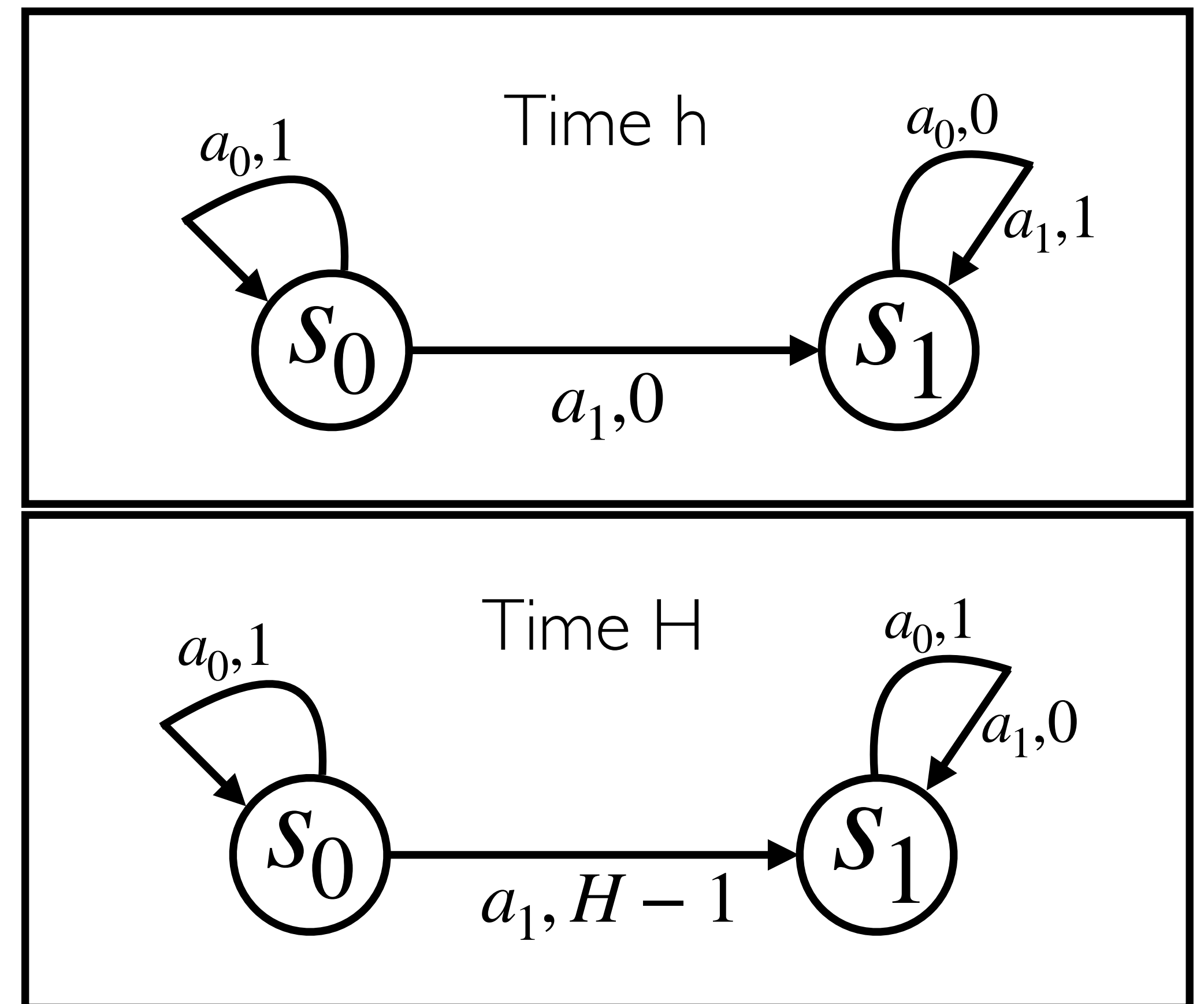- If first state is actually $s_1$ or M receives $a_1$ instead, $\pi^*$ at best gets $1/2$ of its value.

- If states are swapped consistently, $\pi^*$ gets no value!

$$\pi^* \circ \nu$$

| t/S | s_0 | s_1 |
|-----|-----|-----|
| h | a_1 | a_0 |
| H | a_0 | a_1 |

# Security Threats to RL

# Security Threats to RL

- Playing an <u>optimal</u> policy for the <span style="color:red">ideal</span> environment is not always <u>optimal</u> for the <span style="color:red">real</span> environment!

# Security Threats to RL

- Playing an <u>optimal</u> policy for the <span style="color:red">ideal</span> environment is not always <u>optimal</u> for the <span style="color:red">real</span> environment!

- Strategies to compute <span style="color:green">robust</span> policies are needed.

# Security Threats to RL

- Playing an <u>optimal</u> policy for the <span style="color:red">ideal</span> environment is not always <u>optimal</u> for the <span style="color:red">real</span> environment!

- Strategies to compute <span style="color:green">robust</span> policies are needed.

- Inspiration for field of <span style="color:magenta">adversarial RL</span>.

# Adversarial RL

# Adversarial RL

# Adversarial RL

An external attacker can manipulate the interaction.

# Adversarial RL

An external attacker can manipulate the interaction.

# Attack Paradigms

# Attack Paradigms

**Training Time**

# Attack Paradigms

**Training Time**

# Attack Paradigms

**Training Time**

Learn bad $\pi^\dagger$

# Attack Paradigms

**Training Time**



Learn bad $\pi^\dagger$

# Attack Paradigms



**Training Time**

Learn bad $\pi^\dagger$

**Test Time**

$\pi$

# Attack Paradigms

**Training Time**

Learn bad $\pi^{\dagger}$

Test Time

Cause bad outcomes

# Attack Paradigms



**Training Time**

Learn bad $\pi^\dagger$

**Test Time**

$\pi$

Cause bad outcomes

**Trojan**

# Attack Paradigms



**Training Time**

Learn bad $\pi^\dagger$

**Test Time**

$\pi$

Cause bad outcomes

**Trojan**

Hybrid: poison training to make
policy easily test-time attackable

# Panda Example

# Panda Example

In Explaining and Harnessing Adversarial Examples, Goodfellow and his team added a small perturbation to the image of a panda, as seen below. The result was surprising. Not only did the classifier mark the panda as a gibbon, but did so with high confidence.

As you can see, a barely noticeable disturbance that appears normal to us can easily deceive an ML model into predicting an incorrect class.



"panda"

57.7% confidence

"gibbon"

99.3% confidence

*Source: Goodfellow et al, 2014*

# Car Crashing

# Car Crashing

While the panda turned gibbon in the eyes of a machine is a harmless example of an adversarial attack, there are other forms of danger we must watch out for.

For instance, adversarial examples can also be used to hijack the ML models behind autonomous vehicles, causing them to misclassify 'stop' signs as 'yield', as seen below.



*Source: Kumar et al, 2021*

# Attack Surfaces

# Attack Surfaces

# Attack Surfaces

# Attack Surfaces

# Attack Surfaces

# Attack Surfaces

# Attack Surfaces

# Attack Surfaces

# Attack Surfaces

# Attack Surfaces

# Attack Surfaces

# Attack Surfaces

# Attack Surfaces

# Attack Surfaces

# Attack Surfaces



$\pi$

$o^\dagger$

$r^\dagger$

Reward Attack

$a^\dagger$

$s^\dagger$

# Attack Surfaces

# Attack Surfaces

# Attack Surfaces

- **State Attack:** changes the state of $M$ from $s_t$ to $s_t^\dagger$.

# Attack Surfaces

- **State Attack:** changes the state of $M$ from $s_t$ to $s_t^\dagger$.

- **Observation Attack:** changes the agent's observation from $o_t$ to $o_t^\dagger$.

# Attack Surfaces

- **State Attack:** changes the state of $M$ from $s_t$ to $s_t^\dagger$.

- **Observation Attack:** changes the agent's observation from $o_t$ to $o_t^\dagger$.

- **Action Attack:** changes the action $M$ receives from $a_t$ to $a_t^\dagger$.

# Attack Surfaces

- **State Attack:** changes the state of $M$ from $s_t$ to $s_t^{\dagger}$.

- **Observation Attack:** changes the agent's observation from $o_t$ to $o_t^{\dagger}$.

- **Action Attack:** changes the action $M$ receives from $a_t$ to $a_t^{\dagger}$.

- **Reward Attack:** changes the agent's reward from $r_t$ to $r_t^{\dagger}$.

# Attack Surfaces

- **State Attack:** changes the state of $M$ from $s_t$ to $s_t^\dagger$.

- **Observation Attack:** changes the agent's observation from $o_t$ to $o_t^\dagger$.

- **Action Attack:** changes the action $M$ receives from $a_t$ to $a_t^\dagger$.

- **Reward Attack:** changes the agent's reward from $r_t$ to $r_t^\dagger$.

The attacker can manipulate any element of the interaction tuple $(s, a, r)$.

# Maze Environment

# Maze Environment

# Maze Environment

- Green Squares are obstacles.

# Maze Environment

- Green Squares are obstacles.

- Yellow Squares are traversable.

# Maze Environment

- Green Squares are obstacles.

- Yellow Squares are traversable.

- The agent starts at top left corner.

# Maze Environment

- Green Squares are obstacles.

- Yellow Squares are traversable.

- The agent starts at top left corner.

- The agent receives reward only at the bottom right corner.

# Maze Environment

- Green Squares are obstacles.

- Yellow Squares are traversable.

- The agent starts at top left corner.

- The agent receives reward only at the bottom right corner.

- An optimal (shortest path) policy for the agent is in purple.

# Perceived-State Attack

# Perceived-State Attack

- Attacker shows agent $s^\dagger$.

# Perceived-State Attack

- Attacker shows agent $s^\dagger$.

- Agent chooses action $\pi(s^\dagger)$ instead of $\pi(s)$

# Perceived-State Attack

- Attacker shows agent $s^{\dagger}$.

- Agent chooses action $\pi(s^{\dagger})$ instead of $\pi(s)$

# Perceived-State Attack

# Perceived-State Attack

# Perceived-State Attack



Check out Shubham's full paper in Neurips22!
Provable Defense against Backdoor Policies in Reinforcement Learning

# Action Attack

# Action Attack

Attacker intercepts $a = \pi(s)$ and sends $a^\dagger$ to the environment instead.

# Action Attack

Attacker intercepts $a = \pi(s)$ and sends $a^{\dagger}$ to the environment instead.

# True-State Attack

# True-State Attack

Attacker changes the environment's state to $s^\dagger$

# True-State Attack

Attacker changes the environment's state to $s^\dagger$

# Reward Attack

# Reward Attack

Attacker changes the reward the agent
receives to $r^\dagger$

# Reward Attack



Attacker changes the reward the agent receives to $r^\dagger$

# Motivation: Robust Policies

# Motivation: Robust Policies

Optimal policies may be sensitive to noise or attacks.

# Motivation: Robust Policies

Optimal policies may be sensitive to noise or attacks.

# Motivation: Robust Policies

Optimal policies may be sensitive to noise or attacks.

# Motivation: Robust Policies

Optimal policies may be sensitive to noise or attacks.

# What's known?

# What's known?

**Optimal Observation Attacks**

## Robust Deep Reinforcement Learning against Adversarial Perturbations on State Observations

Huan Zhang[*,1]  Hongge Chen[*,2]  Chaowei Xiao[3]

Bo Li[4]  Mingyan Liu[5]  Duane Boning[2]  Cho-Jui Hsieh[1]

[1]UCLA  [2]MIT  [3]NVIDIA  [4]UIUC  [5]University of Michigan

huan@huan-zhang.com, chenhg@mit.edu, chaoweix@nvidia.com,
lbo@illinois.edu,mingyan@umich.edu,boning@mtl.mit.edu,chohsieh@cs.ucla.edu

# What's known?

## Optimal Observation Attacks

**Robust Deep Reinforcement Learning against Adversarial Perturbations on State Observations**

**Huan Zhang**[*,1] **Hongge Chen**[*,2] **Chaowei Xiao**[3]
**Bo Li**[4] **Mingyan Liu**[5] **Duane Boning**[2] **Cho-Jui Hsieh**[1]
[1]UCLA [2] MIT [3]NVIDIA [4]UIUC [5]University of Michigan
huan@huan-zhang.com, chenhg@mit.edu, chaoweix@nvidia.com,
lbo@illinois.edu,mingyan@umich.edu,boning@mtl.mit.edu,chohsieh@cs.ucla.edu

## [Training-time] Action and Reward Attacks

**Understanding the Limits of Poisoning Attacks in Episodic Reinforcement Learning**

**Anshuka Rangi**[1] , **Haifeng Xu**[2] , **Long Tran-Thanh**[3] , **Massimo Franceschetti**[1]
[1] University of California San Diego, USA
[2]University of Virginia, USA
[3]University of Warwick, UK
{arangi, mfranceschetti}@ucsd.edu, hx4ad@virginia.edu, long.tran-thanh@warwick.ac.uk

# What's known?

## Optimal Observation Attacks

### Robust Deep Reinforcement Learning against Adversarial Perturbations on State Observations

**Huan Zhang**[*,1]  **Hongge Chen**[*,2]  **Chaowei Xiao**[3]
**Bo Li**[4]  **Mingyan Liu**[5]  **Duane Boning**[2]  **Cho-Jui Hsieh**[1]
[1]UCLA  [2]MIT  [3]NVIDIA  [4]UIUC  [5]University of Michigan
huan@huan-zhang.com, chenhg@mit.edu, chaoweix@nvidia.com,
lbo@illinois.edu,mingyan@umich.edu,boning@mtl.mit.edu,chohsieh@cs.ucla.edu

## Defense against a specific reward attack algorithm

### Defense Against Reward Poisoning Attacks in Reinforcement Learning

**Kiarash Banihashem**          **Adish Singla**          **Goran Radanovic**
MPI-SWS                         MPI-SWS                   MPI-SWS
kbanihas@mpi-sws.org       adishs@mpi-sws.org       gradanovic@mpi-sws.org

## [Training-time] Action and Reward Attacks

### Understanding the Limits of Poisoning Attacks in Episodic Reinforcement Learning

**Anshuka Rangi**[1] ,  **Haifeng Xu**[2] ,  **Long Tran-Thanh**[3] ,  **Massimo Franceschetti**[1]
[1] University of California San Diego, USA
[2]University of Virginia, USA
[3]University of Warwick, UK
{arangi, mfranceschetti}@ucsd.edu, hx4ad@virginia.edu, long.tran-thanh@warwick.ac.uk

# What's known?

## Optimal Observation Attacks

**Robust Deep Reinforcement Learning against Adversarial Perturbations on State Observations**

**Huan Zhang**[*,1]    **Hongge Chen**[*,2]    **Chaowei Xiao**[3]
**Bo Li**[4]    **Mingyan Liu**[5]    **Duane Boning**[2]    **Cho-Jui Hsieh**[1]
[1]UCLA    [2] MIT    [3]NVIDIA    [4]UIUC    [5]University of Michigan
huan@huan-zhang.com, chenhg@mit.edu, chaoweix@nvidia.com,
lbo@illinois.edu,mingyan@umich.edu,boning@mtl.mit.edu,chohsieh@cs.ucla.edu

## [Training-time] Action and Reward Attacks

**Understanding the Limits of Poisoning Attacks in Episodic Reinforcement Learning**

**Anshuka Rangi**[1] ,  **Haifeng Xu**[2] ,  **Long Tran-Thanh**[3] ,  **Massimo Franceschetti**[1]
[1] University of California San Diego, USA
[2]University of Virginia, USA
[3]University of Warwick, UK
{arangi, mfranceschetti}@ucsd.edu, hx4ad@virginia.edu, long.tran-thanh@warwick.ac.uk

## Defense against a specific reward attack algorithm

**Defense Against Reward Poisoning Attacks in Reinforcement Learning**

**Kiarash Banihashem**         **Adish Singla**         **Goran Radanovic**
MPI-SWS                        MPI-SWS                  MPI-SWS
kbanihas@mpi-sws.org    adishs@mpi-sws.org    gradanovic@mpi-sws.org

Not robust! Attacker can change its algorithm later.

# The Attack Problem

# The Attack Problem

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

# The Attack Problem

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

We decompose the attacked $\pi$-M interaction based on the *flow of information*.

# Adversarial Decomposition

We decompose the attacked $\pi$-M interaction based on the *flow of information*.

# Adversarial Decomposition

We decompose the attacked $\pi$-M interaction based on the *flow of information*.

# Adversarial Decomposition

We decompose the attacked $\pi$-M interaction based on the *flow of information.*

# Adversarial Decomposition

We decompose the attacked $\pi$-M interaction based on the *flow of information*.

# Adversarial Decomposition

We decompose the attacked $\pi$-M interaction based on the *flow of information*.

# Adversarial Decomposition

We decompose the attacked $\pi$-M interaction based on the *flow of information*.

# Adversarial Decomposition

We decompose the attacked $\pi$-M interaction based on the *flow of information*.

# Adversarial Decomposition

We decompose the attacked $\pi$-M interaction based on the *flow of information*.

# Adversarial Decomposition

We decompose the attacked $\pi$-M interaction based on the *flow of information*.

# Adversarial Decomposition

We decompose the attacked $\pi$-M interaction based on the *flow of information*.

# Adversarial Decomposition

We decompose the attacked $\pi$-M interaction based on the *flow of information*.

# Adversarial Decomposition

We decompose the attacked $\pi$-M interaction based on the *flow of information*.

# Adversarial Decomposition

We decompose the attacked $\pi$-M interaction based on the *flow of information*.

# Adversarial Decomposition

We decompose the attacked $\pi$-M interaction based on the *flow of information*.

# Adversarial Decomposition

We decompose the attacked $\pi$-M interaction based on the *flow of information*.

# Adversarial Decomposition

We decompose the attacked $\pi$-M interaction based on the *flow of information*.

# Adversarial Decomposition

We decompose the attacked $\pi$-M interaction based on the *flow of information*.

# Adversarial Decomposition

We decompose the attacked $\pi$-M interaction based on the *flow of information*.

# Adversarial Decomposition

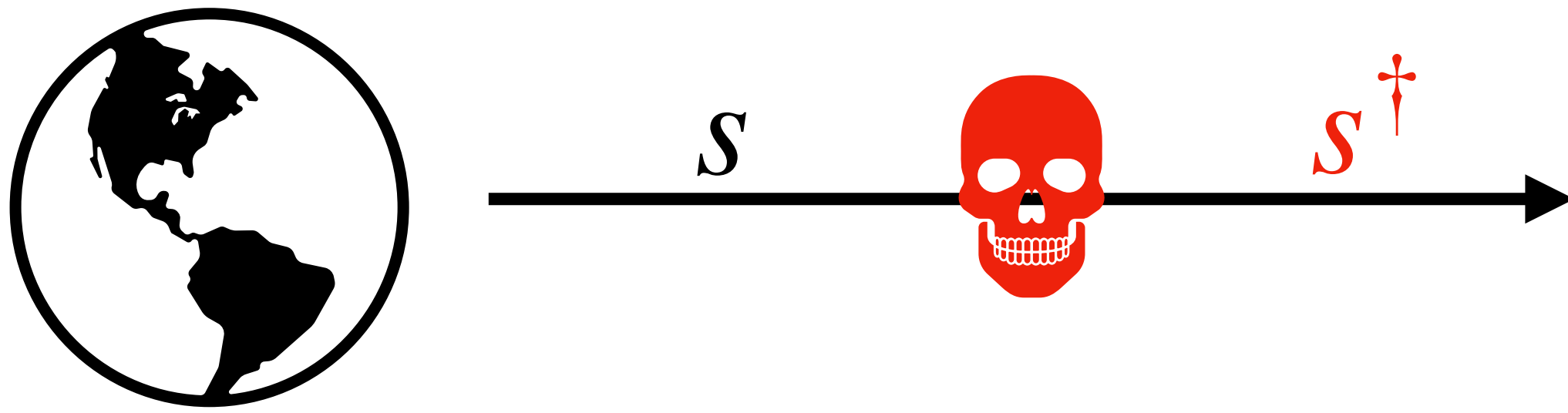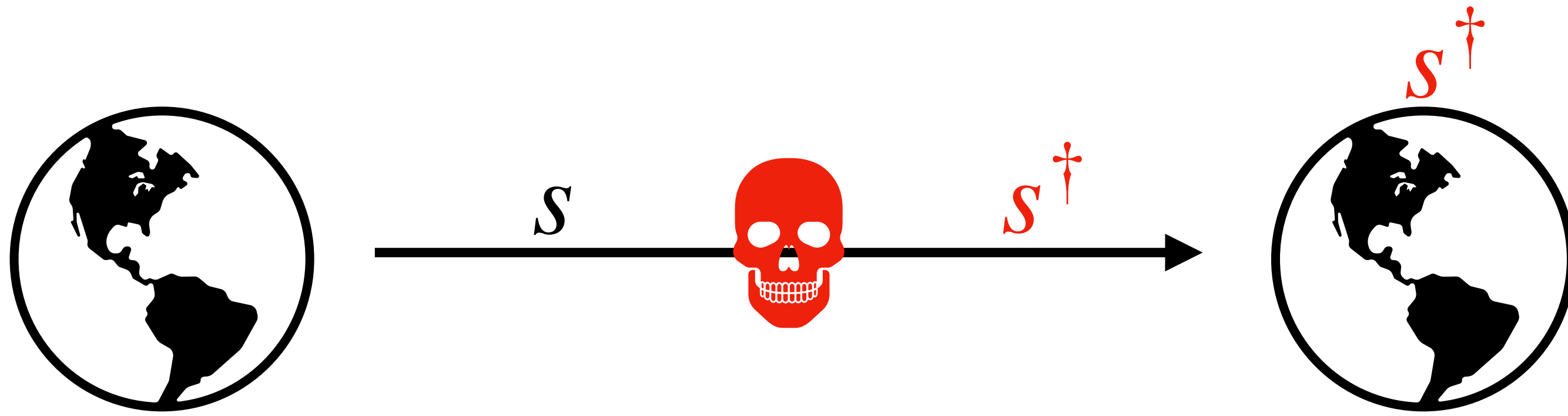We decompose the attacked $\pi$-M interaction based on the *flow of information*.

# Adversarial Decomposition

We decompose the attacked $\pi$-M interaction based on the *flow of information*.

# Adversarial Decomposition

We decompose the attacked $\pi$-M interaction based on the *flow of information*.

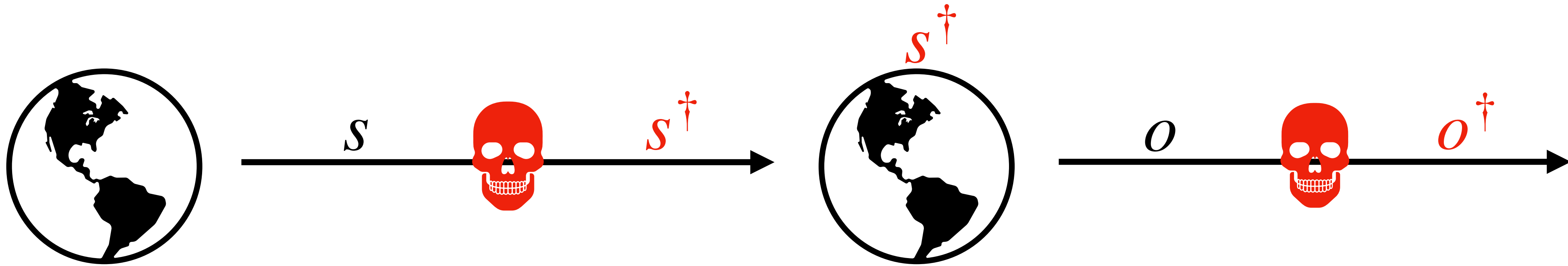# Attacker's Perspective

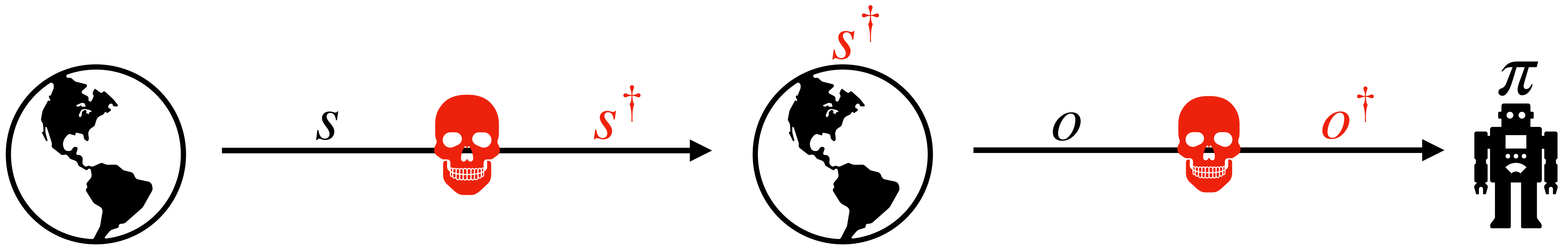

$t_1$

Knowledge: $s$

Attack: $s^{\dagger} \in S$

$t_2$

Knowledge: $s, o$

Attack: $o^{\dagger} \in O$

$t_4$

Knowledge: $s, o, a, r$

Attack: $r^{\dagger} \in R$

$t_3$

Knowledge: $s, o, a$

Attack: $a^{\dagger} \in A$

# Attacker's Perspective



$t_1$

Knowledge: $s$

Attack: $s^\dagger \in S$

$t_2$

Knowledge: $s, o$

Attack: $o^\dagger \in O$

$t_4$

Knowledge: $s, o, a, r$

Attack: $r^\dagger \in R$

$t_3$

Knowledge: $s, o, a$

Attack: $a^\dagger \in A$

Describes MDP $\overline{M}$ for the attacker!

# Meta MDP

# Meta MDP

Attacker's interaction with $\pi$ and $M$ evolves according to MDP $\bar{M}$.

# Meta MDP

Attacker's interaction with $\pi$ and $M$ evolves according to MDP $\bar{M}$.

1. $\bar{S}$ records the attacker's information at any subperiod:

# Meta MDP

Attacker's interaction with $\pi$ and $M$ evolves according to MDP $\bar{M}$.

1. $\bar{S}$ records the attacker's information at any subperiod:

$$\overline{\mathcal{S}} = \mathcal{S} \cup (\mathcal{S} \cup \mathcal{O}) \cup (\mathcal{S} \cup \mathcal{O} \cup \mathcal{A}) \cup (\mathcal{S} \cup \mathcal{O} \cup \mathcal{A} \cup \mathcal{R})$$

# Meta MDP

Attacker's interaction with $\pi$ and $M$ evolves according to MDP $\bar{M}$.

1. $\bar{S}$ records the attacker's information at any subperiod:

$$\overline{\mathcal{S}} = \mathcal{S} \cup (\mathcal{S} \cup \mathcal{O}) \cup (\mathcal{S} \cup \mathcal{O} \cup \mathcal{A}) \cup (\mathcal{S} \cup \mathcal{O} \cup \mathcal{A} \cup \mathcal{R})$$

2. $\bar{A}$ captures the attacks available at any subperiod:

# Meta MDP

Attacker's interaction with $\pi$ and $M$ evolves according to MDP $\bar{M}$.

1. $\bar{S}$ records the attacker's information at any subperiod:

$$\overline{S} = S \cup (S \cup O) \cup (S \cup O \cup A) \cup (S \cup O \cup A \cup R)$$

2. $\bar{A}$ captures the attacks available at any subperiod:

$$\overline{A}(s) \subseteq S, \ \overline{A}(s,o) \subseteq O, \ \overline{A}(s,o,a) \subseteq A, \ \overline{A}(s,o,a,r) \subseteq R$$

# Meta MDP

Attacker's interaction with $\pi$ and $M$ evolves according to MDP $\bar{M}$.

1. $\bar{S}$ records the attacker's information at any subperiod:

$$\overline{\mathcal{S}} = \mathcal{S} \cup (\mathcal{S} \cup \mathcal{O}) \cup (\mathcal{S} \cup \mathcal{O} \cup \mathcal{A}) \cup (\mathcal{S} \cup \mathcal{O} \cup \mathcal{A} \cup \mathcal{R})$$

2. $\bar{A}$ captures the attacks available at any subperiod:

$$\overline{\mathcal{A}}(s) \subseteq \mathcal{S}, \ \overline{\mathcal{A}}(s, o) \subseteq \mathcal{O}, \ \overline{\mathcal{A}}(s, o, a) \subseteq \mathcal{A}, \ \overline{\mathcal{A}}(s, o, a, r) \subseteq \mathcal{R}$$

3. Transitions capture the evolution of information.

# Meta MDP

Attacker's interaction with $\pi$ and $M$ evolves according to MDP $\bar{M}$.

1. $\bar{S}$ records the attacker's information at any subperiod:

$$\overline{S} = S \cup (S \cup O) \cup (S \cup O \cup A) \cup (S \cup O \cup A \cup R)$$

2. $\bar{A}$ captures the attacks available at any subperiod:

$$\overline{A}(s) \subseteq S, \ \overline{A}(s, o) \subseteq O, \ \overline{A}(s, o, a) \subseteq A, \ \overline{A}(s, o, a, r) \subseteq R$$

3. Transitions capture the evolution of information.

**Proposition:** Any optimal policy for $\overline{M}$ is an optimal attack policy.

# Reduction to RL

# Reduction to RL

# Reduction to RL

# Reduction to RL

# Reduction to RL



Optimal attacks can be computed using standard RL techniques!

# Computational Efficiency

# Computational Efficiency

$$|\bar{S}| \leq SOAR \quad \text{and} \quad |\bar{A}| \leq S + O + A + R$$

# Computational Efficiency

$$|\bar{S}| \le SOAR \quad \text{and} \quad |\bar{A}| \le S + O + A + R$$

$\overline{M}$ has only polynomially larger state and action space than $M$.

# Computational Efficiency

$$|\bar{S}| \leq SOAR \quad \text{and} \quad |\bar{A}| \leq S + O + A + R$$

## Attacking RL *efficiently* reduces to RL!

$\bar{M}$ has only polynomially larger state and action space than $M$.

Can we defend against attacks?

# Defense

# 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 interaction under the worst-case attack:

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

Avoids Cat and Mouse Game!

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

# Reduction to MARL

# Reduction to MARL

# Reduction to MARL

Two player MG

# Reduction to MARL



Two player MG

Defense corresponds to a Weak Stackelberg Equilibrium (WSE).

# Challenges

# Challenges

- WSE need not exist.

# Challenges

- WSE need not exist.

- WSE are generally non-Markovian!

# Challenges

- WSE need not exist.

- WSE are generally non-Markovian!

**Proposition:** The defense problem is as hard as solving POMDPs.

Thus, the defense problem is NP-hard to even approximate.

# Special Structure: Sequential Play

# Special Structure: Sequential Play

Key: restrict observation attacks.

# Special Structure: Sequential Play

Key: restrict observation attacks.

P2

$s$

# Special Structure: Sequential Play

Key: restrict observation attacks.

P2

$$s \xrightarrow{\quad s^{\dagger} \quad}$$

# Special Structure: Sequential Play

Key: restrict observation attacks.

P2　　　　　　　P1

$$s \xrightarrow{\textcolor{red}{s^\dagger}} s^\dagger$$

# Special Structure: Sequential Play

Key: restrict observation attacks.

P2       P1

$$s \xrightarrow{\;\textcolor{red}{s^\dagger}\;} s^\dagger \xrightarrow{\;a\;}$$

# Special Structure: Sequential Play

Key: restrict observation attacks.

$$\underset{\text{P2}}{s} \xrightarrow{\textcolor{red}{s^\dagger}} \underset{\text{P1}}{s^\dagger} \xrightarrow{a} \underset{\text{P2}}{(s^\dagger, a)}$$

# Special Structure: Sequential Play

Key: restrict observation attacks.

$$P2 \qquad\qquad P1 \qquad\qquad\qquad P2$$

$$s \quad \xrightarrow{\ \textcolor{red}{s^\dagger}\ } \quad s^\dagger \quad \xrightarrow{\ a\ } \quad (s^\dagger, a) \quad \xrightarrow{\ \textcolor{red}{a^\dagger}\ }$$

# Special Structure: Sequential Play

Key: restrict observation attacks.

$$\underset{s}{\overset{\text{P2}}{s}} \xrightarrow{\textcolor{red}{s^\dagger}} \underset{s^\dagger}{\overset{\text{P1}}{s^\dagger}} \xrightarrow{a} \underset{(s^\dagger, a)}{\overset{\text{P2}}{(s^\dagger, a)}} \xrightarrow{\textcolor{red}{a^\dagger}} \underset{(s^\dagger, a, r)}{\overset{\text{P2}}{(s^\dagger, a, r)}}$$

# Special Structure: Sequential Play

Key: restrict observation attacks.

$$\text{P2} \quad \xrightarrow{\color{red}s^\dagger} \quad \text{P1} \quad \xrightarrow{a} \quad \text{P2} \quad \xrightarrow{\color{red}a^\dagger} \quad \text{P2}$$

$$s \xrightarrow{\ \color{red}s^\dagger\ } s^\dagger \xrightarrow{\ a\ } (s^\dagger, a) \xrightarrow{\ \color{red}a^\dagger\ } (s^\dagger, a, r) \xrightarrow{\ \color{red}r^\dagger\ }$$

# Special Structure: Sequential Play

Key: restrict observation attacks.

P2     P1     P2     P2

$$s \xrightarrow{\ s^\dagger\ } s^\dagger \xrightarrow{\ a\ } (s^\dagger, a) \xrightarrow{\ a^\dagger\ } (s^\dagger, a, r) \xrightarrow{\ r^\dagger\ }$$

Game evolves like a *turn-based* Markov game $\overline{G}$.

# Meta Turn-based Markov Game

1. $\bar{S}$ records the player's information at any subperiod:

$$\overline{\mathcal{S}}_1 = \mathcal{S} \quad \text{and} \quad \overline{\mathcal{S}}_2 = \mathcal{S} \cup (\mathcal{S} \cup \mathcal{A}) \cup (\mathcal{S} \cup \mathcal{A} \cup \mathcal{R})$$

2. $\bar{A}$ captures the actions available at any subperiod:

$$\overline{\mathcal{A}}_1 = \mathcal{A} \quad \text{and} \quad \overline{\mathcal{A}}_2(s) \subseteq \mathcal{S}, \, \overline{\mathcal{A}}_2(s, a) \subseteq \mathcal{A}, \, \overline{\mathcal{A}}_2(s, a, r) \subseteq \mathbb{R}$$

3. Transitions capture the evolution of information.

# Meta Turn-based Markov Game

1. $\bar{S}$ records the player's information at any subperiod:

$$\overline{\mathcal{S}}_1 = \mathcal{S} \quad \text{and} \quad \overline{\mathcal{S}}_2 = \mathcal{S} \cup (\mathcal{S} \cup \mathcal{A}) \cup (\mathcal{S} \cup \mathcal{A} \cup \mathcal{R})$$

2. $\bar{A}$ captures the actions available at any subperiod:

$$\overline{\mathcal{A}}_1 = \mathcal{A} \quad \text{and} \quad \overline{\mathcal{A}}_2(s) \subseteq \mathcal{S}, \ \overline{\mathcal{A}}_2(s, a) \subseteq \mathcal{A}, \ \overline{\mathcal{A}}_2(s, a, r) \subseteq \mathbb{R}$$

3. Transitions capture the evolution of information.

**Proposition:** Any WSE for $\overline{G}$ is an optimal defense policy.

# *Efficient* Reduction to MARL

# *Efficient* Reduction to MARL

$\overline{G}$

# *Efficient* Reduction to MARL

$\overline{G}$ | Zero-sum:

# *Efficient* Reduction to MARL

$\overline{G}$

Zero-sum:

Defense $\longrightarrow$ WSE

# *Efficient* Reduction to MARL

$\overline{G}$ | Zero-sum:

Defense → WSE → MPNE

# *Efficient* Reduction to MARL

$\overline{G}$

Zero-sum:

Defense $\rightarrow$ WSE $\rightarrow$ MPNE

General-sum:

# *Efficient* Reduction to MARL

$\overline{G}$

Zero-sum:

Defense → WSE → MPNE

General-sum:

Defense → WSE

# *Efficient* Reduction to MARL

$\overline{G}$

Zero-sum:

Defense $\rightarrow$ WSE $\rightarrow$ MPNE

General-sum:

Defense $\rightarrow$ WSE $\rightarrow$ MPNE+tiebreak

# Rollback Algorithm

# 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:

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

# Guarantees

# Guarantees

**Theorem:** An optimal defense can be computed or learned in polynomial time if <u>observation attacks are not permitted</u>, and

# Guarantees

**Theorem:** An optimal defense can be computed or learned in polynomial time if <u>observation attacks are not permitted</u>, and

1. $\overline{G}$ is zero-sum, or

# Guarantees

**Theorem:** An optimal defense can be computed or learned in polynomial time if <u>observation attacks are not permitted</u>, and

1. $\overline{G}$ is zero-sum, or

2. $\overline{G}$ has finite-horizon.

# Guarantees

**Theorem:** An optimal defense can be computed or learned in polynomial time if <u>observation attacks are not permitted</u>, and

Complete characterization: hard $\Longleftrightarrow$ observation attacks!

2. $\overline{G}$ has finite-horizon.

# Conclusions

# Conclusions

- Optimal attacks can be efficiently computed for all attack surfaces.

# Conclusions

- Optimal attacks can be efficiently computed for all attack surfaces.

- The defense problem is NP-hard to even approximate.

# Conclusions

- Optimal attacks can be efficiently computed for all attack surfaces.

- The defense problem is NP-hard to even approximate.

- Absent observation attacks, optimal defenses can be efficiently computed.

# Conclusions

- Optimal attacks can be efficiently computed for all attack surfaces.

- The defense problem is NP-hard to even approximate.

- Absent observation attacks, optimal defenses can be efficiently computed.