

# Entanglement of Learning and Security: Strategic Learning for Adaptive Defense

Quanyan Zhu

Tandon School of Engineering  
New York University

CDC Workshop on Learning, Games, and Control for Security of  
Cyber-physical Systems

Dec. 10, 2019



**NYU**

**Systems and Control Group**



# Introduction

- Cyber defense needs to be proactive and adaptive to changing environment.
- There is need for scientific and engineering foundations for online learning for adaptive proactive defense.

- *Moving target defense*: Baseline learning algorithms
- *Proactive defense against APT*: learning of dynamic games with incomplete information
- *Attacker engagement problem*: reinforcement learning
- *Security of reinforcement learning*: deceptive RL

# Game Theory Meets Network Security

## A Tutorial

Stefan Rass and Quanyan Zhu

25th ACM Conference on Computer and Communications Security

October 15, 2018

**UNIVERSITAET  
KLAGENFURT**



<https://arxiv.org/abs/1808.08066>

# Game Theory Meets Network Security

## A Tutorial

Stefan Rass and Quanyan Zhu

# Game Theory for Cyber Deception

---

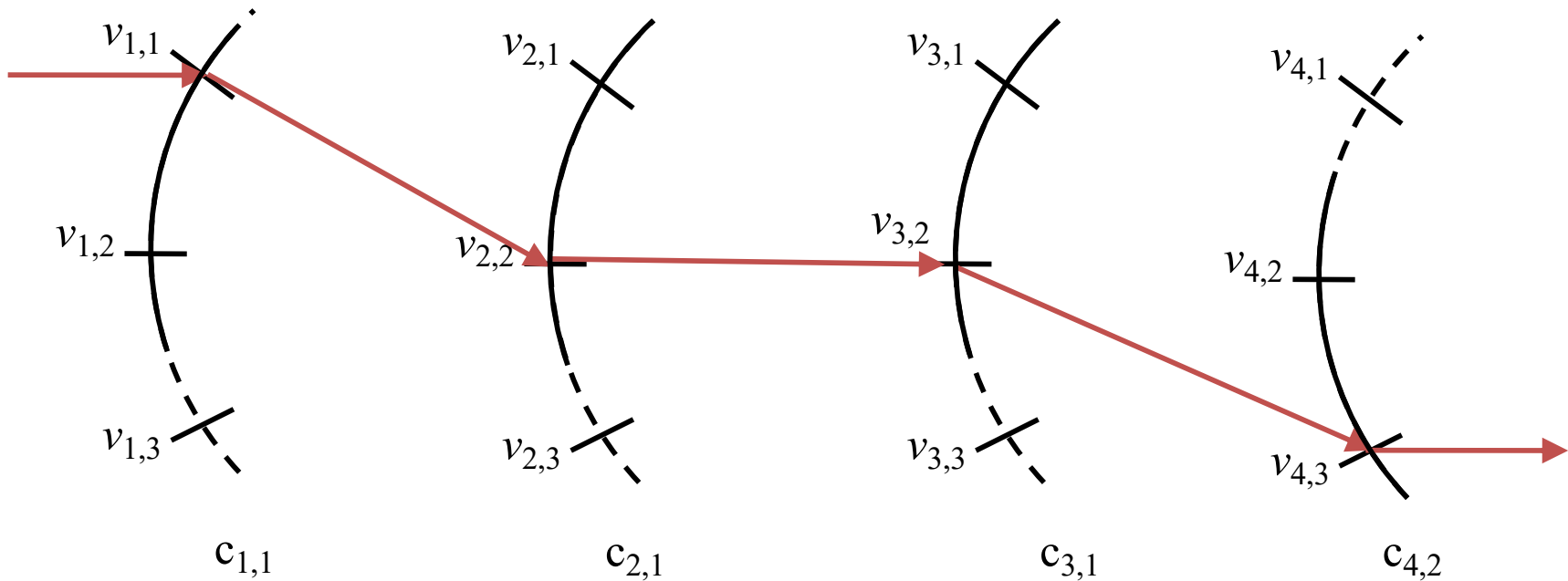
Quanyan Zhu

Symposium and Bootcamp on the Science of Security (HotSoS)

April 2-3, 2019

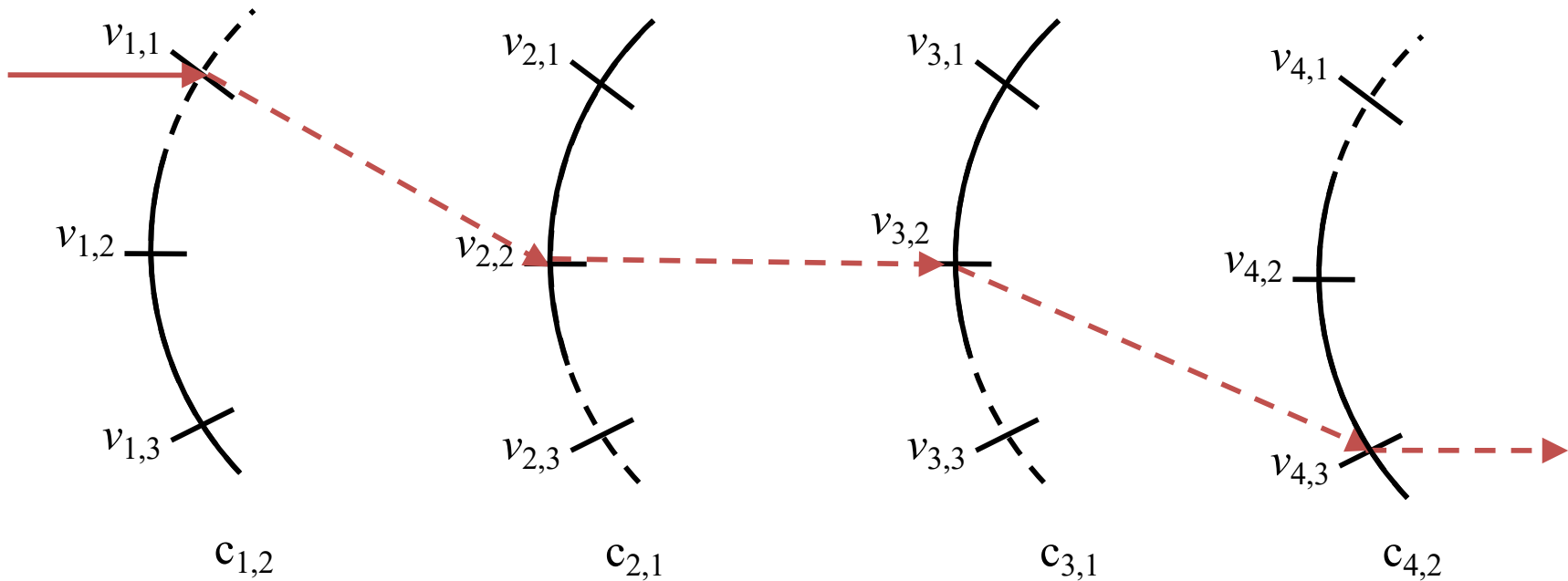
Nashville, Tennessee

# Abstraction



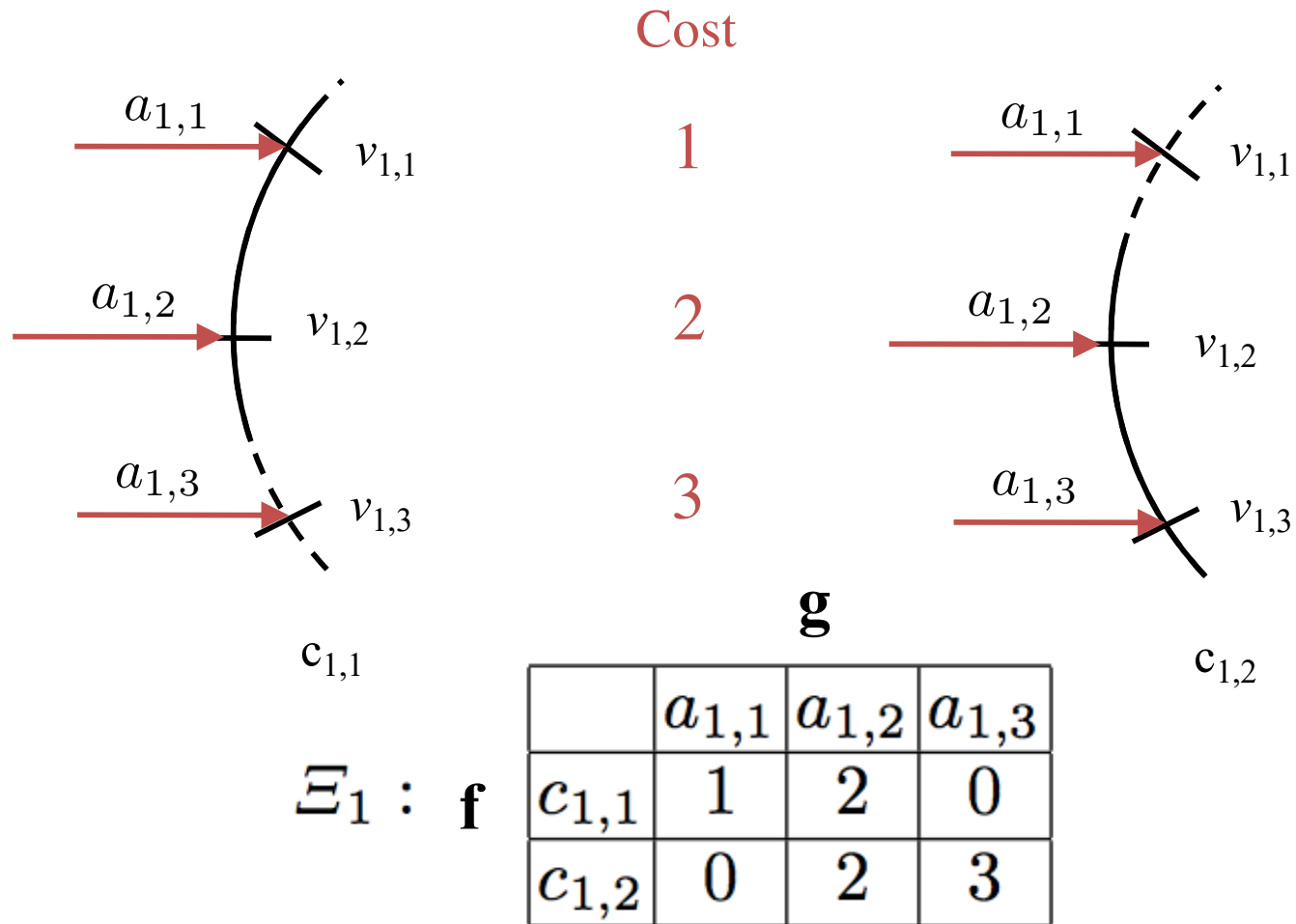
- The system has three vulnerabilities:  $v_{1,1}$ ,  $v_{1,2}$ ,  $v_{1,3}$  at stage 1.
- The system is configured to  $c_{1,1}$ , which exhibits vulnerabilities  $v_{1,1}$ ,  $v_{1,2}$ .
- The *attack surface* of the configuration is the set of vulnerabilities  $\{v_{1,1}, v_{1,2}\}$ .

# Abstraction



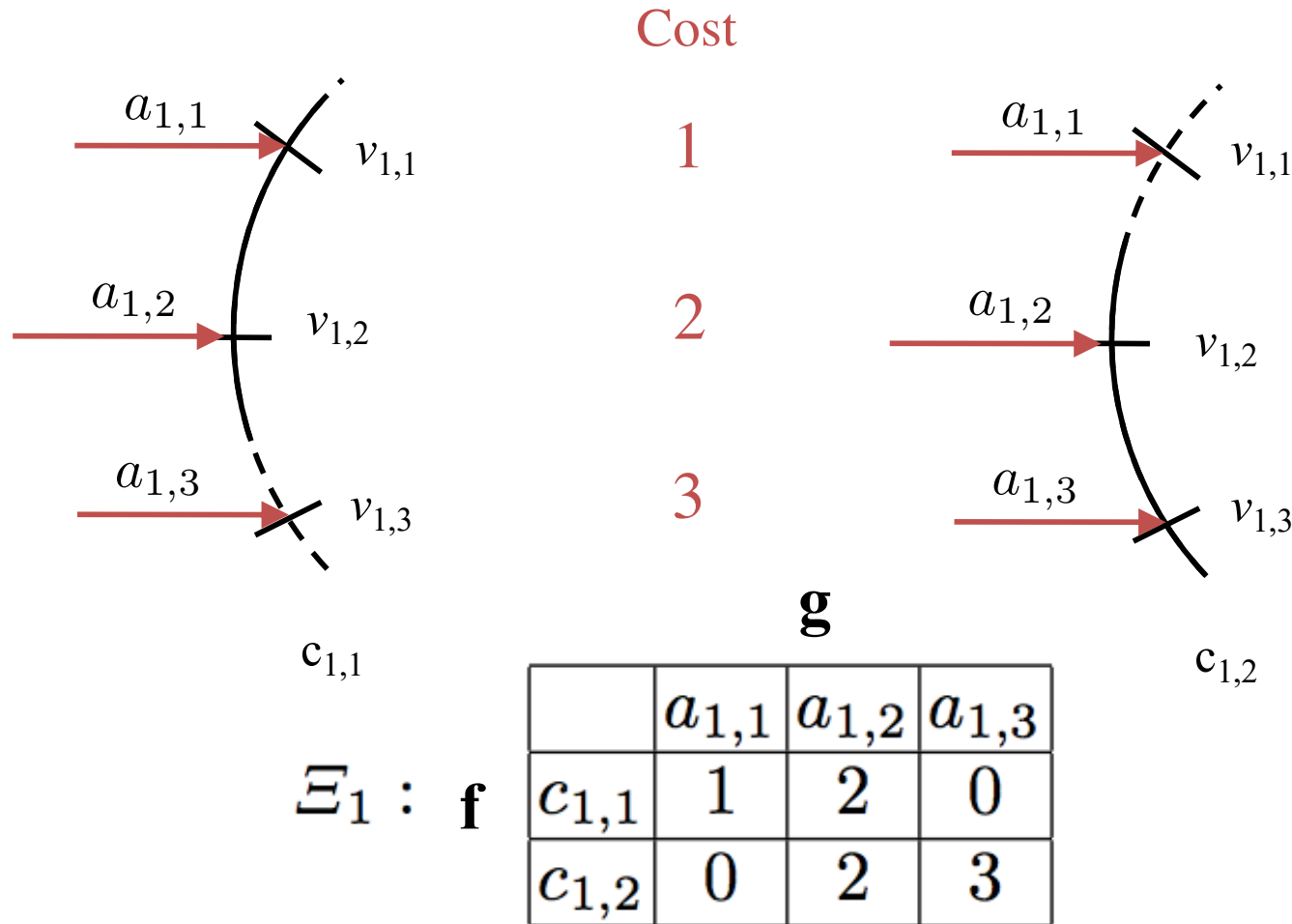
- The system has three vulnerabilities:  $v_{1,1}$ ,  $v_{1,2}$ ,  $v_{1,3}$  at stage 1.
- The system is configured to  $c_{1,2}$ , which exhibits vulnerabilities  $v_{1,2}$ ,  $v_{1,3}$ .
- The *attack surface* of the configuration is the set of vulnerabilities  $\{v_{1,2}, v_{1,3}\}$ .

# Complete Information Game-Theoretic Model



- $\mathbf{f}$  and  $\mathbf{g}$  are mixed strategies. Mixed strategies are *randomizing* strategies.
- Nash equilibrium to the zero-sum game exists, and yields the worst-case defense strategy.

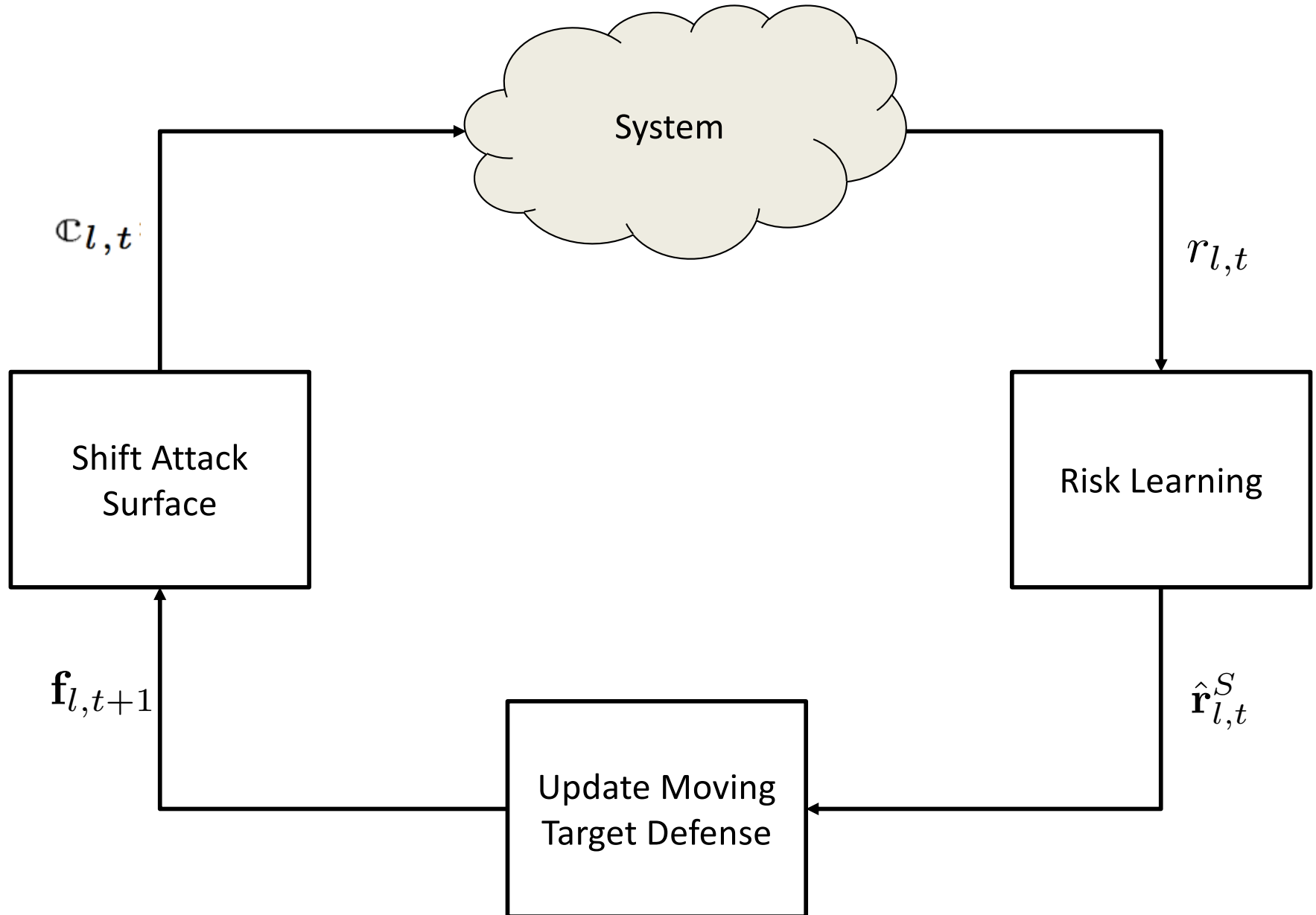
# Not Sufficient Yet!

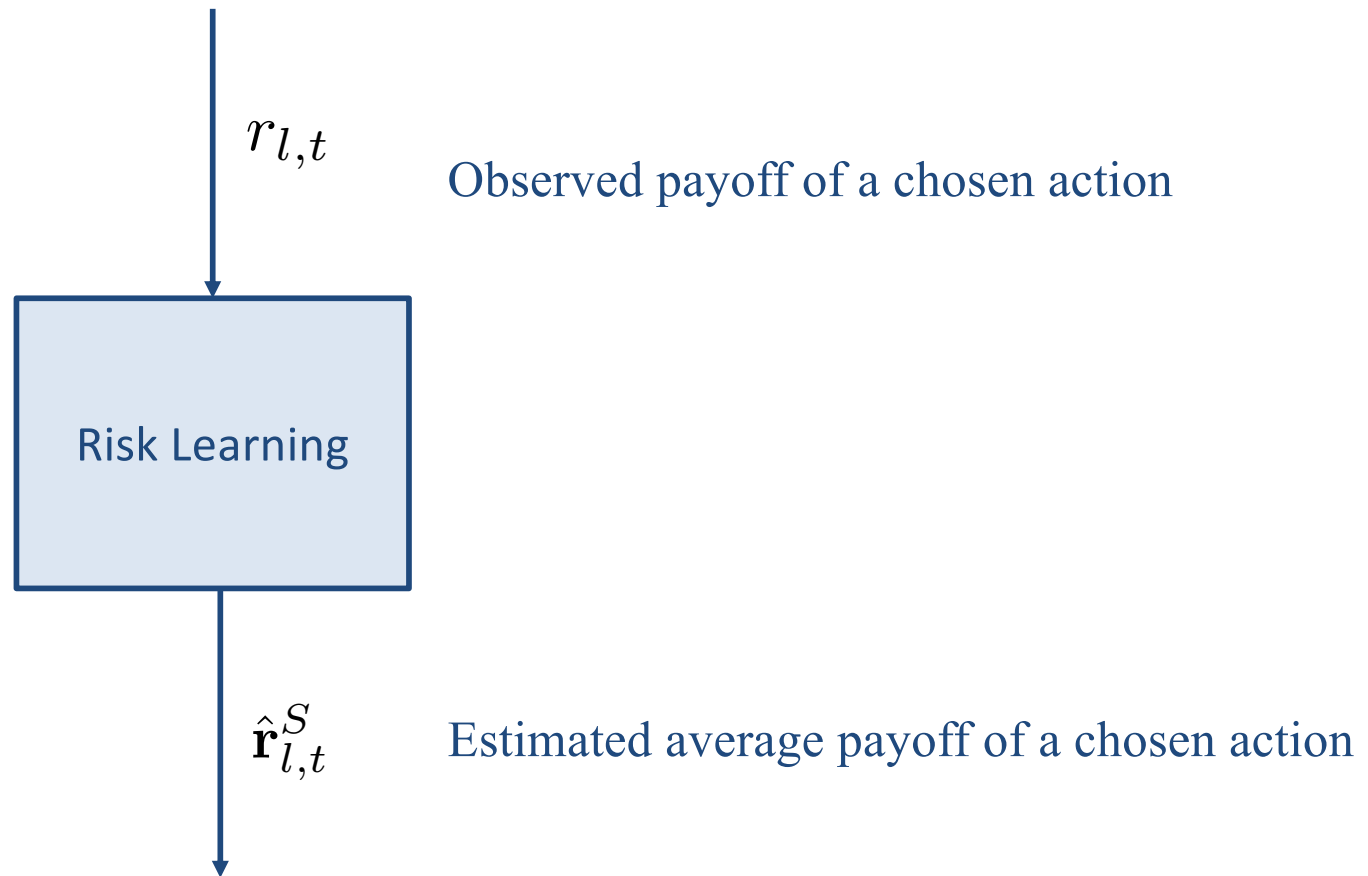


- $\mathbf{f}$  and  $\mathbf{g}$  are mixed strategies. Mixed strategies are *randomizing* strategies.
- Nash equilibrium to the zero-sum game exists, and yields the worst-case defense strategy.



# A Feedback System Model





$$\hat{r}_{l,t+1}^S(c_{l,h}) = \hat{r}_{l,t}^S(c_{l,h}) + \mu_t^S \mathbb{1}_{\{c_{l,t}=c_{l,h}\}}(r_{l,t} - \hat{r}_{l,t}^S(c_{l,h})),$$

$$\hat{r}_{l,t+1}^A(a_{l,h}) = \hat{r}_{l,t}^A(a_{l,h}) + \mu_t^A \mathbb{1}_{\{a_{l,t}=a_{l,h}\}}(r_{l,t} - \hat{r}_{l,t}^A(a_{l,h})).$$

- $\mu_t^S$  and  $\mu_t^A$  are learning rates.
- Averaging over observed payoffs.



$$(\text{SP}) \quad \sup_{\mathbf{f}_{l,t+1} \in \mathcal{F}_l} \underbrace{\langle \mathbf{f}_{l,t+1}, -\hat{\mathbf{r}}_{l,t}^S \rangle}_{\text{Average risk to be minimized}} - \epsilon_{l,t}^S \underbrace{\sum_{h=1}^{m_l} f_{l,h,t+1} \ln \left( \frac{f_{l,h,t+1}}{f_{l,h,t}} \right)}_{\text{Relative entropy: Distance between two distributions}}.$$

Average risk to be minimized

Relative entropy: Distance between two distributions

Cost on changing the strategy (usability)

$$(SP) \quad \sup_{\mathbf{f}_{l,t+1} \in \mathcal{F}_l} \underbrace{\langle \mathbf{f}_{l,t+1}, -\hat{\mathbf{r}}_{l,t}^S \rangle}_{\text{Average risk to be minimized}} - \underbrace{\epsilon_{l,t}^S \sum_{h=1}^{m_l} f_{l,h,t+1} \ln \left( \frac{f_{l,h,t+1}}{f_{l,h,t}} \right)}_{\text{Relative entropy: Distance between two distributions}}.$$

Average risk to be minimized

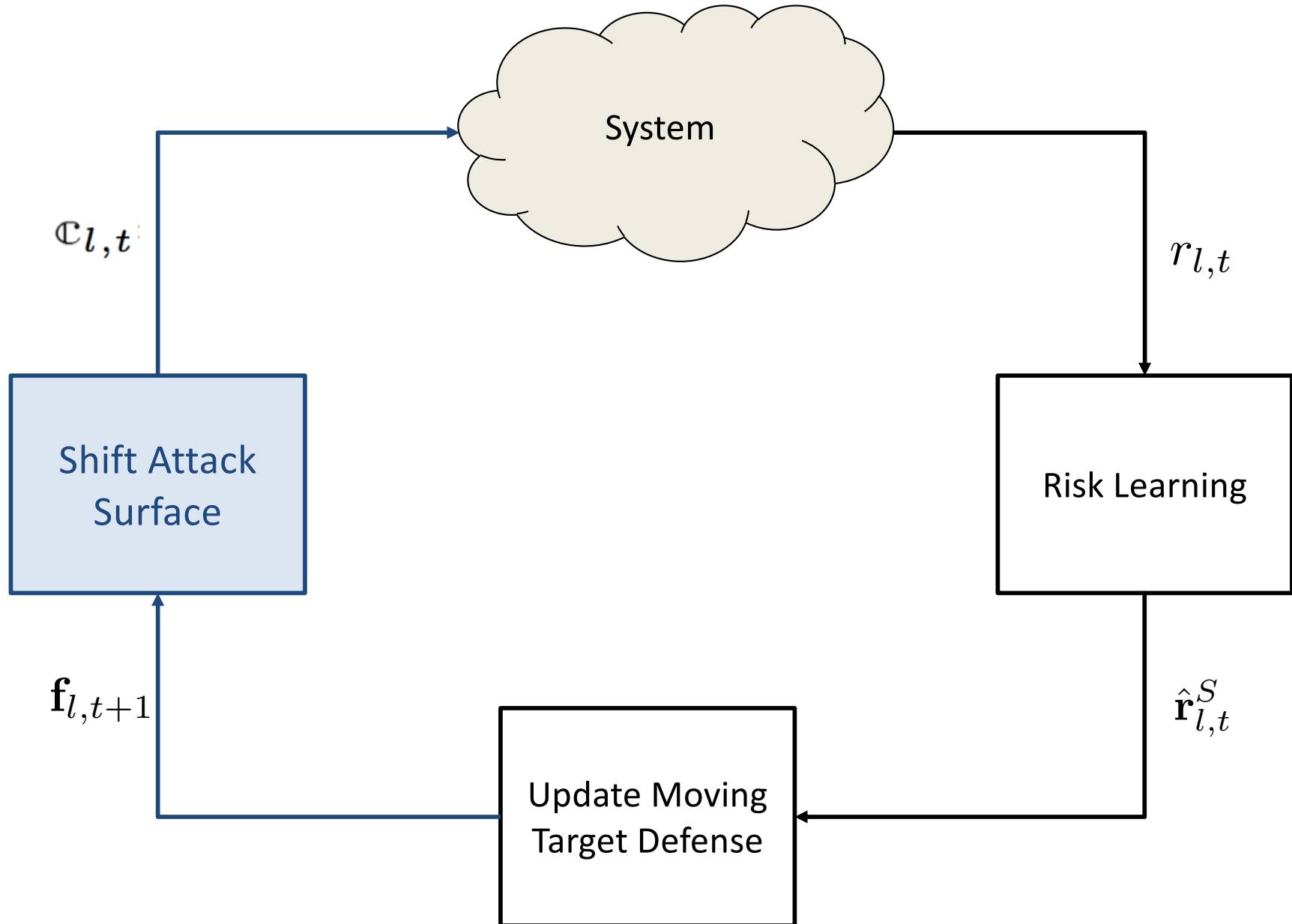
Relative entropy: Distance between two distributions

Cost on changing the strategy (usability)

$$f_{l,h,t+1} = \frac{f_{l,h,t} e^{-\frac{\hat{r}_{l,t}(c_{l,h})}{\epsilon_{l,t}^S}}}{\sum_{h'=1}^{m_l} f_{l,h',t} e^{-\frac{\hat{r}_{l,t}(c_{l,h'})}{\epsilon_{l,t}^S}}}$$

- High risk  $\rightarrow$  low probability
- High  $\epsilon \rightarrow$  Costly to change  $\rightarrow f_{l,h,t}$  (Lower rationality)
- Low  $\epsilon \rightarrow$  Less costly to change (High Rationality)

# A Feedback System Model



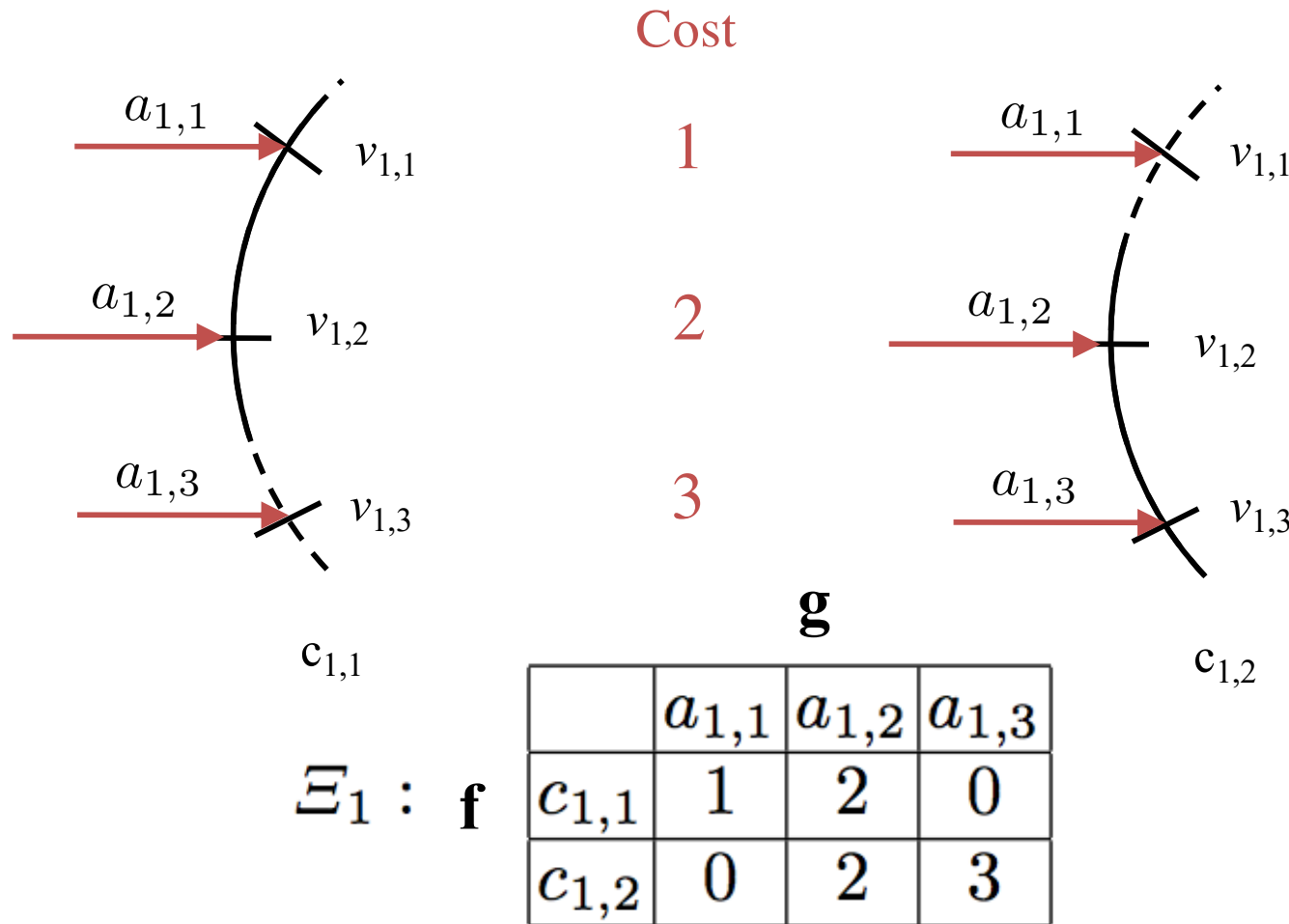
# Mathematical Analysis of the Feedback System

$$\hat{r}_{l,t+1}^S(c_{l,h}) = \hat{r}_{l,t}^S(c_{l,h}) + \mu_t^S \mathbb{1}_{\{c_{l,t}=c_{l,h}\}} (r_{l,t} - \hat{r}_{l,t}^S(c_{l,h}))$$

$$f_{l,h,t+1} = (1 - \lambda_{l,t}^S) f_{l,h,t} + \lambda_{l,t}^S \left( \frac{f_{l,h,t} e^{-\frac{\hat{r}_{l,t}(c_{l,h})}{\epsilon_{l,t}^S}}}{\sum_{h'=1}^{m_l} f_{l,h',t} e^{-\frac{\hat{r}_{l,t}(c_{l,h'})}{\epsilon_{l,t}^S}}} \right)$$

- Use stochastic approximation to show the convergence to an ordinary differential equation (ODE).
- Use ODE to show the convergence of the coupled dynamics to the equilibrium

# Game-Theoretic Model Revisited



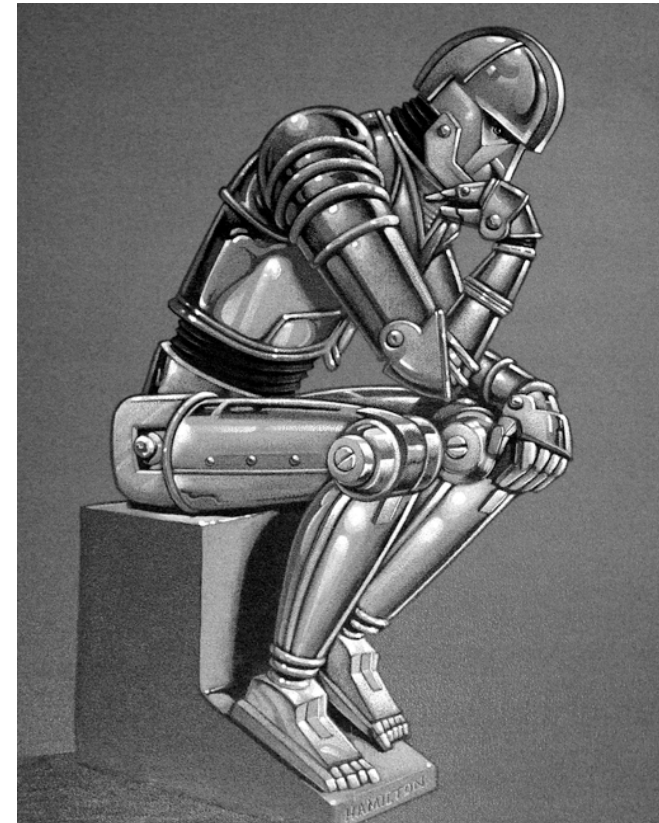
- $\mathbf{f}$  and  $\mathbf{g}$  are mixed strategies. Mixed strategies are *randomizing* strategies.
- Nash equilibrium to the zero-sum game exists, and yields the worst-case defense strategy.

# Boltzmann Learning

$$\begin{cases} \mathbf{f}_{t+1} &= (1 - \lambda_{1,t})\mathbf{f}_t + \lambda_{1,t}\tilde{\beta}_{1,\epsilon}(\hat{\mathbf{u}}_{1,t}) \\ \hat{\mathbf{u}}_{1,t+1} &= \hat{\mathbf{u}}_{1,t} + \frac{\mu_{1,t}}{f_t(a_1)} \mathbb{1}_{\{a_{1,t}=a_1\}} (U_{1,t} - \hat{\mathbf{u}}_{1,t}) \\ \mathbf{g}_{t+1} &= (1 - \lambda_{2,t})\mathbf{g}_t + \lambda_{2,t}\tilde{\beta}_{2,\epsilon}(\hat{\mathbf{u}}_{2,t}) \\ \hat{\mathbf{u}}_{2,t+1} &= \hat{\mathbf{u}}_{2,t} + \frac{\mu_{2,t}}{g_t(a_2)} \mathbb{1}_{\{a_{2,t}=a_2\}} (U_{2,t} - \hat{\mathbf{u}}_{2,t}) \end{cases}$$

**Soft-max function**

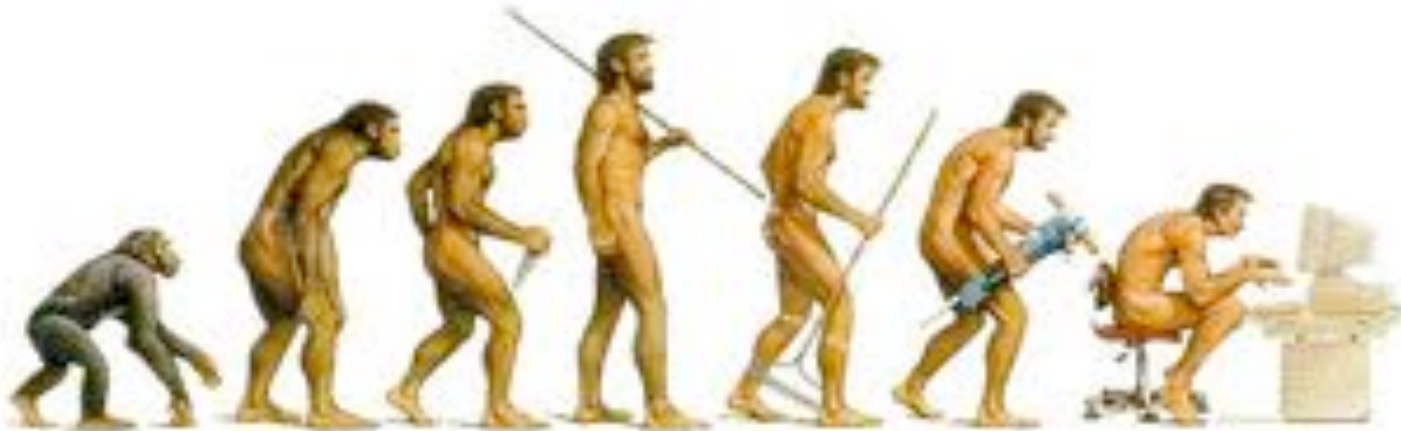
$$\tilde{\beta}_{i,\epsilon}(\hat{\mathbf{u}}_{i,t})(a_i) = \frac{e^{\frac{1}{\epsilon}\hat{u}_{i,t}(a_i)}}{\sum_{a'_i} e^{\frac{1}{\epsilon}\hat{u}_{i,t}(a'_i)}}$$



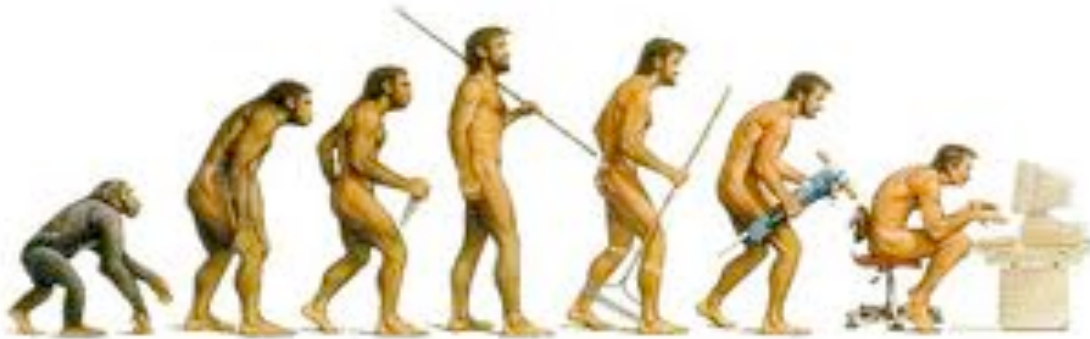


# Replicator Dynamics

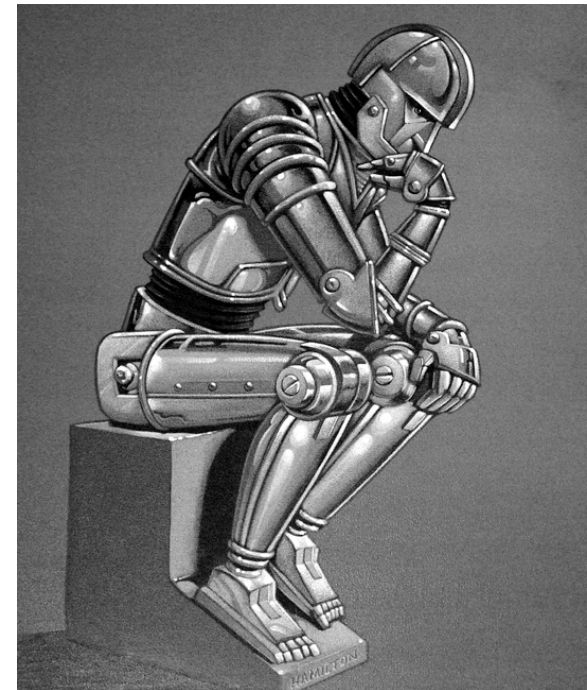
$$\left\{ \begin{array}{l} \mathbf{f}_{t+1} \\ \hat{\mathbf{u}}_{1,t+1} \\ \mathbf{g}_{t+1} \\ \hat{\mathbf{u}}_{2,t+1} \end{array} \right. = \begin{array}{l} \mathbf{f}_t + \lambda_{1,t} U_{1,t} \cdot (\mathbb{1}_{\{a_{1,t}=a_1\}} - \mathbf{f}_t) \\ \hat{\mathbf{u}}_{1,t} + \mu_{1,t} \mathbb{1}_{\{a_{1,t}=a_1\}} (U_{1,t} - \hat{\mathbf{u}}_{1,t}) \\ \mathbf{g}_t + \lambda_{2,t} U_{2,t} \cdot (\mathbb{1}_{\{a_{2,t}=a_2\}} - \mathbf{g}_t) \\ \hat{\mathbf{u}}_{2,t} + \mu_{2,t} \mathbb{1}_{\{a_{2,t}=a_2\}} (U_{2,t} - \hat{\mathbf{u}}_{2,t}) \end{array}$$



# Boltzmann Learning vs. Replicator Dynamics



Robust yet inefficient

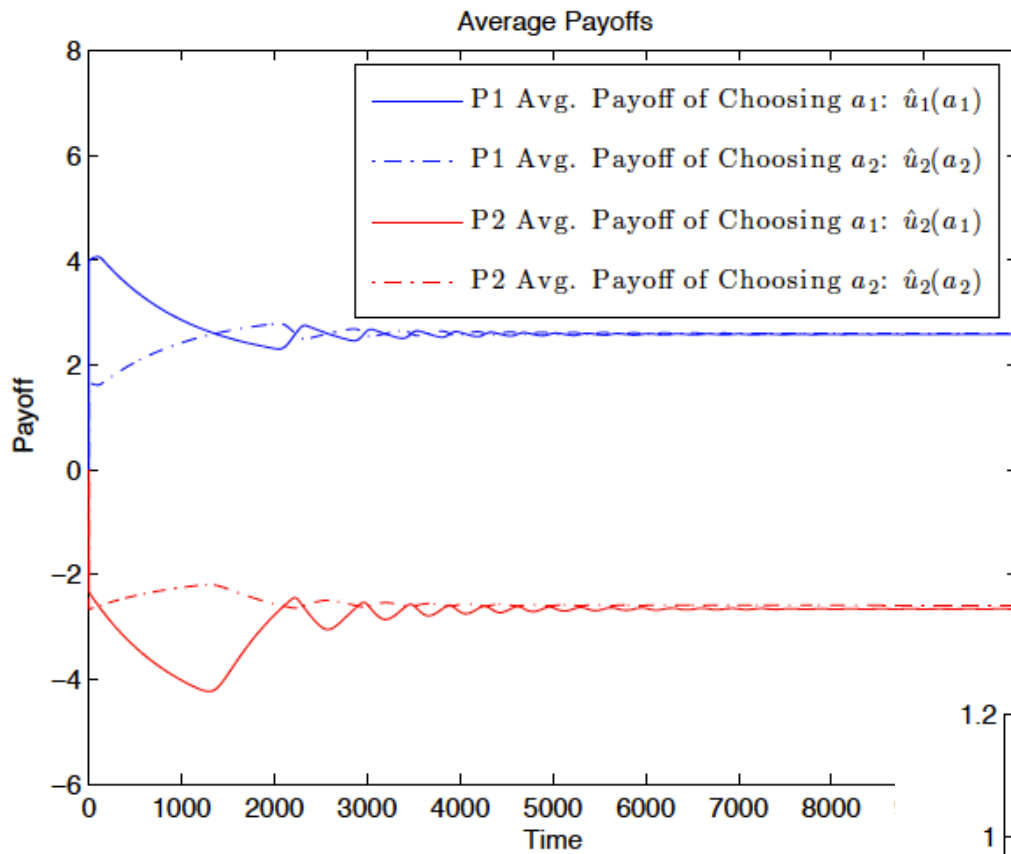


Fragile yet efficient

# Heterogeneous and Hybrid Learning

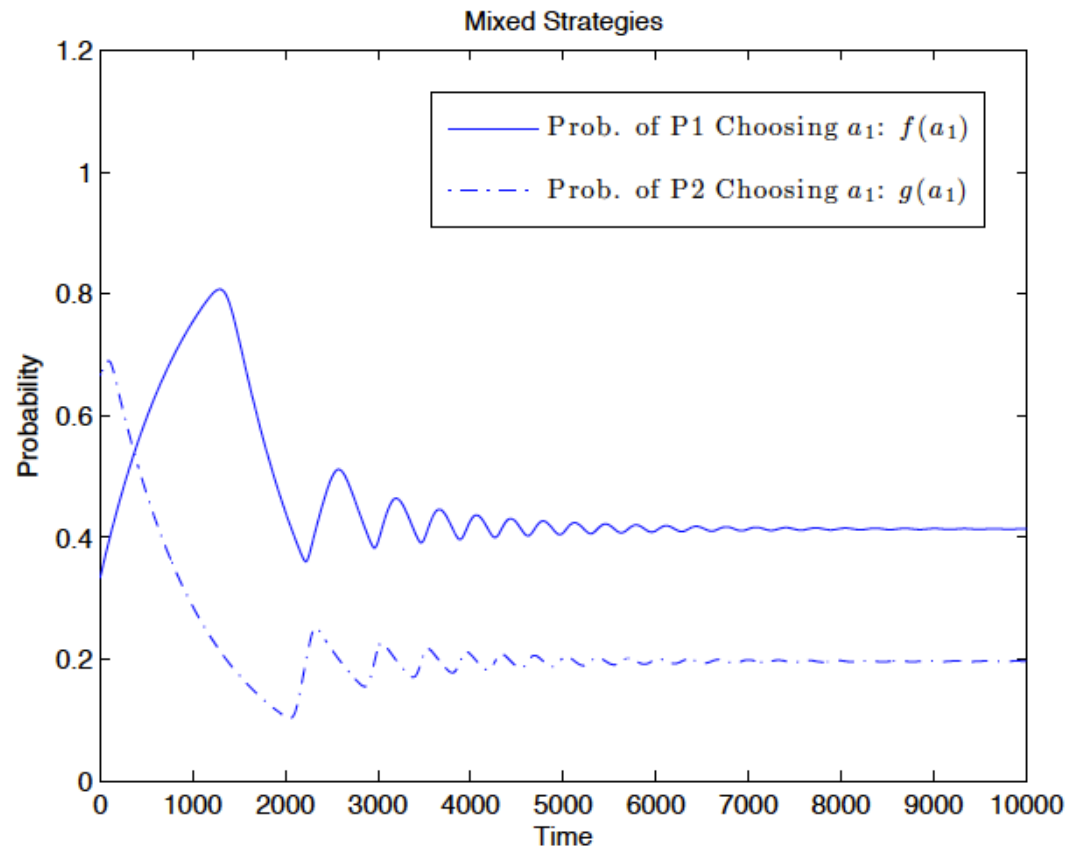
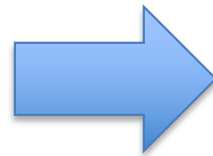
- **Heterogeneous** learning: Different players adopt different learning schemes.
- **Hybrid** learning: Players adopt different learning schemes at different times.

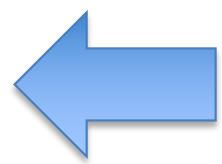
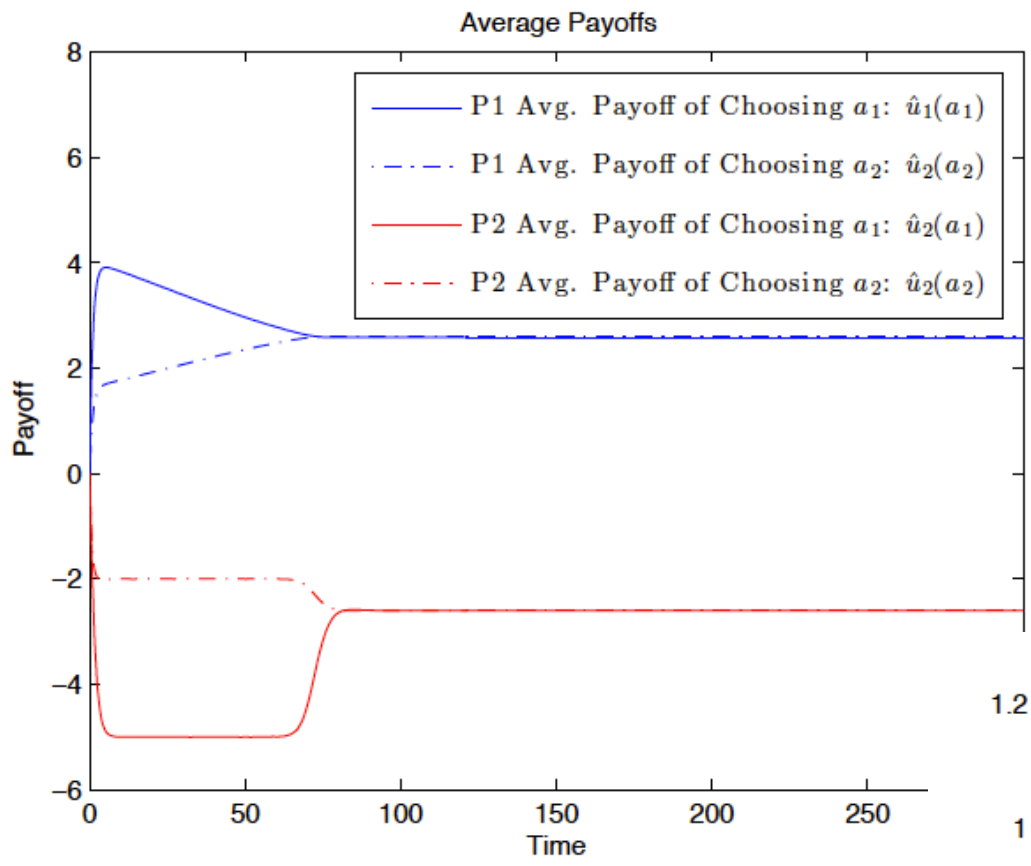
Q. Zhu, H. Tembine and T. Basar, “Hybrid learning in stochastic games and its application in network security,” In F. L. Lewis and D. Liu (Eds.), Reinforcement Learning and Approximate Dynamic Programming for Feedback Control, IEEE Press Computational Intelligence Series, 2012.



**Average Payoffs  
Boltzmann Learning**

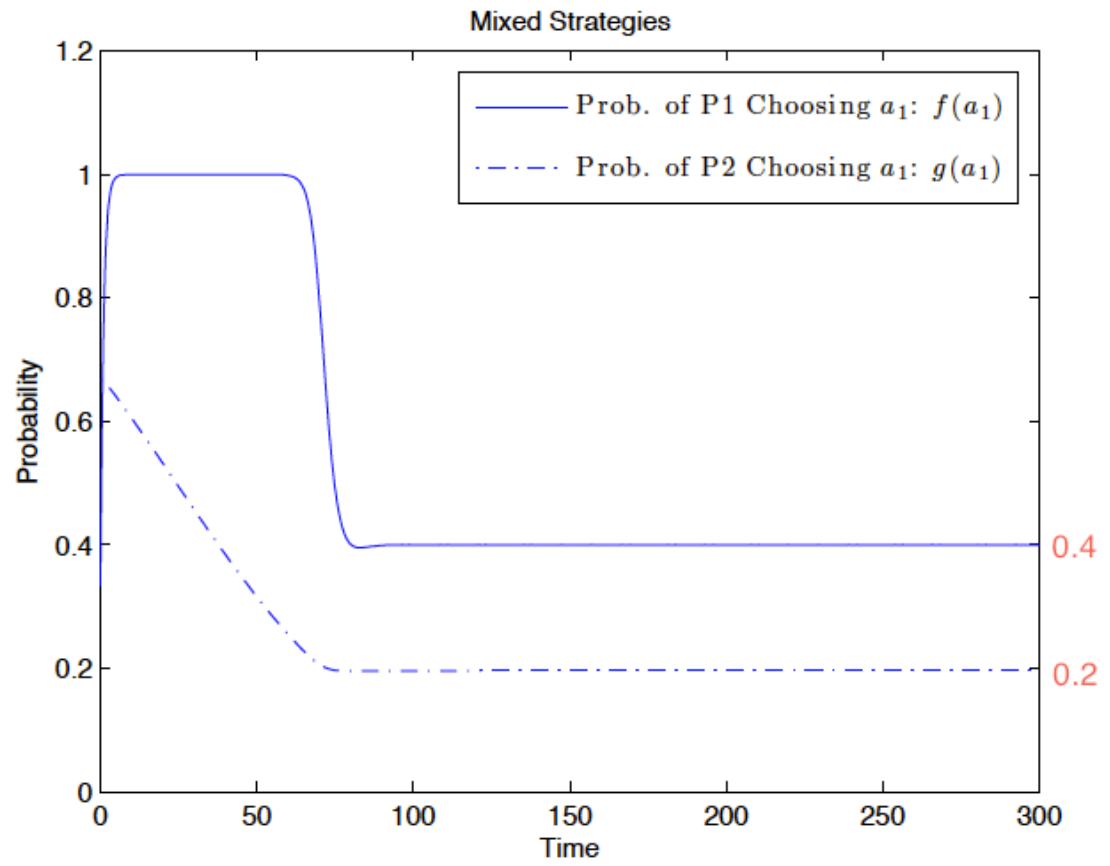
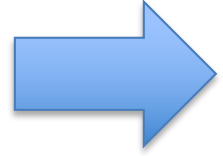
**Mixed Strategies  
Boltzmann Learning**



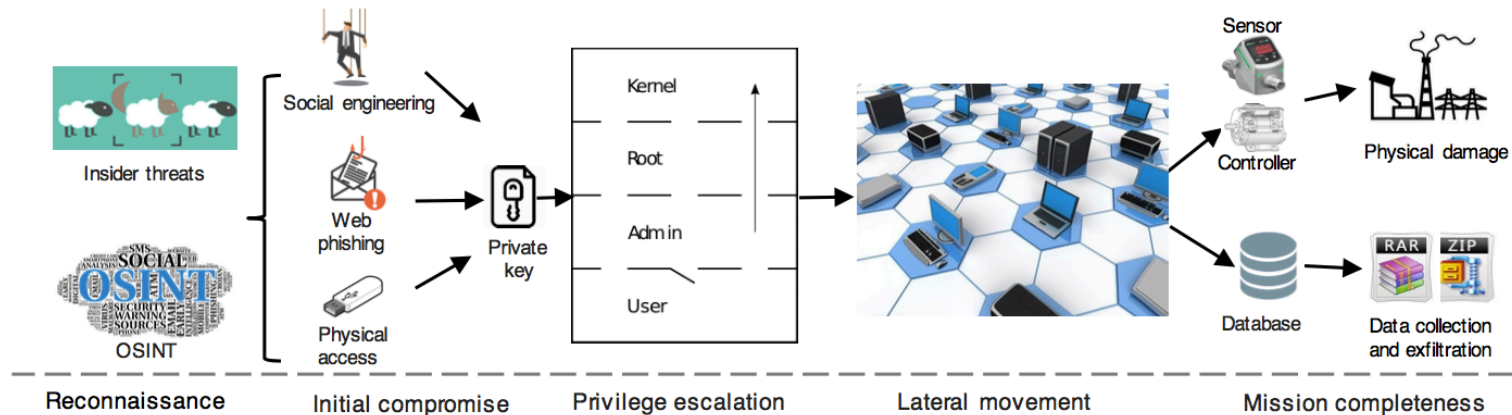


**Average Payoffs  
Boltzmann with Replicator**

**Mixed Strategies  
Boltzmann with Replicator**



# Proactive Defense Against Advanced Persistent Threat



## APT features

- Targeted attack: reconnaissance
- Persistent: firm and patient
- Advanced: technically and strategically

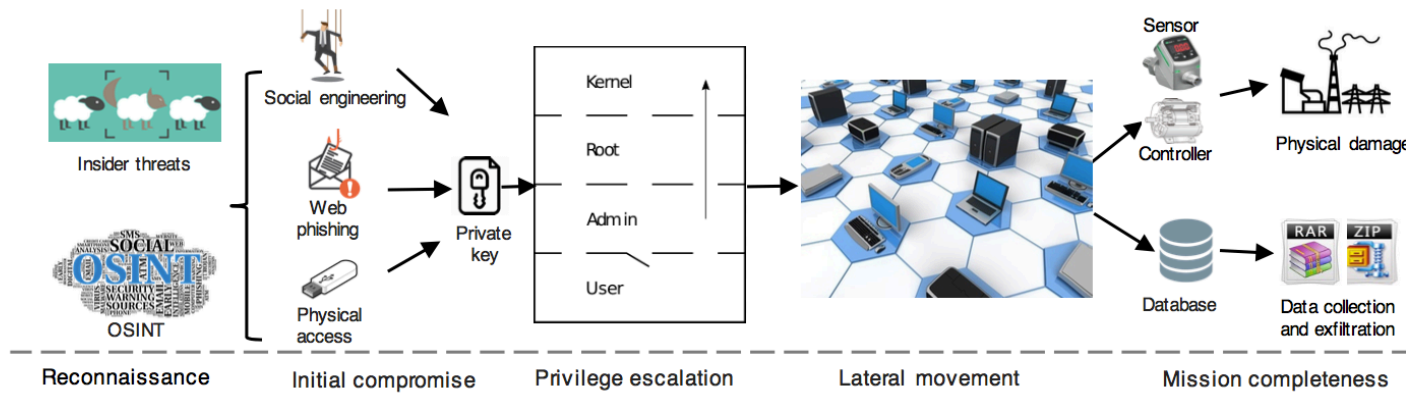


## Cyber Deception

- Stealthy
- Deceptive

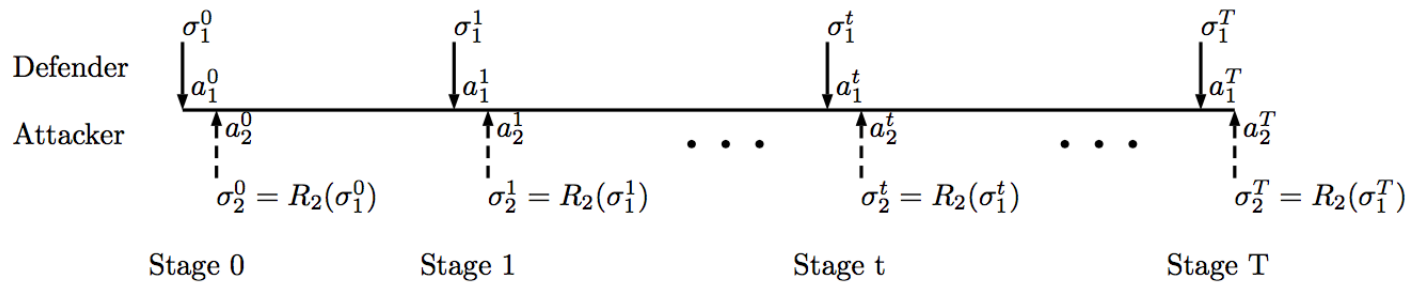
# Multi-Stage Game Framework

Proactive defense: active response prior to the attack



Cross-layer Defense

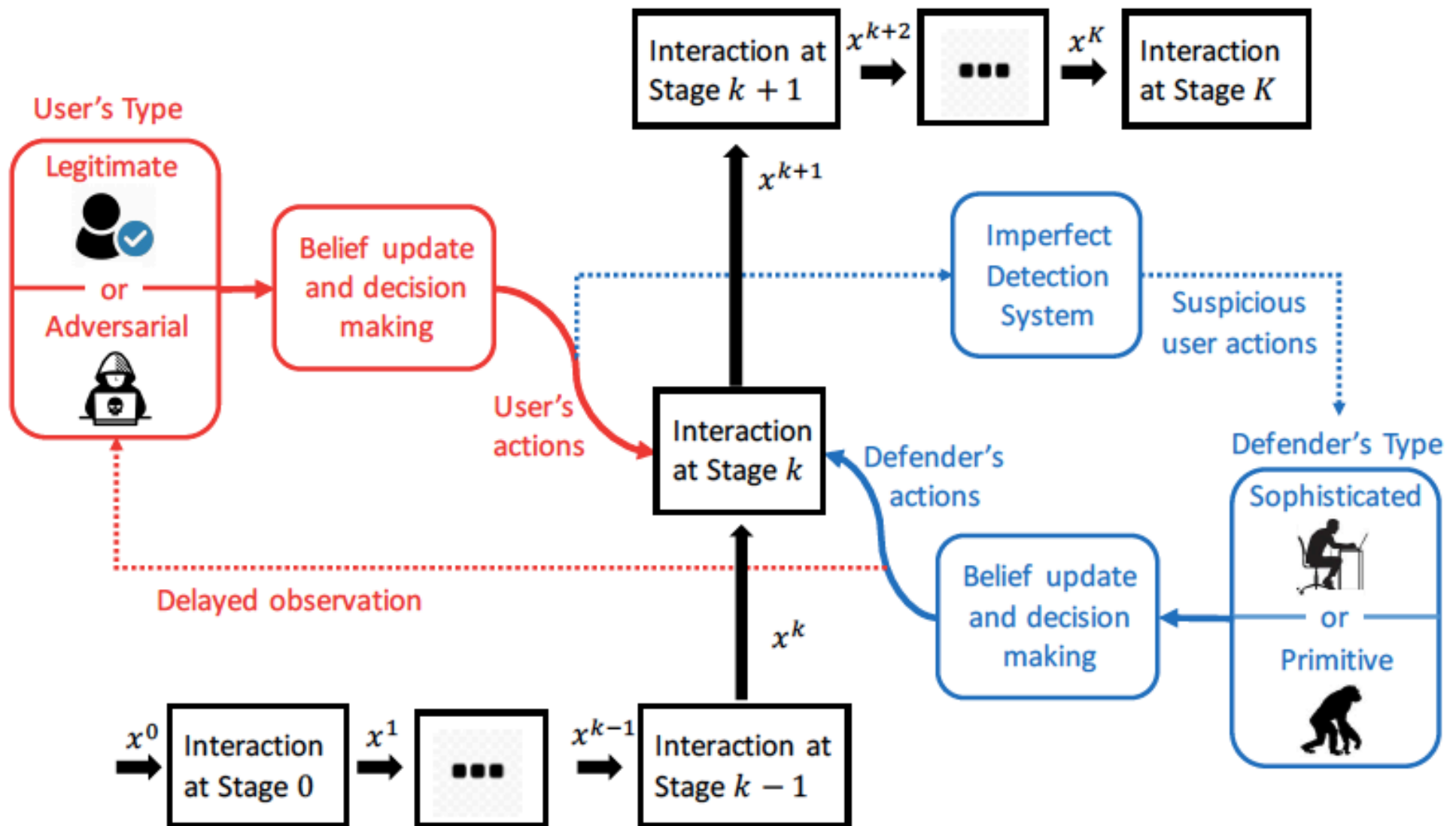
Online learning



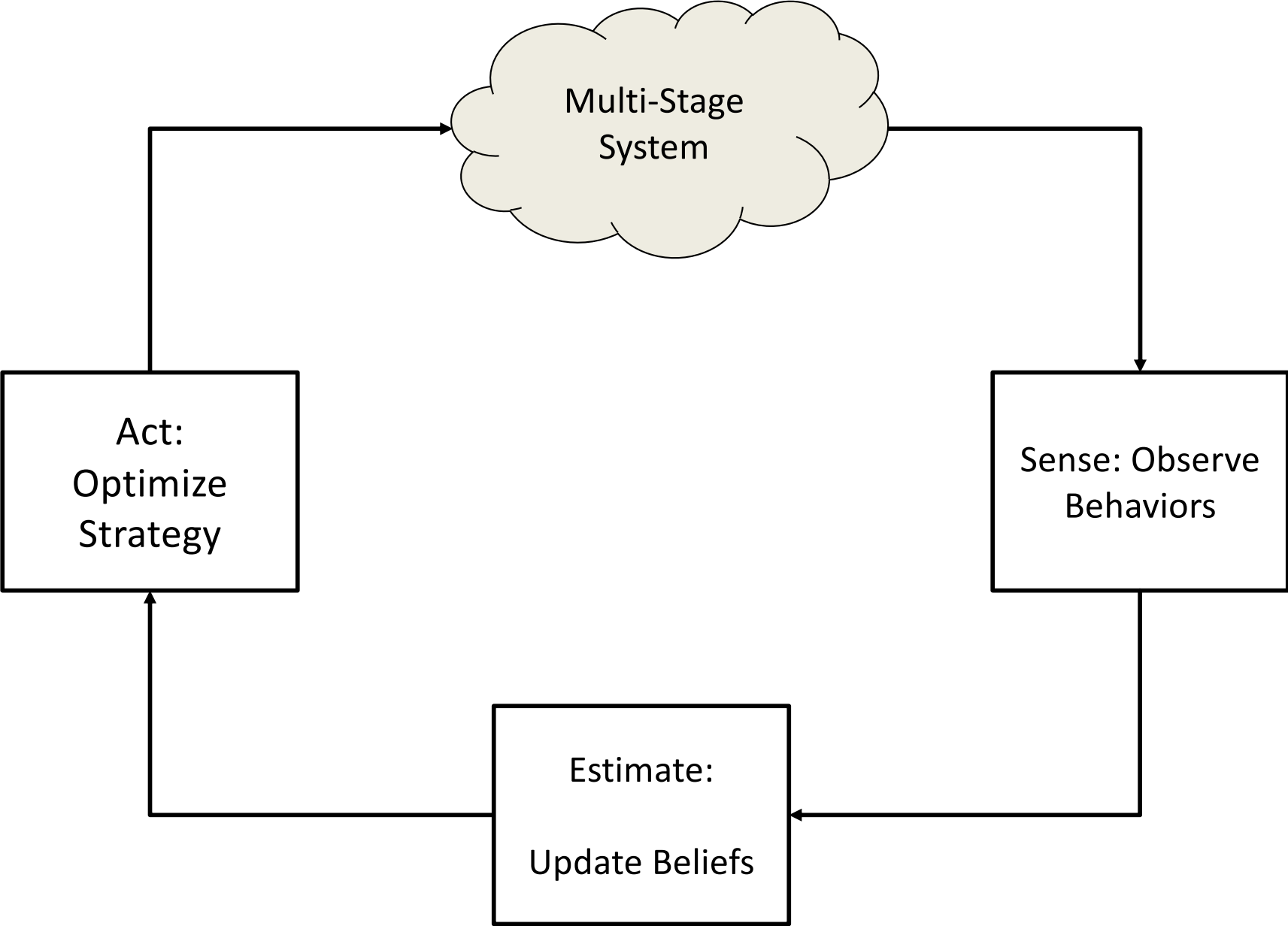
Observable History

$$\mathbf{h}^0 = \emptyset \quad \mathbf{h}^1 = \{a_1^0, a_2^0\} \quad \mathbf{h}^t = \{\mathbf{h}^{t-1}, a_1^{t-1}, a_2^{t-1}\} \quad \mathbf{h}^T = \{a_1^0, \dots, a_1^{T-1}, a_2^0, \dots, a_2^{T-1}\}$$

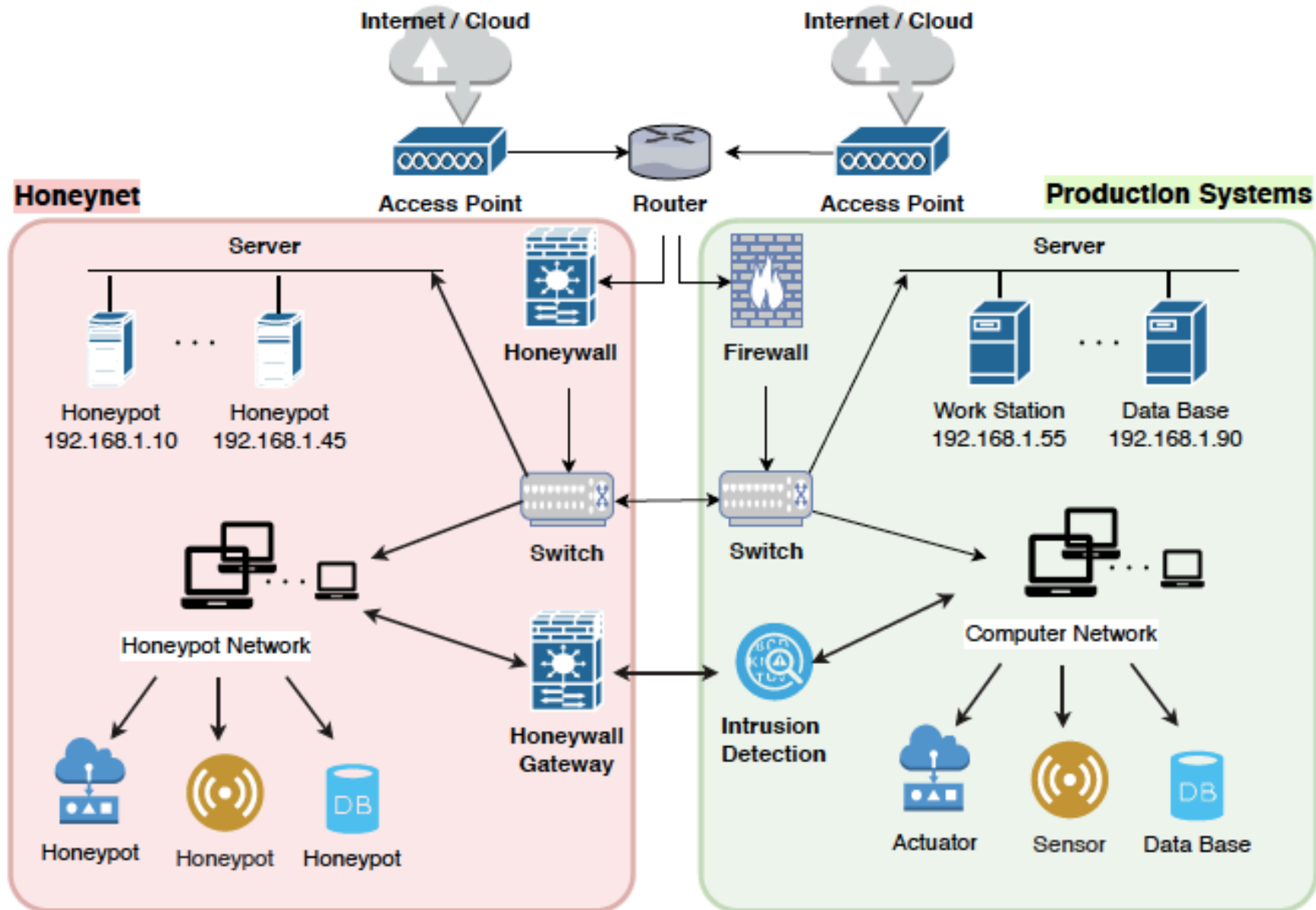
# Multi-Stage Co-Learning

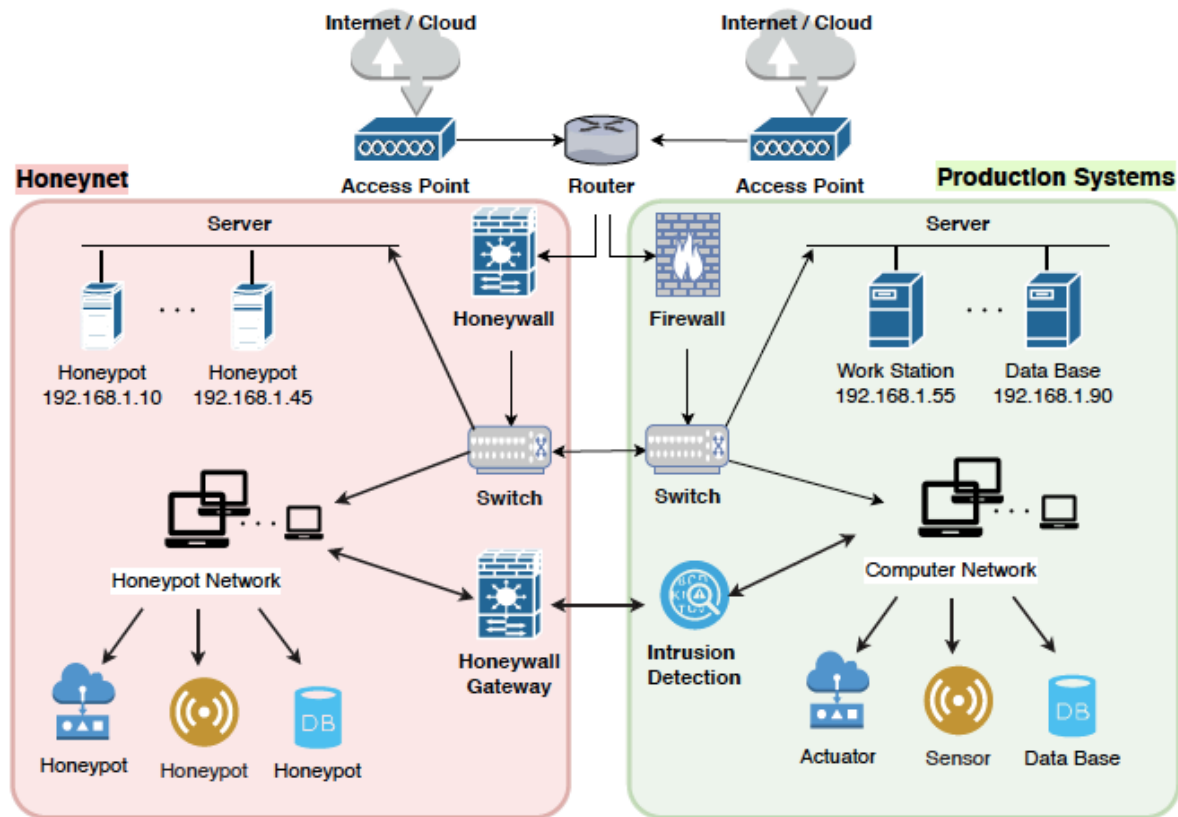






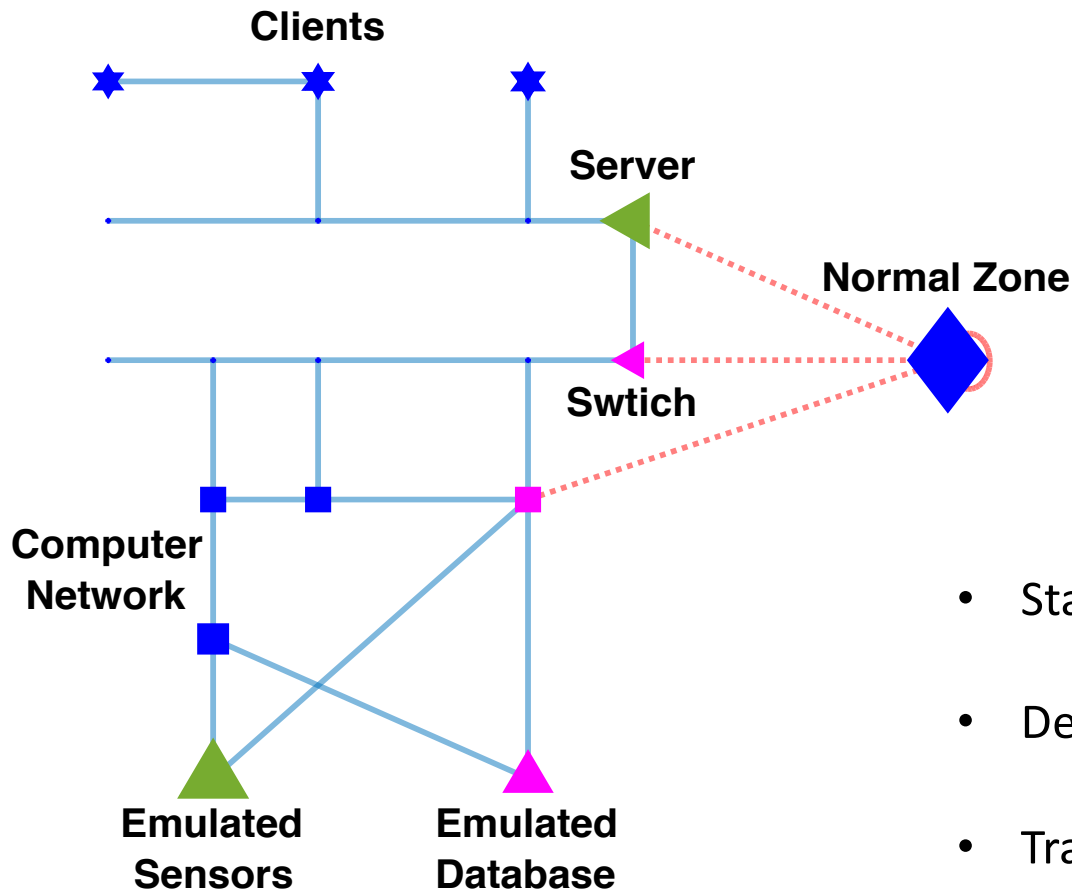
# Attacker Engagement in Honeypot





- Use a honeynet to emulate a production system.
- Interact with rather than directly eject attackers.
- Quickly attract attackers to target honeypots and engage them for a desired time.
- Grant attackers proper degree of freedom to avoid the escape risk and the identification risk.

# Abstraction: Semi-Markov Decision Process



Red: eject the attacker

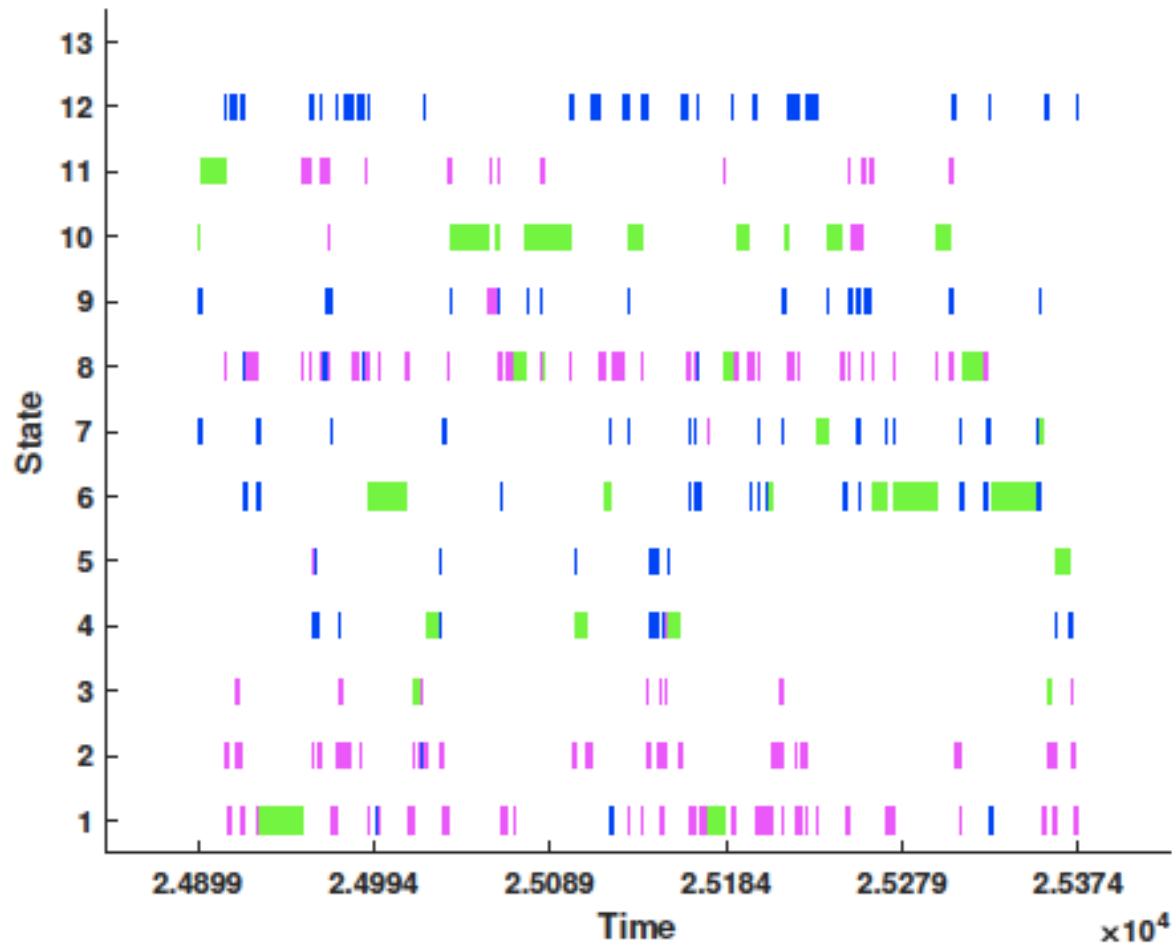
Blue: pure recording of attacker's activity

Purple: low-interaction with the attacker

Green: high-interaction with the attacker

- State at time  $t : s^t \in S$
- Defender's action at state  $s : a \in A(s)$
- Transition probability:  $p(s^{t+1} | s^t, a)$
- Sojourn time distribution  $q(t | s^t, s^{t+1}, a)$

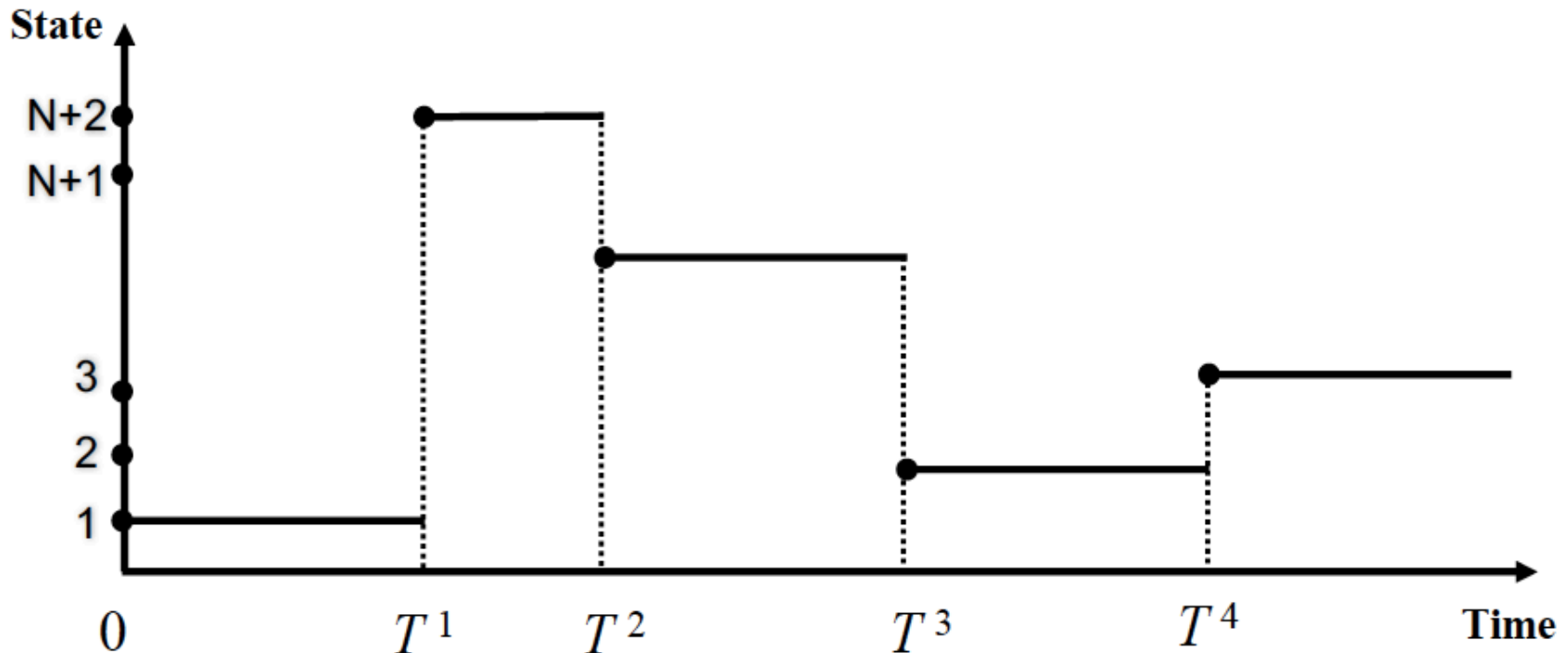
# Attacker's Footprint



- Treat transition kernel and sojourn distribution as threat intelligence.
- Characterize the escape risk and the identification risk.

# Semi-Markov Decision Process and Learning

- $T^k$  is the time of the  $k^{\text{th}}$  transition, which is a random variable.
- Defender receives a reward  $r$  at time  $\tau$ ,  $T^k \leq \tau \leq T^{k+1}$  if the next state is  $s'$  and the duration time at current  $s$ ,  $a$  is  $T^{k+1} - T^k$ .



# Q-Learning for SMDP

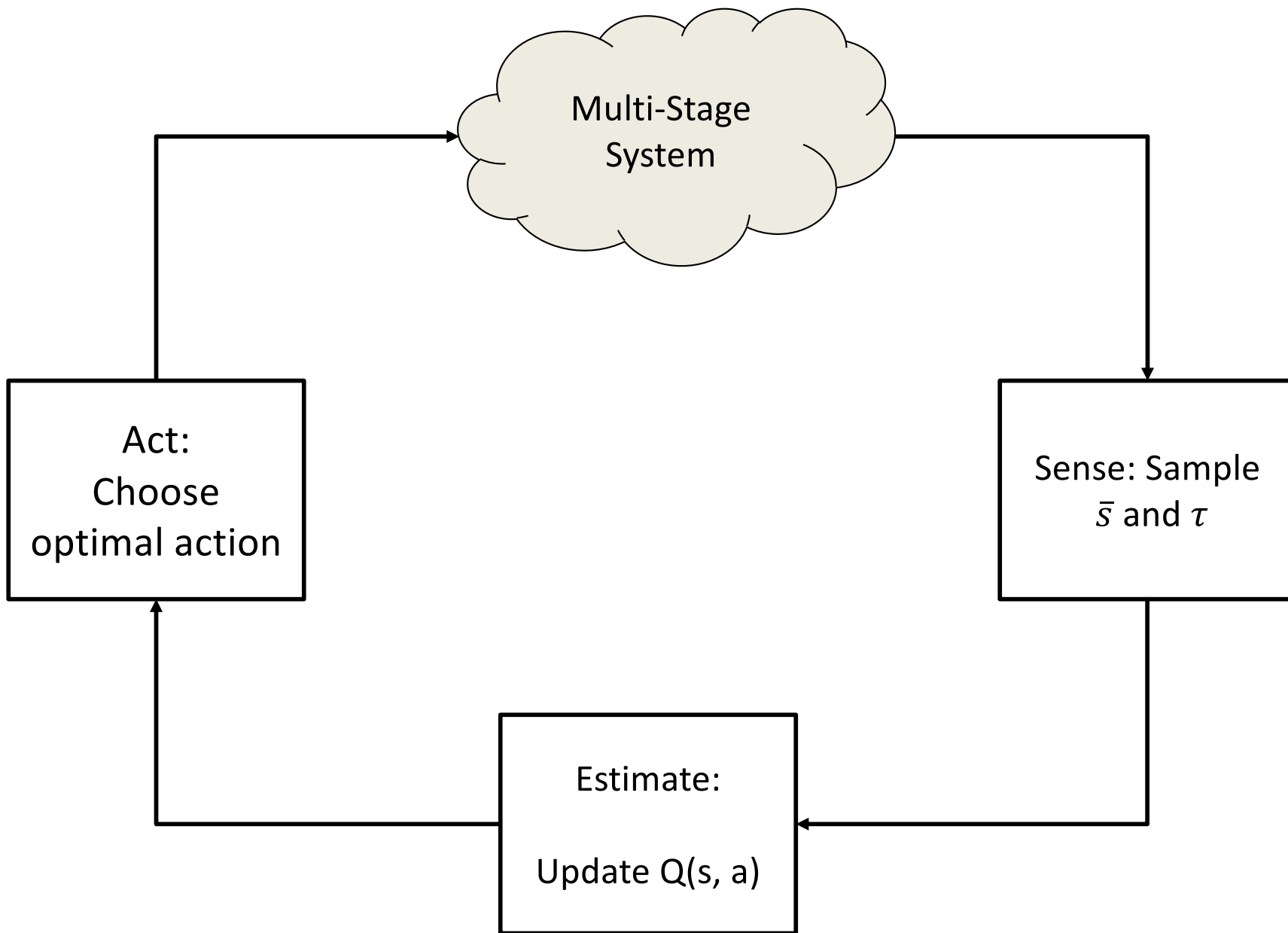
- Defender's stationary policy:  $\pi: S \rightarrow A(s)$

- DP representation

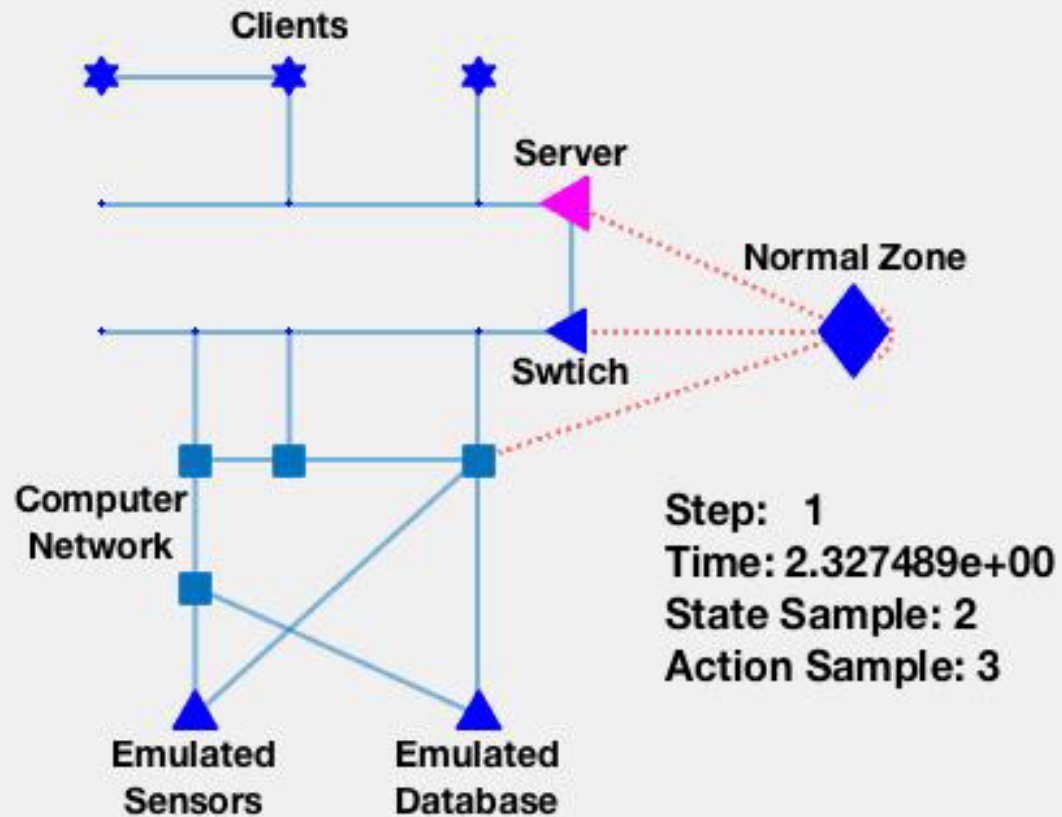
$$v(s) = \max_a \sum_{s'} p(s'|s, a) (r^\gamma(s, a, s') + \tilde{q}(s, a, s')v(s'))$$

- Q-Learning: sampled state  $\bar{s}$  and the duration time  $\tau$

$$Q^{k+1}(s, a) = (1 - \alpha^k)Q^k(s, a) + \alpha^k (r_1(s, a, \bar{s}) + \frac{(1 - e^{-\gamma\tau})r_2(s, a, \bar{s})}{\gamma} + e^{-\gamma\tau} \max_{a'} Q^k(\bar{s}, a'))$$





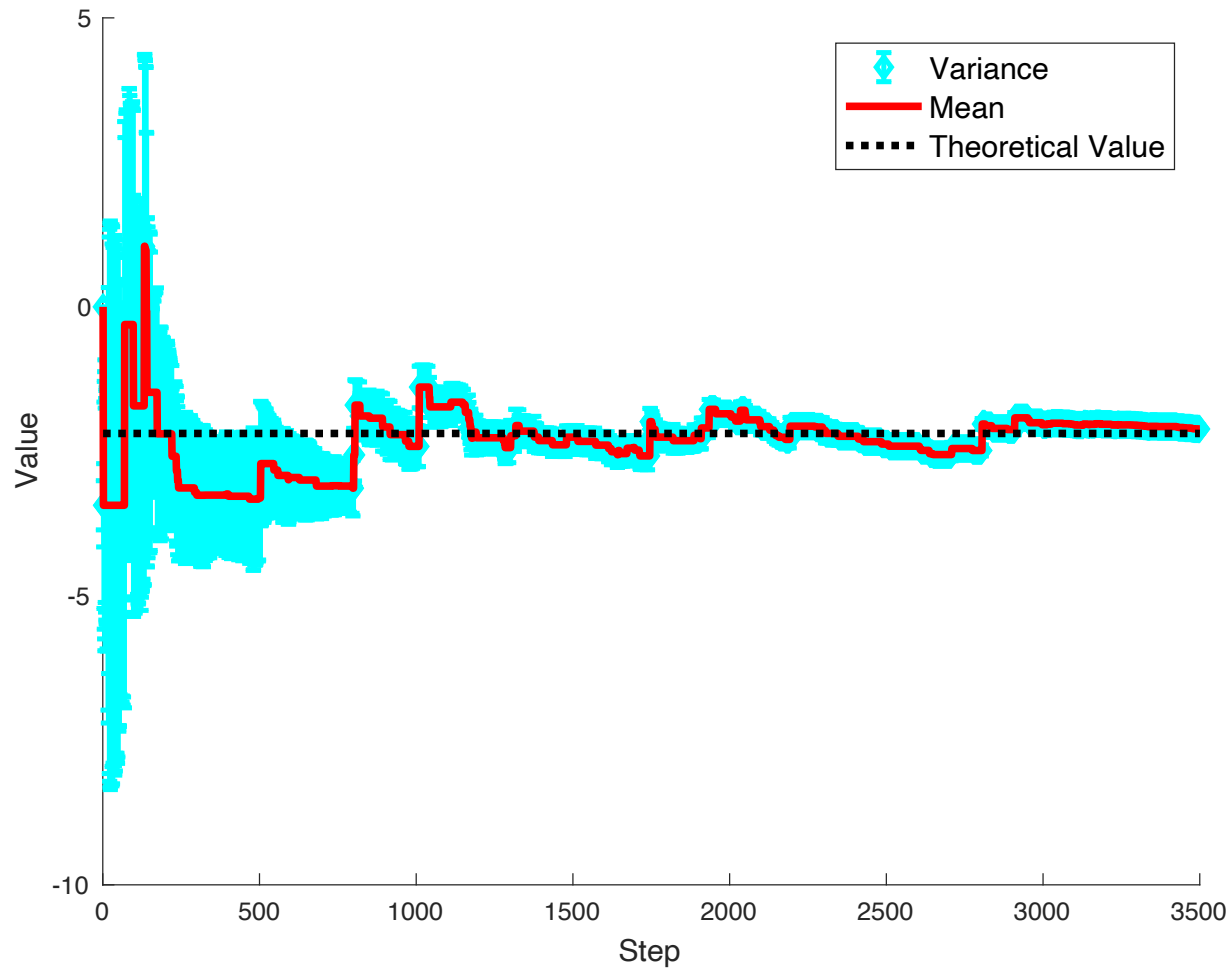


Red: eject the attacker

Purple: low-interaction with the attacker

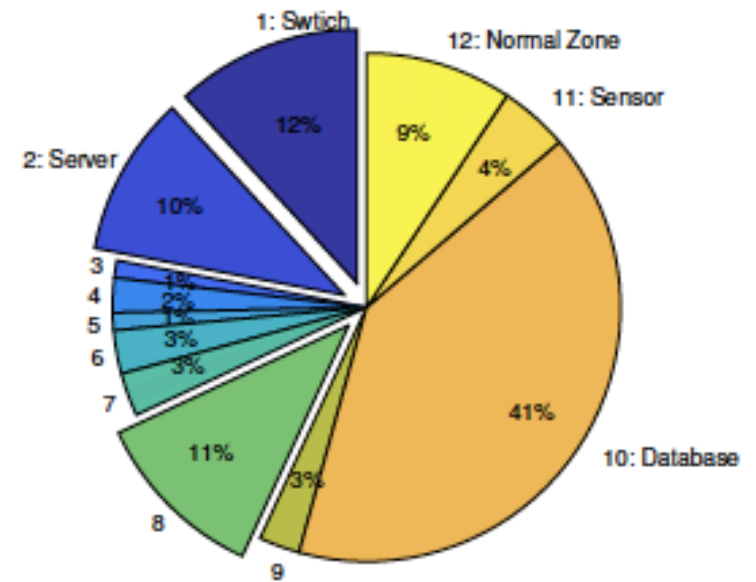
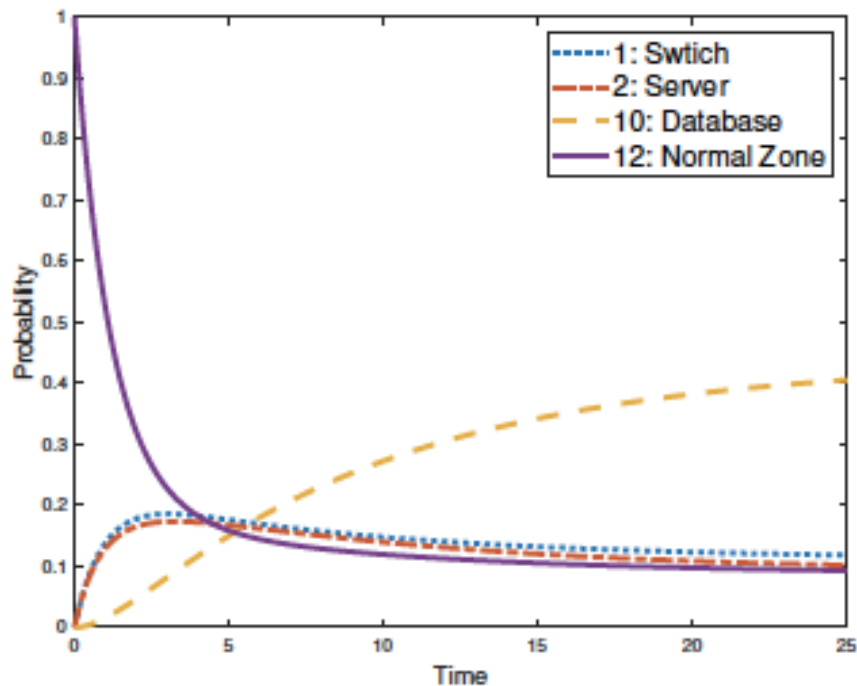
Blue: pure recording of attacker's activity

Green: high-interaction with the attacker



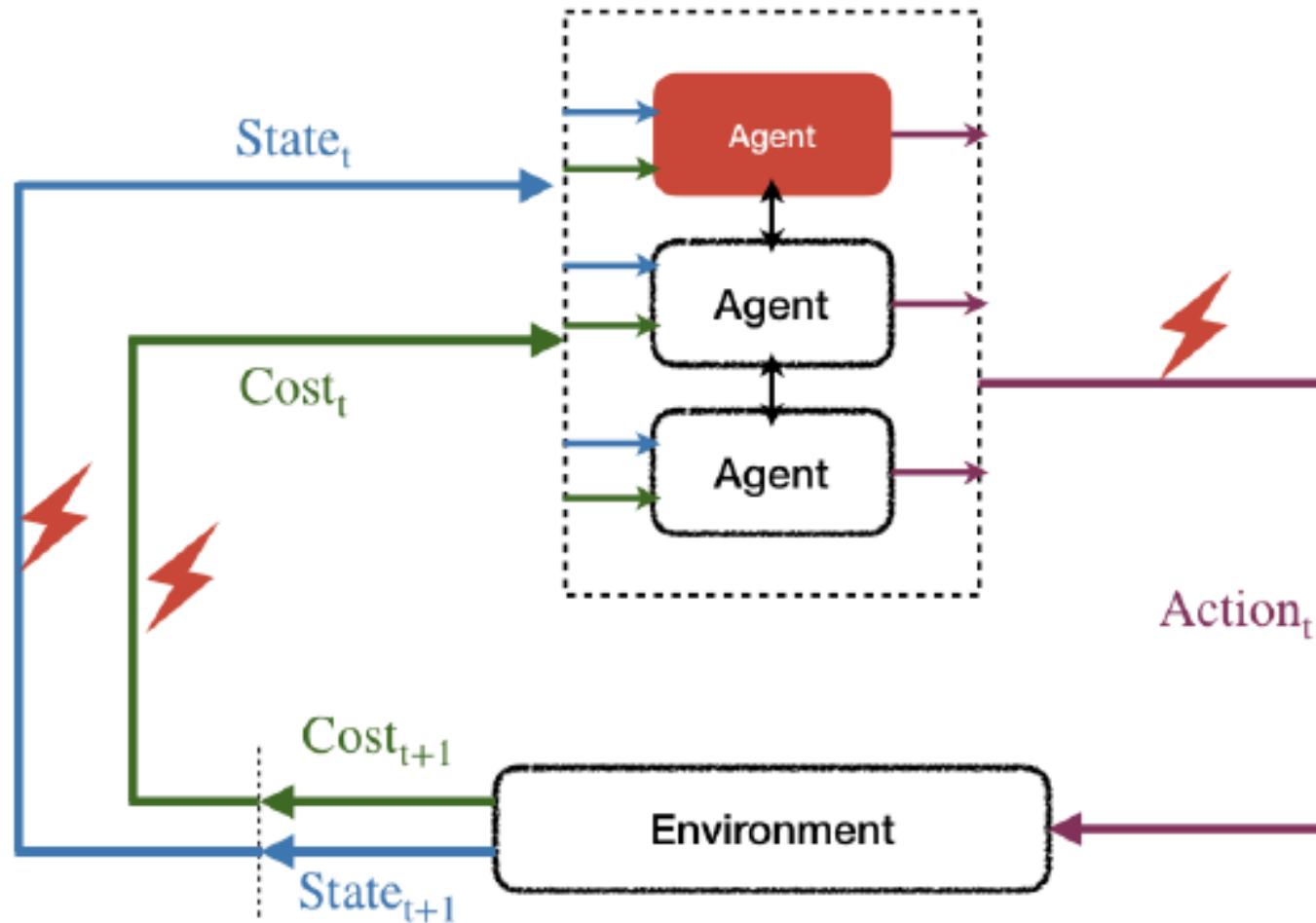
L. Huang and Q. Zhu, "Strategic Learning for Active, Adaptive, and Autonomous Cyber Defense," Adaptive Autonomous Secure Cyber Systems, Jajodia, S., Cybenko, G., Subrahmanian, V.S., Swarup, V., Wang, C., Wellman, M. (Eds.), 2020.

# Security Metrics to Evaluate Engagement



L. Huang, Q. Zhu, "Adaptive Honeypot Engagement through Reinforcement Learning of Semi-Markov Decision Processes," Conference on Decision and Game Theory for Security (GameSec), Oct. 30 - Nov. 1, 2019, Stockholm, Sweden.

# RL in Adversarial Environment



Y. Huang, Q. Zhu, "Deceptive Reinforcement Learning Under Adversarial Manipulations on Cost Signals," Conference on Decision and Game Theory for Security (GameSec), Oct. 30 - Nov. 1, 2019, Stockholm, Sweden.

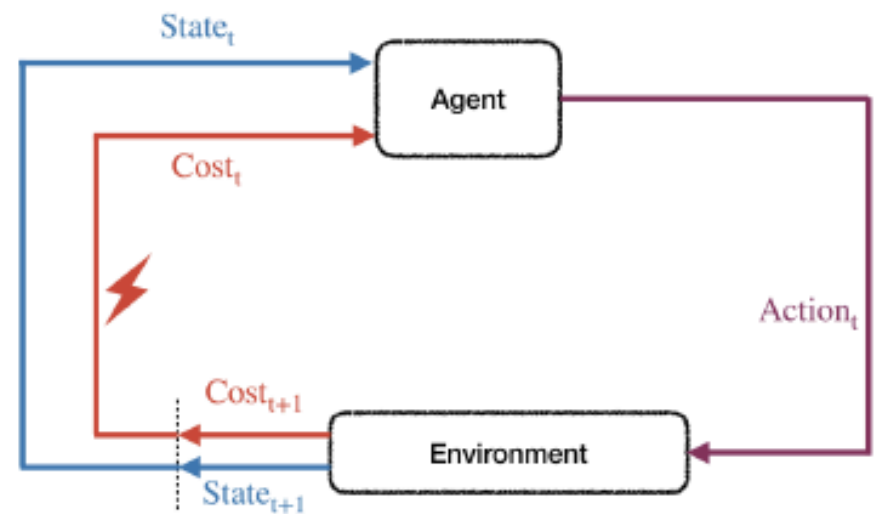
- Drones with RL techniques can be indirectly trained and weaponized by terrorists.
- Self-driving car can be misled to collisions by receiving false feedback.



# RL with Manipulated Cost Signals

Consider a Markov Decision Process (MDP) with **manipulated cost signals**, denoted by  $\langle \mathcal{S}, \mathcal{A}, c, \tilde{c}, \mathcal{P}, \beta \rangle$ .

- State space  $\mathcal{S} := \{1, 2, \dots, S\}$ .
- Action space  $\mathcal{A} := \{a_1, \dots, a_A\}$ .
- Cost function  $c : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ .
- Manipulated cost  $\tilde{c}$ .
- Transition probability kernel  $\mathcal{P}$
- Discounted factor  $\beta$



- The agent aims for **the optimal policy**  $w : \mathcal{S} \rightarrow \mathcal{A}$  minimizing

$$J(i, \mathbf{Z}) := \mathbf{E}\left[\sum_{t=0}^{\infty} \beta^t c(\Phi(t), Z(t)) \mid \Phi(0) = i\right],$$

where  $\Phi(\tau)$  is the state at time  $\tau$ ,  $Z(\tau)$  is the action chosen at time  $\tau$  by  $Z(\tau) = w(\Phi(\tau))$ .

- To find the optimal policy  $w^*$ , the agent implements **Q-learning algorithm**, i.e., for  $i = \Phi(t)$ ,  $a = Z(t)$ ,

$$Q_{t+1}(i, a) = Q_t(i, a) + s(t) \times \left[ \beta \min_b Q_t(\Psi_{t+1}(i, a), b) + c(i, a) - Q_t(i, a) \right],$$

where  $\Psi_{n+1}(i, a)$  is a realized  $\mathcal{S}$ -valued random variable with law  $p(i, \cdot, a)$ .

- Under proper conditions [V.Borkar; SIAM JCO; 2000],  $Q_t \rightarrow Q^*$  a.s., as  $t \rightarrow \infty$ . And  $Q^*$  satisfies

$$Q^*(i, a) = c(i, a) + \beta \sum_j p(i, j, a) \min_b Q^*(j, b), \quad i \in \mathcal{S}, a \in \mathcal{A}. \quad (1)$$

- We use  $Q^* = F(Q^*)$  to capture relation (1), where  $F : \mathbb{R}^{\mathcal{S} \times \mathcal{A}} \rightarrow \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ .
- The agent can find an optimal policy  $w^*$  by

$$w^*(i) = \arg \min_{a \in \mathcal{A}} Q^*(i, a), \quad i \in \mathcal{S}.$$



The Q-learning algorithm with manipulated cost signals:

$$Q_{t+1}(i, a) = Q_t(i, a) + s(t) \times \left[ \beta \min_b Q_t(\Psi_{t+1}(i, a), b) + \tilde{c}_t - Q_t(i, a) \right].$$

**Question:** How should the attacker manipulate cost signals?

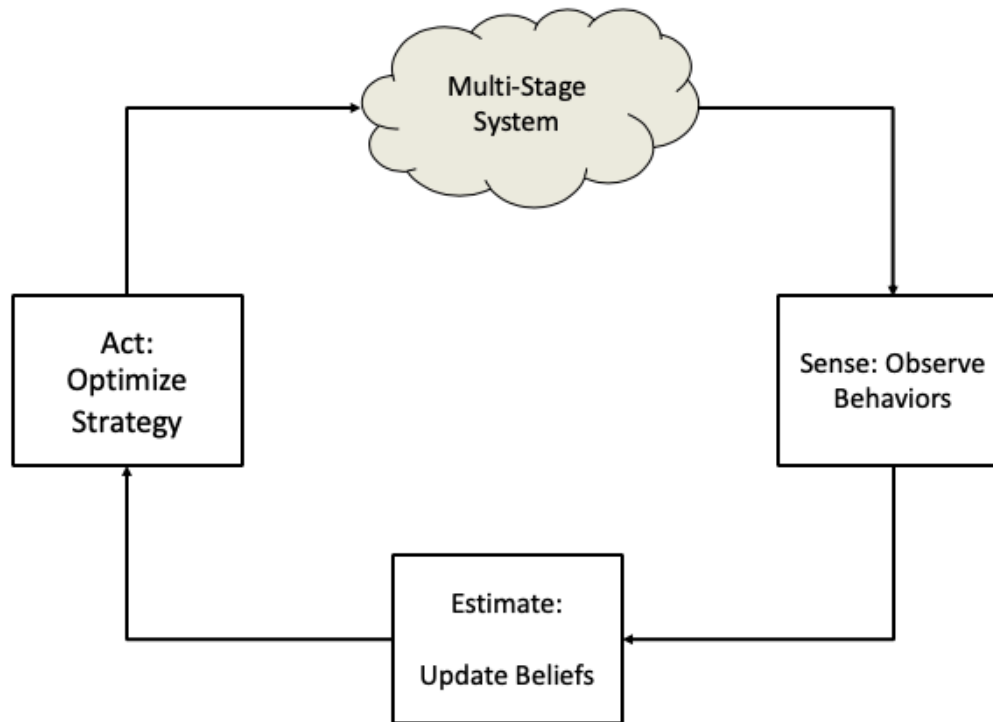
### Theorem (Manipulation Rule)

Let  $\tilde{Q}^* \in \mathbb{R}^{S \times A}$  be the Q-values learned under the falsified cost  $\tilde{c} \in \mathbb{R}^{S \times A}$ . Then  $\tilde{Q}^* \in \mathcal{V}_{w^\dagger}$  if and only if the falsified cost signals  $\tilde{c}$  designed by the adversary satisfy the following conditions

$$\tilde{c}(i, a) > (\mathbf{1}_i - \beta P_{ia})^T (I - \beta P_{w^\dagger})^{-1} \tilde{c}_{w^\dagger}. \quad (2)$$

for all  $i \in \mathcal{S}$ ,  $a \in \mathcal{A} \setminus \{w^\dagger(i)\}$ .

# Conclusion: Feedbacks and learning are fundamental for adaptive defense.



- *Moving target defense*: Baseline learning algorithms
- *Proactive defense against APT*: learning of dynamic games with incomplete information
- *Attacker engagement problem*: reinforcement learning
- *Security of reinforcement learning*: deceptive RL

# References

- Q. Zhu, H. Tembine and T. Basar, “Hybrid learning in stochastic games and its application in network security,” In F. L. Lewis and D. Liu (Eds.), Reinforcement Learning and Approximate Dynamic Programming for Feedback Control, IEEE Press Computational Intelligence Series, 2012. DOI: 10.1002/9781118453988.ch14
- K Horak, Q. Zhu and B. Bosansky, “Manipulating Adversary's Belief: A Dynamic Game Approach to Deception by Design for Proactive Network Security, " 8th Conference on Decision and Game Theory for Security (GameSec), Oct. 23-25, 2017, Vienna, Austria.
- L. Huang and Q. Zhu, “Strategic Learning for Active, Adaptive, and Autonomous Cyber Defense, "Adaptive Autonomous Secure Cyber Systems, Jajodia, S., Cybenko, G., Subrahmanian, V.S., Swarup, V., Wang, C., Wellman, M. (Eds.), 2020.