# Strategic Learning for Active, Adaptive, and Autonomous Cyber Defense

Linan Huang and Quanyan Zhu

**Abstract** The increasing instances of advanced attacks call for a new defense paradigm that is active, autonomous, and adaptive, named as the '3A' defense paradigm. This chapter introduces three defense schemes that actively interact with attackers to increase the attack cost and gather threat information, i.e., defensive deception for detection and counter-deception, feedback-driven Moving Target Defense (MTD), and adaptive honeypot engagement. Due to the cyber deception, external noise, and the absent knowledge of the other players' behaviors and goals, these schemes possess three progressive levels of information restrictions, i.e., from the parameter uncertainty, the payoff uncertainty, to the environmental uncertainty. To estimate the unknown and reduce the uncertainty, we adopt three different strategic learning schemes that fit the associated information restrictions. All three learning schemes share the same feedback structure of sensation, estimation, and actions so that the most rewarding policies get reinforced and converge to the optimal ones in autonomous and adaptive fashions. This work aims to shed lights on proactive defense strategies, lay a solid foundation for strategic learning under incomplete information, and quantify the tradeoff between the security and costs.

## 1 Introduction

Recent instances of `WannaCry` ransomware, `Petya` cyberattack, and `Stuxnet` malware have demonstrated the trends of modern attacks and the corresponding new security challenges as follows.

———————————

Linan Huang

Department of Electrical and Computer Engineering, New York University, 2 MetroTech Center, Brooklyn, NY, 11201, USA, e-mail: lh2328@nyu.edu

Quanyan Zhu

Department of Electrical and Computer Engineering, New York University, 2 MetroTech Center, Brooklyn, NY, 11201, USA, e-mail: qz494@nyu.edu

- **Advanced**: Attackers leverage sophisticated attack tools to invalidate the off-the-shelf defense schemes such as the firewall and intrusion detection systems.
- **Targeted**: Unlike automated probes, targeted attacks conduct thorough research in advance to expose the system architecture, valuable assets, and defense schemes.
- **Persistent**: Attackers can restrain the adversary's behaviors and bide their times to launch critical attacks. They are persistent in achieving the goal.
- **Adaptive**: Attackers can learn the defense strategies and unpatched vulnerabilities during the interaction with the defender and tailor their strategies accordingly.
- **Stealthy and Deceptive**: Attackers conceal their true intentions and disguise their claws to evade detection. The adversarial cyber deception endows attackers an information advantage over the defender.

Thus, defenders are urged to adopt active, adaptive, and autonomous defense paradigms to deal with the above challenges and proactively protect the system prior to the attack damages rather than passively compensate for the loss. In analogy to the classical **Kerckhoffs's principle** in the 19th century that attackers know the system, we suggest a new security principle for modern cyber systems as follows:
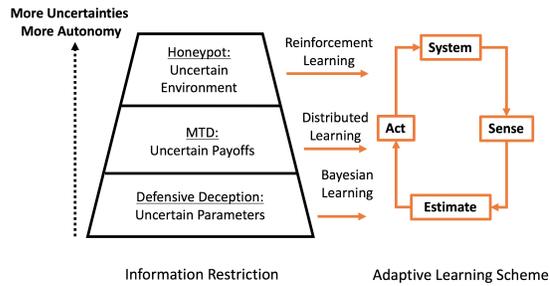
**Principle of 3A Defense**: A cyber defense paradigm is considered to be insufficiently secure if its effectiveness relies on

- Rule-abiding human behaviors.
- A perfect protection against vulnerabilities and a perfect prevention from system penetration.
- A perfect knowledge of attacks.

Firstly, 30% of data breaches are caused by privilege misuse and error by insiders according to Verizon's data breach report in 2019 [1]. Security administration does not work well without the support of technology, and autonomous defense strategies are required to deal with the increasing volume of sophisticated attacks. Secondly, systems always have undiscovered vulnerabilities or unpatched vulnerabilities due to the long supply chain of uncontrollable equipment providers [2] and the increasing complexities in the system structure and functionality. Thus, an effective paradigm should assume a successful infiltration and pursue strategic securities through interacting with intelligent attackers. Finally, due to adversarial deception techniques and external noises, the defender cannot expect a perfect attack model with predicable behaviors. The defense mechanism should be robust under incomplete information and adaptive to the evolution of attacks.

In this chapter, we illustrate three active defense schemes in our previous works, which are designed based on the new cyber security principle. They are defensive deception for detection and counter-deception [3, 4, 5] in Section 2, feedback-driven Moving Target Defense (MTD) [6] in Section 3, and adaptive honeypot engagement [7] in Section 4. All three schemes is of incomplete information, and we arrange them based on three progressive levels of information restrictions as shown in the left part of Fig. 1.

**Fig. 1** The left part of the figure describes the degree of information restriction. From bottom to top, the defense scheme becomes more autonomous and relies less on an exact attack model, which also result in more uncertainties. The right part is the associated feedback learning schemes.



The first scheme in Section 2 considers the obfuscation of characteristics of known attacks and systems through a random parameter called the player's *type*. The only uncertainty origins from the player's *type*, and the mapping from the type to the utility is known deterministically. The MTD scheme in Section 3 considers unknown attacks and systems whose utilities are completely uncertain, while the honeypot engagement scheme in Section 4 further investigates environmental uncertainties such as the transition probability, the sojourn distribution, and the investigation reward.

To deal with these uncertainties caused by different information structures, we suggest three associated learning schemes as shown in the right part of Fig. 1, i.e., Bayesian learning for the parameter estimation, distributed learning for the utility acquisition without information sharing, and reinforcement learning for the optimal policy obtainment under the unknown environment. All three learning methods form a feedback loop that strategically incorporates the samples generated during the interaction between attackers and defenders to persistently update the beliefs of known and then take actions according to current optimal decision strategies. The feedback structure makes the learning adaptive to behavioral and environmental changes.

Another common point of these three schemes is the quantification of the tradeoff between security and the different types of cost. In particular, the costs result from the attacker's identification of the defensive deception, the system usability, and the risk of attackers penetrating production systems from the honeynet, respectively.

## 1.1 Literature

The idea of using deceptions defensively to detect and deter attacks has been studied theoretically as listed in the taxonomic survey [8], implemented to the Adversarial Tactics, Techniques and Common Knowledge (ATT&CK$^{TM}$) adversary model system [9], and tested in the real-time cyber-wargame experiment [10]. Many previous works imply the similar idea of type obfuscation, e.g., creating social network avatars (fake personas) on the major social networks [11], implementing honey files

for ransomware actions [12], and disguising a production system as a honeypot to scare attackers away [13].

Moving target defense (MTD) allows dynamic security strategies to limit the exposure of vulnerabilities and the effectiveness of the attacker's reconnaissance by increasing complexities and costs of attacks [14]. To achieve an effective MTD, [15] proposes the instruction set and the address space layout randomization, [16] studies the deceptive routing against jamming in multi-hop relay networks, and [17] uses the Markov chain to model the MTD process and discusses the optimal strategy to balance the defensive benefit and the network service quality.

The previous two methods use the defensive deception to protect the system and assets. To further gather threat information, the defender can implement honeypots to lure attackers to conduct adversarial behaviors and reveal their TTPs in a controlled and monitored environment. Previous works [18, 19] have investigated the adaptive honeypot deployment to effectively engage attackers without their notices. The authors in recent work [20] proposes a continuous-state Markov Decision Process (MDP) model and focuses on the optimal timing of the attacker ejection.

Game-theoretic models are natural frameworks to capture the multistage interaction between attackers and defenders. Recently, game theory has been applied to different sets of security problems, e.g., Stackelberg and signaling games for deception and proactive defenses [21, 6, 22, 23, 24, 16, 25, 26, 27], network games for cyber-physical security [28, 29, 30, 31, 32, 33, 34, 35, 36, 37], dynamic games for adaptive defense [38, 39, 40, 41, 3, 42, 43, 44, 45, 46], and mechanism design theory for security [47, 48, 49, 50, 51, 52, 53, 54, 55].

Information asymmetry among the players in network security is a challenge to deal with. The information asymmetry can be either leveraged or created by the attacker or the defender for achieving a successful cyber deception. For example, techniques such as honeynets [56, 22], moving target defense [6, 14, 3], obfuscation [57, 58, 59, 60], and mix networks [61] have been introduced to create difficulties for attackers to map out the system information.

To overcome the created or inherent uncertainties of networks, many works have studied the strategic learning in security games, e.g., Bayesian learning for unknown adversarial strategies [62], heterogeneous and hybrid distributed learning [46, 63], multiagent reinforcement learning for intrusion detection [64]. Moreover, these learning schemes are combined to achieve better properties, e.g., distributed Bayesian learning [65], Bayesian reinforcement learning [66], and distributed reinforcement learning [67].

## 1.2 Notation

Throughout the chapter, we use calligraphic letter $\mathcal{A}$ to define a set and $|\mathcal{A}|$ as the cardinality of the set. Let $\triangle\mathcal{A}$ represent the set of probability distributions over $\mathcal{A}$. If set $\mathcal{A}$ is discrete and finite, $\triangle\mathcal{A} := \{p : \mathcal{A} \mapsto R_+ | \sum_{a \in \mathcal{A}} p(a) = 1\}$, otherwise, $\triangle\mathcal{A} := \{p : \mathcal{A} \mapsto R_+ | \int_{a \in \mathcal{A}} p(a) = 1\}$. Row player $P_1$ is the defender (pronoun

'she') and $P_2$ (pronoun 'he') is the user (or the attacker) who controls the column of the game matrix. Both players want to maximize their own utilities. The indicator function $\mathbf{1}_{\{x=y\}}$ equals one if $x = y$, and zero if $x \neq y$.

## 1.3 Organization of the Chapter

The rest of the paper is organized as follows. In Section 2, we elaborate defensive deception as a countermeasure of the adversarial deception under a multistage setting where Bayesian learning is applied for the parameter uncertainty. Section 3 introduces a multistage MTD framework and the uncertainties of payoffs result in distributed learning schemes. Section 4 further considers reinforcement learning for environmental uncertainties under the honeypot engagement scenario. The conclusion and discussion are presented in Section 5.
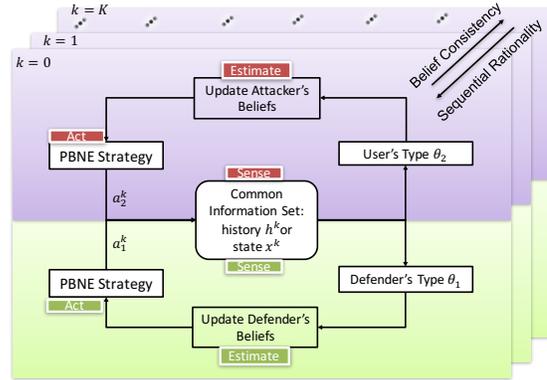
## 2 Bayesian Learning for Uncertain Parameters

Under the mild restrictive information structure, each player' utility is completely governed by a finite group of parameters which form his/her *type*. Each player's *type* characterizes all the uncertainties about this player during the game interaction, e.g., the physical outcome, the payoff function, and the strategy feasibility, as an equivalent utility uncertainty without loss of generality [68]. Thus, the revelation of the type value directly results in a game of complete information. In the cyber security scenario, a discrete type can distinguish either systems with different kinds of vulnerabilities or attackers with different targets. The type can also be a continuous random variable representing either the threat level or the security awareness level [3, 4]. Since each player $P_i$ takes actions to maximize his/her own type-dependent utility, the other player $P_j$ can form a belief to estimate $P_i$'s type based on the observation of $P_i$'s action history. The utility optimization under the beliefs results in the Perfect Bayesian Nash Equilibrium (PBNE) which generates new action samples and updates the belief via the Bayesian rule. We plot the feedback Bayesian learning process in Fig. 2 and elaborate each element in the following subsections based on our previous work [5].
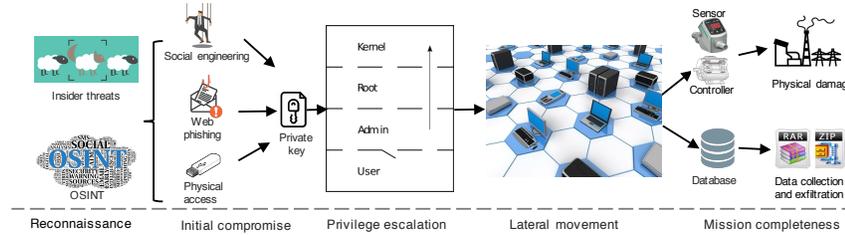
## 2.1 Type and Multistage Transition

Through adversarial deception techniques, attackers can disguise their subversive actions as legitimate behaviors so that the defender $P_1$ cannot judge whether a user $P_2$'s type $\theta_2 \in \Theta_2 := \{\theta_2^g, \theta_2^b\}$ is legitimate $\theta_2^g$ or adversarial $\theta_2^b$. As a countermeasure, the defender can introduce the defensive deception so that the attacker cannot

**Fig. 2** The feedback loop of the Bayesian learning from the initial stage $k = 1$ to the terminal stage $k = K$. Each player forms a belief of the other player's type and persistently updates the belief based on the actions resulted from the PBNE strategies which are the results of Bayesian learning.



distinguish between a primitive system $\theta_1^L$ and a sophisticated system $\theta_1^H$, i.e., the defender has a binary type $\theta_1 \in \Theta_1 := \{\theta_1^H, \theta_1^L\}$. A sophisticated system is costly yet deters attacks and causes damages to attackers. Thus, a primitive system can disguise as a sophisticated one to draw the same threat level to attackers yet avoid the implementation cost of sophisticated defense techniques.

Many cyber networks contain hierarchical layers, and up-to-date attackers such as Advanced Persistent Threats (APTs) aim to penetrate these layers and reach specific targets at the final stage as shown in Fig. 3.



Fig. 3: The multistage structure of APT kill chain is composed of reconnaissance, initial compromise, privilege escalation, lateral movement, and mission execution.

At stage $k \in \{0, 1, \cdots, K\}$, $P_i$ takes an action $a_i^k \in \mathcal{A}_i^k$ from a finite and discrete set $\mathcal{A}_i^k$. Both players' actions become fully observable after applied and each action does not directly reveal the private type. For example, both legitimate and adversarial users can choose to access the sensor, and both primitive and sophisticated defenders can choose to monitor the sensor. Both players' actions up to stage $k$ constitute the *history* $h^k = \{a_1^0, \cdots, a_1^{k-1}, a_2^0, \cdots, a_2^{k-1}\} \in \mathcal{H}^k := \prod_{i=1}^2 \prod_{\bar{k}=0}^{k-1} \mathcal{A}_i^{\bar{k}}$. Given history $h^k$ at the current stage $k$, players at stage $k + 1$ obtain an updated history $h^{k+1} = h^k \cup \{a_1^k, a_2^k\}$ after the observation $a_1^k, a_2^k$. A state $x^k \in \mathcal{X}^k$ at each stage $k$ is the smallest set of quantities that summarize information about actions in previous stages so that

the initial state $x^0 \in \mathcal{X}^0$ and the history at stage $k$ uniquely determine $x^k$ through a known state transition function $f^k$, i.e., $x^{k+1} = f^k(x^k, a_1^k, a_2^k)$, $\forall k \in \{0, 1, \cdots, K-1\}$. The state can represent the location of the user in the attack graph, and also other quantities such as users' privilege levels and status of sensor failures.

A behavioral strategy $\sigma_i^k \in \Sigma_i^k : \mathcal{I}_i^k \mapsto \triangle(\mathcal{A}_i^k)$ maps $P_i$'s information set $\mathcal{I}_i^k$ at stage $k$ to a probability distribution over the action space $\mathcal{A}_i^k$. At the initial stage 0, since the only information available is the player's type realization, the information set $\mathcal{I}_i^0 = \Theta_i$. The action is a realization of the behavioral strategy, or equivalently, a sample drawn from the probability distribution $\sigma_i^k(\cdot|I_i^k)$. With a slight abuse of notation, we denote $\sigma_i^k(a_i^k|I_i^k)$ as the probability of $P_i$ taking action $a_i^k \in \mathcal{A}_i^k$ given the available information $I_i^k \in \mathcal{I}_i^k$.

## 2.2 Bayesian Update under Two Information Structure

Since the other player's type is of private information, $P_i$ forms a belief $b_i^k : \mathcal{I}_i^k \mapsto \triangle(\Theta_j)$, $j \neq i$, on $P_j$'s type using the available information $\mathcal{I}_i^k$. Likewise, given information $I_i^k \in \mathcal{I}_i^k$ at stage $k$, $P_i$ believes with a probability $b_i^k(\theta_j|I_i^k)$ that $P_j$ is of type $\theta_j \in \Theta_j$. The initial belief $b_i^0 : \Theta_i \mapsto \triangle\Theta_j$, $\forall i, j \in \{1, 2\}$, $j \neq i$, is formed based on an imperfect detection, side-channel information or the statistic estimation resulted from past experiences.

If the system has a *perfect recall* $\mathcal{I}_i^k = \mathcal{H}^k \times \Theta_i$, then players can update their beliefs according to the Bayesian rule:

$$b_i^{k+1}(\theta_j|h^k \cup \{a_i^k, a_j^k\}, \theta_i) = \frac{\sigma_i^k(a_i^k|h^k, \theta_i)\sigma_j^k(a_j^k|h^k, \theta_j)b_i^k(\theta_j|h^k, \theta_i)}{\sum_{\bar{\theta}_j \in \Theta_j} \sigma_i^k(a_i^k|h^k, \theta_i)\sigma_j^k(a_j^k|h^k, \bar{\theta}_j)b_i^k(\bar{\theta}_j|h^k, \theta_i)}. \quad (1)$$

Here, $P_i$ updates the belief $b_i^k$ based on the observation of the action $a_i^k, a_j^k$. When the denominator is 0, the history $h^{k+1}$ is not reachable from $h^k$, and a Bayesian update does not apply. In this case, we let $b_i^{k+1}(\theta_j|h^k \cup \{a_i^k, a_j^k\}, \theta_i) := b_i^0(\theta_j|\theta_i)$.

If the information set is taken to be $\mathcal{I}_i^k = \mathcal{X}^k \times \Theta_i$ with the Markov property that $\Pr(x^{k+1}|\theta_j, x^k, \cdots, x^1, x^0, \theta_i) = \Pr(x^{k+1}|\theta_j, x^k, \theta_i)$, then the Bayesian update between two consequent states is

$$b_i^{k+1}(\theta_j|x^{k+1}, \theta_i) = \frac{\Pr(x^{k+1}|\theta_j, x^k, \theta_i)b_i^k(\theta_j|x^k, \theta_i)}{\sum_{\bar{\theta}_j \in \Theta_j} \Pr(x^{k+1}|\bar{\theta}_j, x^k, \theta_i)b_i^k(\bar{\theta}_j|x^k, \theta_i)}. \quad (2)$$

The Markov belief update (2) can be regarded as an approximation of (1) using action aggregations. Unlike the history set $\mathcal{H}^k$, the dimension of the state set $|\mathcal{X}^k|$ does not grow with the number of stages. Hence, the Markov approximation significantly reduces the memory and computational complexity.

## 2.3 Utility and PBNE

At each stage $k$, $P_i$'s stage utility $\bar{J}_i^k : \mathcal{X}^k \times \mathcal{A}_1^k \times \mathcal{A}_2^k \times \theta_1 \times \theta_2 \times \mathcal{R} \mapsto \mathcal{R}$ depends on both players' types and actions, the current state $x^k \in \mathcal{X}^k$, and an external noise $w_i^k \in \mathcal{R}$ with a known probability density function $\varpi_i^k$. The noise term models unknown or uncontrolled factors that can affect the value of the stage utility. Denote the expected stage utility as $J_i^k(x^k, a_1^k, a_2^k, \theta_1, \theta_2) :=$ $E_{w_i^k \sim \varpi_i^k} \bar{J}_i^k(x^k, a_1^k, a_2^k, \theta_1, \theta_2, w_i^k), \forall x^k, a_1^k, a_2^k, \theta_1, \theta_2$.

Given the type $\theta_i \in \Theta_i$, the initial state $x^{k_0} \in \mathcal{X}^{k_0}$, and both players' strategies $\sigma_i^{k_0:K} := [\sigma_i^k(a_i^k | x^k, \theta_i)]_{k=k_0, \cdots, K} \in \prod_{k=k_0}^{K} \Sigma_i^k$ from stage $k_0$ to $K$, we can determine the expected cumulative utility $U_i^{k_0:K}$ for $P_i, i \in \{1, 2\}$, by taking expectations over the mixed-strategy distributions and the $P_i$'s belief on $P_j$'s type, i.e.,

$$U_i^{k_0:K}(\sigma_i^{k_0:K}, \sigma_j^{k_0:K}, x^{k_0}, \theta_i) := \sum_{k=k_0}^{K} E_{\theta_j \sim b_i^k, a_i^k \sim \sigma_i^k, a_j^k \sim \sigma_j^k} J_i^k(x^k, a_1^k, a_2^k, \theta_1, \theta_2). \quad (3)$$

The attacker and the defender use the Bayesian update to reduce their uncertainties on the other player's type. Since their actions affect the belief update, both players at each stage should optimize their expected cumulative utilities concerning the updated beliefs, which leads to the solution concept of PBNE in Definition 1.

**Definition 1** Consider the two-person $K$-stage game with a double-sided incomplete information, a sequence of beliefs $b_i^k, \forall k \in \{0, \cdots, K\}$, an expected cumulative utility $U_i^{0:K}$ in (3), and a given scalar $\varepsilon \geq 0$. A sequence of strategies $\sigma_i^{*,0:K} \in \prod_{k=0}^{K} \Sigma_i^k$ is called $\varepsilon$-perfect Bayesian Nash equilibrium for player $i$ if the following two conditions are satisfied.

C1: Belief consistency: under the strategy pair $(\sigma_1^{*,0:K}, \sigma_2^{*,0:K})$, each player's belief $b_i^k$ at each stage $k = 0, \cdots, K$ satisfies (2).

C2: Sequential rationality: for all given initial state $x^{k_0} \in \mathcal{X}^{k_0}$ at every initial stage $k_0 \in \{0, \cdots, K\}, \forall \sigma_1^{k_0:K} \in \prod_{k=0}^{K} \Sigma_1^k, \forall \sigma_2^{k_0:K} \in \prod_{k=0}^{K} \Sigma_2^k$,

$$
\begin{aligned}
U_1^{k_0:K}(\sigma_1^{*,k_0:K}, \sigma_2^{*,k_0:K}, x^{k_0}, \theta_1) + \varepsilon &\geq U_1^{k:K}(\sigma_1^{k_0:K}, \sigma_2^{*,k_0:K}, x^{k_0}, \theta_1), \\
U_2^{k_0:K}(\sigma_1^{*,k_0:K}, \sigma_2^{*,k_0:K}, x^{k_0}, \theta_2) + \varepsilon &\geq U_2^{k:K}(\sigma_1^{*,k_0:K}, \sigma_2^{k_0:K}, x^{k_0}, \theta_2).
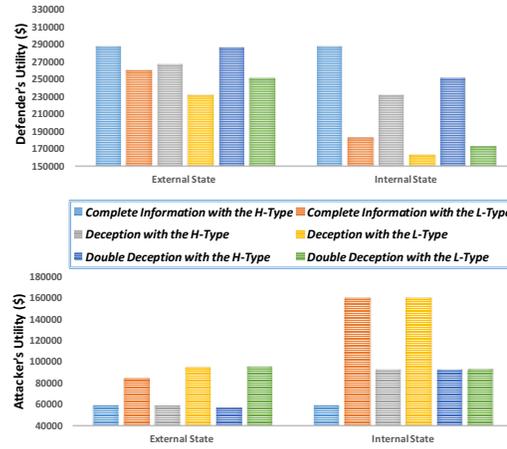\end{aligned}
\quad (4)
$$

When $\varepsilon = 0$, the equilibrium is called Perfect Bayesian Nash Equilibrium (**PBNE**).

Solving PBNE is challenging. If the type space is discrete and finite, then given each player's belief at all stages, we can solve the equilibrium strategy satisfying condition C2 via dynamic programming and a bilinear program. Next, we update the belief at each stage based on the computed equilibrium strategy. We iterate the above update on the equilibrium strategy and belief until they satisfy condition C1 as demonstrated in [5]. If the type space is continuous, then the Bayesian update can be simplified into a parametric update under the conjugate prior assumption. Next,

the parameter after each belief update can be assimilated into the backward dynamic programming of equilibrium strategy with an expanded state space [4]. Although no iterations are required, the infinite dimension of continuous type space limits the computation to two by two game matrices.

We apply the above framework and analysis to a case study of Tennessee Eastman (TE) process and investigate both players' multistage utilities under the adversarial and the defensive deception in Fig. 4. Some insights are listed as follows.

**Fig. 4** The cumulative utilities of the attacker and the defender under the complete information, the adversarial deception, and the defensive deception. The complete information refers to the scenario where both players know the other player's type. The deception with the H-type or the L-type means that the attacker knows the defender's type to be $\theta_1^H$ or $\theta_1^L$, respectively, yet the defender has no information about the user's type. The double-sided deception indicates that both players do not know the other player's type.



First, the defender's payoffs under type $\theta_1^H$ can increase as much as 56% than those under type $\theta_1^L$. Second, the defender and the attacker receive the highest and the lowest payoff, respectively, under the complete information. When the attacker introduces deceptions over his type, the attacker's utility increases and the system utility decreases. Third, when the defender adopts defensive deceptions to introduce double-sided incomplete information, we find that the decrease of system utilities is reduced by at most 64%, i.e., the decrease of system utilities changes from $55,570 to $35,570 under the internal state and type $\theta_1^H$. The double-sided incomplete information also brings lower utilities to the attacker than the one-sided adversarial deception. However, the system utility under the double-sided deception is still less than the complete information case, which concludes that acquiring complete information of the adversarial user is the most effective defense. However, if the complete information cannot be obtained, the defender can mitigate her loss by introducing defensive deceptions.

# 3 Distributed Learning for Uncertain Payoffs

In the previous section, we study known attacks and systems that adopt cyber deception to conceal their types. We assume common knowledge of the prior probability distribution of the unknown type, and also a common observation of either the action history or the state at each stage. Thus, each player can use Bayesian learning to reduce the other player's type uncertainty.

In this section, we consider unknown attacks in the MTD game stated in [6] where each player has no information on the past actions of the other player, and the payoff functions are subject to noises and disturbances with unknown statistical characteristics. Without information sharing between players, the learning is distributed.

## 3.1 Static Game Model of MTD

We consider a system of $N$ layers yet focus on the static game at layer $l \in \{1, 2, \cdots, N\}$ because the technique can be employed at each layer of the system independently. At layer $l$, $\mathcal{V}_l := \{v_{l,1}, v_{l,2}, \cdots, v_{l,n_l}\}$ is the set of $n_l$ system vulnerabilities that an attacker can exploit to compromise the system. Instead of a static configuration at layer $l$, the defender can choose to change her configuration from a finite set of $m_l$ feasible configurations $C_l := \{c_{l,1}, c_{l,2}, \cdots, c_{l,m_l}\}$. Different configurations result in different subsets of vulnerabilities among $\mathcal{V}_l$, which are characterized by the vulnerability map $\pi_l : C_l \rightarrow 2^{\mathcal{V}_l}$. We call $\pi_l(c_{l,j})$ the attack surface at stage $l$ under configuration $c_{l,j}$.

Suppose that for each vulnerability $v_{l,j}$, the attacker can take a corresponding attack $a_{l,j} = \gamma_l(v_{l,j})$ from the action set $\mathcal{A}_l := \{a_{l,1}, a_{l,2}, \cdots, a_{l,n_l}\}$. Attack action $a_{l,j}$ is only effective and incurs a bounded cost $D_{ij} \in \mathbb{R}_+$ when the vulnerability $v_{l,j} = \gamma_l^{-1}(a_{l,j})$ exists in the current attack surface $\pi_l(c_{l,k})$. Thus, the damage caused by the attacker at stage $l$ can be represented as

$$r_l(a_{l,j}, c_{l,i}) = \begin{cases} D_{ij}, & \gamma_l^{-1}(a_{l,j}) \in \pi_l(c_{l,k}) \\ 0, & \text{otherwise} \end{cases}. \tag{5}$$

Since vulnerabilities are inevitable in a modern computing system, we can randomize the configuration and make it difficult for the attacker to learn and locate the system vulnerability, which naturally leads to the mixed strategy equilibrium solution concept of the game. At layer $l$, the defender's strategy $\mathbf{f}_l = \{f_{l,1}, f_{l,2}, \cdots, f_{l,m_l}\} \in \triangle C_l$ assigns probability $f_{l,j} \in [0, 1]$ to configuration $c_{l,j}$ while the attacker's strategy $\mathbf{g}_l := \{g_{l,1}, g_{l,2}, \cdots, g_{l,n_l}\} \in \triangle \mathcal{A}_l$ assigns probability $g_{l,i} \in [0, 1]$ to attack action $a_{l,i}$. The zero-sum game possesses a mixed strategy saddle-point equilibrium (SPE) $(\mathbf{f}_l^* \in \triangle C_l, \mathbf{g}_l^* \in \triangle \mathcal{A}_l)$, and a unique game value $\mathbb{r}(\mathbf{f}_l^*, \mathbf{g}_l^*)$, i.e.,

$$\mathbb{r}_l(\mathbf{f}_l^*, \mathbf{g}_l) \leq \mathbb{r}_l(\mathbf{f}_l^*, \mathbf{g}_l^*) \leq \mathbb{r}_l(\mathbf{f}_l, \mathbf{g}_l^*), \forall \mathbf{f}_l \in \triangle C_l, \mathbf{g}_l \in \triangle \mathcal{A}_l, \tag{6}$$

where the expected cost $\mathbb{r}_l$ is given by

$$\mathbb{r}_l(\mathbf{f}_l, \mathbf{g}_l) := \mathbb{E}_{\mathbf{f}_l, \mathbf{g}_l} r_l = \sum_{k=1}^{n_l} \sum_{h=1}^{m_l} f_{l,h} g_{l,k} r_l(a_{l,k}, c_{l,h}). \tag{7}$$

We illustrate the multistage MTD game in Fig. 5 and focus on the first layer with two available configurations $C_1 := \{c_{1,1}, c_{1,2}\}$ in the blue box. Configuration $c_{1,1}$ in Fig. 5a has an attack surface $\pi_1(c_{1,1}) = \{v_{1,1}, v_{1,2}\}$ while configuration $c_{1,2}$ in Fig. 5b reveals two vulnerabilities $v_{1,2}, v_{1,3} \in \pi_1(c_{1,2})$. Then, if the attacker takes action $a_{1,1}$ and the defender changes the configuration from $c_{1,1}$ to $c_{1,2}$, the attack is deterred at the first layer.



(a) Attack surface $\pi_1(c_{1,1}) = \{v_{1,1}, v_{1,2}\}$.        (b) Attack surface $\pi_1(c_{1,2}) = \{v_{1,2}, v_{1,3}\}$.
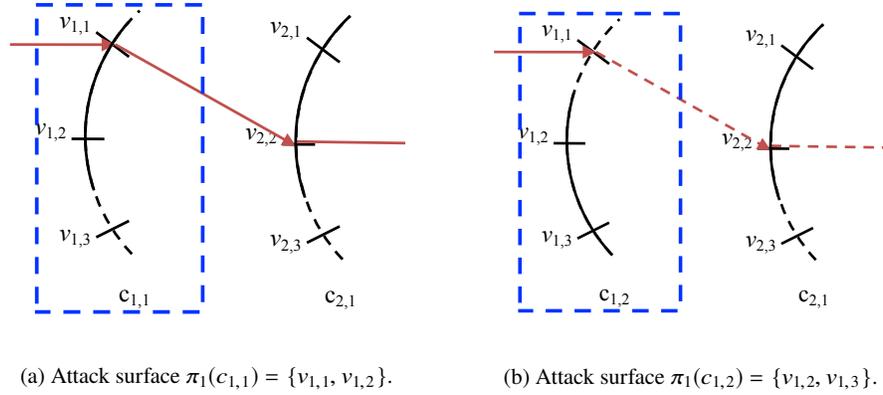
Fig. 5: Given a static configuration $c_{1,1}$, an attacker can succeed in reaching the resources at deeper layers by forming an attack path $v_{1,1} \rightarrow v_{2,2} \rightarrow \cdots$. A change of configuration to $c_{1,2}$ can thwart the attacker at the first layer.
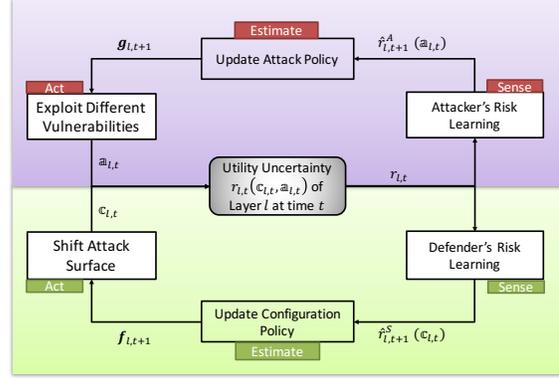
## 3.2 Distributed Learning

In practical cybersecurity domain, the payoff function $r_l$ is subjected to noises of unknown distributions. Then, each player reduces the payoff uncertainty by repeatedly observing the payoff realizations during the interaction with the other player. We use subscript $t$ to denote the strategy or cost at time $t$.

There is no communication at any time between two agents due to the non-cooperative environment, and the configuration and attack action are kept private, i.e., each player cannot observe the other player's action. Thus, each player independently chooses action $\mathbb{c}_{l,t} \in C_l$ or $\mathbb{a}_{l,t} \in \mathcal{A}_l$ to estimate the average risk of the system $\hat{r}_{l,t}^S : C_l \rightarrow \mathbb{R}_+$ and $\hat{r}_{l,t}^A : \mathcal{A}_l \rightarrow \mathbb{R}_+$ at layer $l$. Based on the estimated average risk

$\hat{r}_{l,t}^S$ and the previous policy $\mathbf{f}_{l,t}$, the defender can obtain her updated policy $\mathbf{f}_{l,t+1}$. Likewise, the attacker can also update his policy $\mathbf{g}_{l,t+1}$ based on $\hat{r}_{l,t}^A$ and $\mathbf{g}_{l,t}$. The new policy pair $(\mathbf{f}_{l,t+1}, \mathbf{g}_{l,t+1})$ determines the next payoff sample. The entire distributed learning feedback loop is illustrated in Fig. 6 where we distinguish the adversarial and defensive learning in red and green, respectively.



**Fig. 6** The distributed learning of the multistage MTD game at layer $l$. Adversarial learning in red does not share information with defensive learning in green. The distributed learning fashion means that the learning rule does not depend on the other player's action, yet the observed payoff depends on both players' actions.

In particular, players update their estimated average risks based on the payoff sample $r_{l,t}$ under the chosen action pair $(\mathbb{c}_{l,t}, \mathbb{a}_{l,t})$ as follows. Let $\mu_t^S$ and $\mu_t^A$ be the payoff learning rate for the system and attacker, respectively.

$$
\begin{aligned}
\hat{r}_{l,t+1}^S(c_{l,h}) &= \hat{r}_{l,t}^S(c_{l,h}) + \mu_t^S \mathbf{1}_{\{\mathbb{c}_{l,t}=c_{l,h}\}}(r_{l,t} - \hat{r}_{l,t}^S(c_{l,h})), \\
\hat{r}_{l,t+1}^A(a_{l,h}) &= \hat{r}_{l,t}^A(a_{l,h}) + \mu_t^A \mathbf{1}_{\{\mathbb{a}_{l,t}=a_{l,h}\}}(r_{l,t} - \hat{r}_{l,t}^A(a_{l,h})).
\end{aligned}
\tag{8}
$$

The indicators in (8) mean that both players only update the estimate average risk of the current action.

### 3.2.1 Security versus Usability

Frequent configuration changes may achieve the complete security yet also decrease the system usability. To quantify the tradeoff between the security and the usability, we introduce the switching cost of policy from $\mathbf{f}_{l,t}$ to $\mathbf{f}_{l,t+1}$ as their entropy:

$$
R_{l,t}^S := \sum_{h=1}^{m_l} f_{l,h,t+1} \ln\left(\frac{f_{l,h,t+1}}{f_{l,h,t}}\right).
\tag{9}
$$

Then, the total cost at time $t$ combines the expected cost with the entropy penalty in a ratio of $\epsilon_{l,t}^S$. When $\epsilon_{l,t}^S$ is high, the policy changes less and is more usable, yet may cause a large loss and be less rational.

$$(\text{SP}): \sup_{\mathbf{f}_{l,t+1} \in \triangle C_l} - \sum_{h=1}^{m_l} f_{l,h,t+1} \hat{r}_{l,t}^S(c_{l,h}) - \epsilon_{l,t}^S R_{l,t}^S. \tag{10}$$

A similar learning cost is introduced for the attacker:

$$(\text{AP}): \sup_{\mathbf{g}_{l,t+1} \in \triangle \mathcal{A}_l} - \sum_{h=1}^{n_l} g_{l,h,t+1} \hat{r}_{l,t}^A(a_{l,h}) - \epsilon_{l,t}^A \sum_{h=1}^{n_l} g_{l,h,t+1} \ln\left(\frac{g_{l,h,t+1}}{g_{l,h,t}}\right). \tag{11}$$

At any time $t+1$, we are able to obtain the equilibrium strategy $(f_{l,h,t+1}, g_{l,h,t+1})$ and game value $(W_{l,t}^S, W_{l,t}^A)$ in closed form of the previous strategy and the estimated average risk at time $t$ as follows.

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

$$W_{l,t}^S = \epsilon_{l,t}^S \ln\left(\sum_{h=1}^{m_l} f_{l,h,t} e^{-\frac{\hat{r}_{l,t}(c_{l,h})}{\epsilon_{l,t}^S}}\right), \qquad W_{l,t}^A = \epsilon_{l,t}^A \ln\left(\sum_{h=1}^{n_l} g_{l,h,t} e^{-\frac{\hat{r}_{l,t}(a_{l,h})}{\epsilon_{l,t}^A}}\right). \tag{12}$$

### 3.2.2 Learning Dynamics and ODE Counterparts

The closed form of policy leads to the following learning dynamics with learning rates $\lambda_{l,t}^S, \lambda_{l,t}^A \in [0, 1]$.

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

$$g_{l,h,t+1} = (1 - \lambda_{l,t}^A) g_{l,h,t} + \lambda_{l,t}^A \frac{g_{l,h,t} e^{-\frac{\hat{r}_{l,t}(a_{l,h})}{\epsilon_{l,t}^A}}}{\sum_{h'=1}^{n_l} g_{l,h',t} e^{-\frac{\hat{r}_{l,t}(a_{l,h'})}{\epsilon_{l,t}^A}}}. \tag{13}$$

If $\lambda_{l,t}^S = 1, \lambda_{l,t}^A = 1$, (13) is the same as (12). According to the stochastic approximation theory, the convergence of the policy and the average risk requires the learning rates $\lambda_{l,t}^A, \lambda_{l,t}^S, \mu_{l,t}^A, \mu_{l,t}^S$ to satisfy the regular condition of convergency in Definition 2.

**Definition 2** A number sequence $\{x_t\}, t = 1, 2, \cdots$, is said to satisfy the regular condition of convergency if

$$\sum_{t=1}^{\infty} x_t = +\infty, \quad \sum_{t=1}^{\infty} (x_t)^2 < +\infty. \tag{14}$$

The coupled dynamics of the payoff learning (8) and policy learning(13) converge to their Ordinary Differential Equations (ODEs) counterparts in system dynamics (15) and attacker dynamics (16), respectively. Let $e_{c_{l,h}} \in \triangle C_l, e_{a_{l,h}} \in \triangle \mathcal{A}_l$ be vectors of proper dimensions with the $h$-th entry being 1 and others being 0.

$$\frac{d}{dt} f_{l,h,t} = f_{l,h,t} \left( \frac{e^{-\frac{\hat{r}_{l,t}(c_{l,h})}{\epsilon_{l,t}^S}}}{\sum_{h'=1}^{m_l} f_{l,h',t} e^{-\frac{\hat{r}_{l,t}(c_{l,h'})}{\epsilon_{l,t}^S}}} - 1 \right), h = 1, 2, \cdots, m_l,$$

$$\frac{d}{dt} \hat{r}_{l,t}^S(c_{l,h}) = -\mathbb{r}_{l,t}(e_{c_{l,h}}, \mathbf{g}_{l,t}) - \hat{r}_{l,t+1}^S(c_{l,h}), c_{l,h} \in C_l. \tag{15}$$

$$\frac{d}{dt} g_{l,h,t+1} = g_{l,h,t} \left( \frac{e^{-\frac{\hat{r}_{l,t}(a_{l,h})}{\epsilon_{l,t}^A}}}{\sum_{h'=1}^{n_l} g_{l,h',t} e^{-\frac{\hat{r}_{l,t}(a_{l,h'})}{\epsilon_{l,t}^A}}} - 1 \right), h = 1, 2, \cdots, n_l,$$

$$\frac{d}{dt} \hat{r}_{l,t+1}^A(a_{l,h}) = \mathbb{r}_{l,t}(\mathbf{f}_{l,t}, e_{a_{l,h}}) - \hat{r}_{l,t+1}^A(a_{l,h}), a_{l,h} \in \mathcal{A}_l. \tag{16}$$

We can show that the SPE of the game is the steady state of the ODE dynamics in (15), (16), and the interior stationary points of the dynamics are the SPE of the game [6].
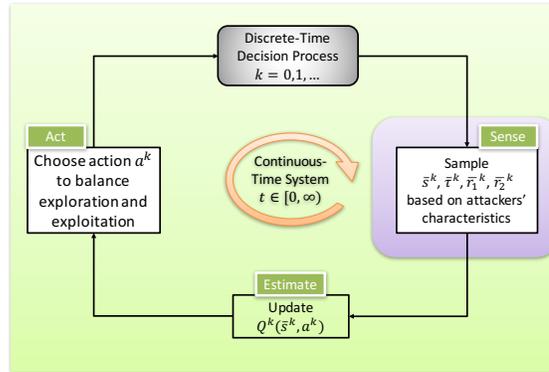
### 3.2.3 Heterogeneous and Hybrid Learning

The entropy regulation terms in (10) and (11) result in a closed form of strategies and learning dynamics in (13). Without the closed form, distributed learners can adopt general learning schemes which combine the payoff and the strategy update as stated in [46]. Specifically, algorithm CRL0 mimics the replicator dynamics and updates the strategy according to the current sample value of the utility. On the other hand, algorithm CRL1 updates the strategy according to a soft-max function of the estimated utilities so that the most rewarding policy get reinforced and will be picked with a higher probability. The first algorithm is robust yet inefficient, and the second one is fragile yet efficient. Moreover, players are not obliged to adopt the same learning scheme at different time. The heterogeneous learning focuses on different players adopting different learning schemes [46], while hybrid learning means that players can choose different learning schemes at different times based on their rationalities and preferences [63]. According to stochastic approximation

techniques, these learning schemes with random updates can be studied using their deterministic ODE counterparts.

# 4 Reinforcement Learning for Uncertain Environments

This section considers uncertainties on the entire environment, i.e., the state transition, the sojourn time, and the investigation payoff, in the active defense scenario of the honeypot engagement [7]. We use the Semi-Markov Decision Process (SMDP) to capture these environmental uncertainties in the continuous time system. Although the attacker's duration time is continuous at each honeypot, the defender's engagement action is applied at a discrete time epoch. Based on the observed samples at each decision epoch, the defender can estimate the environment elements determined by attackers' characteristics, and use reinforcement learning methods to obtain the optimal policy. We plot the entire feedback learning structure in Fig. 7. Since the attacker should not identify the existence of the honeypot and the defender's engagement actions, he will not take actions to jeopardize the learning.

**Fig. 7** The feedback structure of reinforcement learning methods on SMDP. The red background means that the attacker's characteristics determine the environmental uncertainties and the samples observed in the honeynet. The attacker is not involved in parts of the green background. The learning scheme in Fig. 7 extends the one in Section 3 to consider a continuous time elapse and multistage transitions.



## 4.1 Honeypot Network and SMDP Model

The honeypots form a network to emulate a production system. From an attacker's viewpoint, two network structures are the same as shown in Fig. 8. Based on the network topology, we introduce the continuous-time infinite-horizon discounted SMDPs, which can be summarized by the tuple $\{t \in [0, \infty), \mathcal{S}, \mathcal{A}(s_j), tr(s_l|s_j, a_j), z(\cdot|s_j, a_j, s_l), r^\gamma(s_j, a_j, s_l), \gamma \in [0, \infty)\}$. We illustrate each element of the tuple through a 13-state example in Fig. 9.
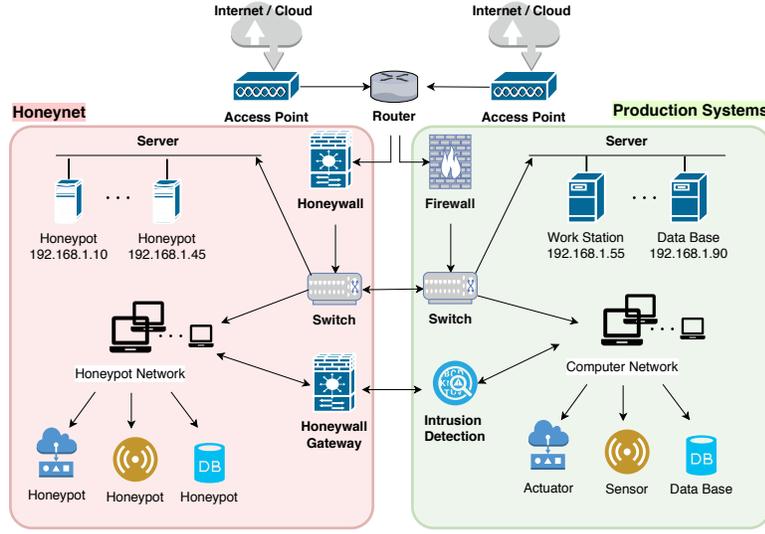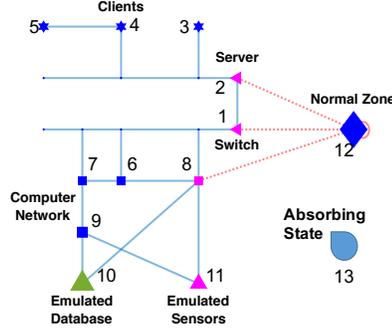
Fig. 8: The honeynet in red emulates and shares the same structure as the targeted production system in green.

Each node in Fig. 9 represents a state $s_i \in \mathcal{S}, i \in \{1, 2, \cdots, 13\}$. At time $t \in [0, \infty)$, the attacker is either at one of the honeypot node denoted by state $s_i \in \mathcal{S}, i \in \{1, 2, \cdots, 11\}$, at the normal zone $s_{12}$, or at a virtual absorbing state $s_{13}$ once attackers are ejected or terminate on their own. At each state $s_i \in \mathcal{S}$, the defender can choose an action $a_i \in \mathcal{A}(s_i)$. For example, at honeypot nodes, the defender can conduct action $a_E$ to eject the attacker, action $a_P$ to purely record the attacker's activities, low-interactive action $a_L$, or high-interactive action $a_H$, i.e., $\mathcal{A}(s_i) := \{a_E, a_P, a_L, a_H\}, i \in \{1, \cdots, N\}$. The high-interactive action is costly to implement yet can both increases the probability of a longer sojourn time at honeypot $n_i$, and reduces the probability of attackers penetrating the normal system from $n_i$ if connected. If the attacker resides in the normal zone either from the beginning or later through the pivot honeypots, the defender can choose either action $a_E$ to eject the attacker immediately, or action $a_A$ to attract the attacker to the honeynet by generating more deceptive inbound and outbound traffics in the honeynet, i.e., $\mathcal{A}(s_{12}) := \{a_E, a_A\}$. Based on the current state $s_j \in \mathcal{S}$ and the defender's action $a_j \in \mathcal{A}(s_j)$, the attacker transits to state $s_l \in \mathcal{S}$ with probability $tr(s_l|s_j, a_j)$ and the sojourn time at state $s_j$ is a continuous random variable with probability density $z(\cdot|s_j, a_j, s_l)$. Once the attacker arrives at a new honeypot $n_i$, the defender dynamically applies an interaction action at honeypot $n_i$ from $\mathcal{A}(s_i)$ and keeps interacting with the attacker until she transits to the next honeypot. If the defender changes the action before the transition, the attacker may be able to detect the change and become

**Fig. 9** Honeypots emulate different components of the production system. Actions $a_E, a_P, a_L, a_H$ are denoted in red, blue, purple, and green, respectively. The size of node $n_i$ represents the state value $v(s_i), i \in \{1, 2, \cdots, 11\}$.



aware of the honeypot. Since the decision is made at the time of transition, we can transform the above continuous time model on horizon $t \in [0, \infty)$ into a discrete decision model at decision epoch $k \in \{0, 1, \cdots, \infty\}$. The time of the attacker's $k^{th}$ transition is denoted by a random variable $T^k$, the landing state is denoted as $s^k \in \mathcal{S}$, and the adopted action after arriving at $s^k$ is denoted as $a^k \in \mathcal{A}(s^k)$.

The defender gains an investigation reward by engaging and analyzing the attacker in the honeypot. To simplify the notation, we segment the investigation reward during time $t \in [0, \infty)$ into ones at discrete decision epochs $T^k, k \in \{0, 1, \cdots, \infty\}$. When $\tau \in [T^k, T^{k+1}]$ amount of time elapses at stage $k$, the defender's investigation reward $r(s^k, a^k, s^{k+1}, T^k, T^{k+1}, \tau) = r_1(s^k, a^k, s^{k+1})\mathbf{1}_{\{\tau=0\}} + r_2(s^k, a^k, T^k, T^{k+1}, \tau)$, at time $\tau$ of stage $k$, is the sum of two parts. The first part is the immediate cost of applying engagement action $a^k \in \mathcal{A}(s^k)$ at state $s^k \in \mathcal{S}$ and the second part is the reward rate of threat information acquisition minus the cost rate of persistently generating deceptive traffics. Due to the randomness of the attacker's behavior, the information acquisition can also be random, thus the actual reward rate $r_2$ is perturbed by an additive zero-mean noise $w_r$. As the defender spends longer time interacting with attackers, investigating their behaviors and acquires better understandings of their targets and TTPs, less new information can be extracted. In addition, the same intelligence becomes less valuable as time elapses due to the timeliness. Thus, we use a discounted factor of $\gamma \in [0, \infty)$ to penalize the decreasing value of the investigation reward as time elapses.

The defender aims at a policy $\pi \in \Pi$ which maps state $s^k \in \mathcal{S}$ to action $a^k \in \mathcal{A}(s^k)$ to maximize the long-term expected utility starting from state $s^0$, i.e.,

$$u(s^0, \pi) = E[\sum_{k=0}^{\infty} \int_{T^k}^{T^{k+1}} e^{-\gamma(\tau+T^k)}(r(S^k, A^k, S^{k+1}, T^k, T^{k+1}, \tau) + w_r)d\tau]. \quad (17)$$

At each decision epoch, the value function $v(s^0) = \sup_{\pi \in \Pi} u(s^0, \pi)$ can be represented by dynamic programming, i.e.,

$$v(s^0) = \sup_{a^0 \in \mathcal{A}(s^0)} E[\int_{T^0}^{T^1} e^{-\gamma(\tau+T^0)} r(s^0, a^0, S^1, T^0, T^1, \tau) d\tau + e^{-\gamma T^1} v(S^1)]. \quad (18)$$

We assume a constant reward rate $r_2(s^k, a^k, T^k, T^{k+1}, \tau) = \bar{r}_2(s^k, a^k)$ for simplicity. Then, (18) can be transformed into an equivalent MDP form, i.e., $\forall s^0 \in \mathcal{S}$,

$$v(s^0) = \sup_{a^0 \in \mathcal{A}(s^0)} \sum_{s^1 \in \mathcal{S}} tr(s^1|s^0, a^0)(r^\gamma(s^0, a^0, s^1) + z^\gamma(s^0, a^0, s^1)v(s^1)), \quad (19)$$

where $z^\gamma(s^0, a^0, s^1) := \int_0^\infty e^{-\gamma\tau} z(\tau|s^0, a^0, s^1) d\tau \in [0, 1]$ is the Laplace transform of the sojourn probability density $z(\tau|s^0, a^0, s^1)$ and the equivalent reward $r^\gamma(s^0, a^0, s^1)$ $:= r_1(s^0, a^0, s^1) + \frac{\bar{r}_2(s^0, a^0)}{\gamma}(1 - z^\gamma(s^0, a^0, s^1)) \in [-m_c, m_c]$ is assumed to be bounded by a constant $m_c$.

**Definition 3** There exists constants $\theta \in (0, 1)$ and $\delta > 0$ such that

$$\sum_{s^1 \in \mathcal{S}} tr(s^1|s^0, a^0)z(\delta|s^0, a^0, s^1) \leq 1 - \theta, \forall s^0 \in \mathcal{S}, a^0 \in \mathcal{A}(s^0). \quad (20)$$

The right-hand side of (18) is a contraction mapping under the regulation condition in Definition 3. Then, we can find the unique optimal policy $\pi^* = arg \max_{\pi \in \Pi} u(s^0, \pi)$ by value iteration, policy iteration or linear programming. Fig. 9 illustrates the optimal policy and the state value by the color and the size of the node, respectively. In the example scenario, the honeypot of database $n_{10}$ and sensors $n_{11}$ are the main and secondary targets of the attacker, respectively. Thus, defenders can obtain a higher investigation reward when they manage to engage the attacker in these two honeypot nodes with a larger probability and for a longer time. However, instead of naively adopting high interactive actions, a savvy defender also balances the high implantation cost of $a_H$. Our quantitative results indicate that the high interactive action should only be applied at $n_{10}$ to be cost-effective. On the other hand, although the bridge nodes $n_1, n_2, n_8$ which connect to the normal zone $n_{12}$ do not contain higher investigation rewards than other nodes, the defender still takes action $a_L$ at these nodes. The goal is to either increase the probability of attracting attackers away from the normal zone or reduce the probability of attackers penetrating the normal zone from these bridge nodes.

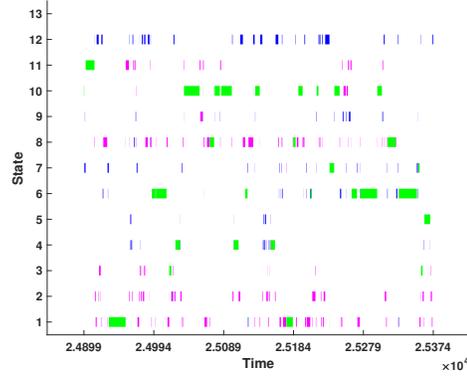## 4.2 Reinforcement Learning of SMDP

The absent knowledge of the attacker's characteristics results in environmental uncertainty of the investigation reward, the attacker's transition probability, and the sojourn distribution. We use $Q$-learning algorithm to obtain the optimal engagement policy based on the actual experience of the honeynet interactions, i.e., $\forall \bar{s}^k \in \mathcal{S}, \forall a^k \in \mathcal{A}(\bar{s}^k)$,

$$Q^{k+1}(\bar{s}^k, a^k) = (1 - \alpha^k(\bar{s}^k, a^k))Q^k(\bar{s}^k, a^k) + \alpha^k(\bar{s}^k, a^k)[\bar{r}_1(\bar{s}^k, a^k, \bar{s}^{k+1})$$

$$+ \bar{r}_2(\bar{s}^k, a^k)\frac{(1 - e^{-\gamma \bar{\tau}^k})}{\gamma} - e^{-\gamma \bar{\tau}^k} \max_{a' \in \mathcal{A}(\bar{s}^{k+1})} Q^k(\bar{s}^{k+1}, a')], \quad (21)$$

where $\alpha^k(\bar{s}^k, a^k) \in (0, 1)$ is the learning rate, $\bar{s}^k, \bar{s}^{k+1}$ are the observed states at stage $k$ and $k + 1$, $\bar{r}_1, \bar{r}_2$ is the observed investigation rewards, and $\bar{\tau}^k$ is the observed sojourn time at state $s^k$. When the learning rate satisfies the condition of convergency in Definition 2, i.e., $\sum_{k=0}^{\infty} \alpha^k(s^k, a^k) = \infty, \sum_{k=0}^{\infty}(\alpha^k(s^k, a^k))^2 < \infty, \forall s^k \in \mathcal{S}, \forall a^k \in \mathcal{A}(s^k)$, and all state-action pairs are explored infinitely, $\max_{a' \in \mathcal{A}(s^k)} Q^{\infty}(s^{\infty}, a')$, in (21) converges to value $v(s^k)$ with probability 1.

At each decision epoch $k \in \{0, 1, \cdots\}$, the action $a^k$ is chosen according to the $\epsilon$-greedy policy, i.e., the defender chooses the optimal action $arg \max_{a' \in \mathcal{A}(s^k)} Q^k(s^k, a')$ with a probability $1 - \epsilon$, and a random action with a probability $\epsilon$. Note that the exploration rate $\epsilon \in (0, 1]$ should not be too small to guarantee sufficient samples of all state-action pairs. The $Q$-learning algorithm under a pure exploration policy $\epsilon = 1$ still converges yet at a slower rate.



**Fig. 10** One instance of $Q$-learning on SMDP where the $x$-axis shows the sojourn time and the $y$-axis represents the state transition. The chosen actions $a_E, a_P, a_L, a_H$ are denoted in red, blue, purple, and green, respectively.
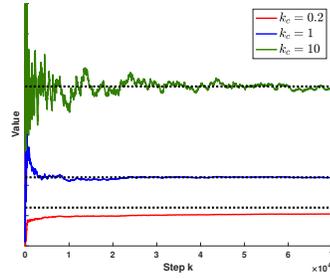
In our scenario, the defender knows the reward of ejection action $a_A$ and $v(s_{13}) = 0$, thus does not need to explore action $a_A$ to learn it. We plot one learning trajectory of the state transition and sojourn time under the $\epsilon$-greedy exploration policy in Fig. 10, where the chosen actions $a_E, a_P, a_L, a_H$ are denoted in red, blue, purple, and green, respectively. If the ejection reward is unknown, the defender should be restrictive in exploring $a_A$ which terminates the learning process. Otherwise, the defender may need to engage with a group of attackers who share similar behaviors to obtain sufficient samples to learn the optimal engagement policy.

In particular, we choose $\alpha^k(s^k, a^k) = \frac{k_c}{k_{\{s^k, a^k\}} - 1 + k_c}, \forall s^k \in \mathcal{S}, \forall a^k \in \mathcal{A}(s^k)$, to guarantee the asymptotic convergence, where $k_c \in (0, \infty)$ is a constant parameter and $k_{\{s^k, a^k\}} \in \{0, 1, \cdots\}$ is the number of visits to state-action pair $\{s^k, a^k\}$ up to stage $k$. We need to choose a proper value of $k_c$ to guarantee a good numerical

performance of convergence in finite steps as shown in Fig. 11a. We shift the green and blue lines vertically to avoid the overlap with the red line and represent the corresponding theoretical values in dotted black lines. If $k_c$ is too small as shown in the red line, the learning rate decreases so fast that new observed samples hardly update the $Q$-value and the defender may need a long time to learn the right value. However, if $k_c$ is too large as shown in the green line, the learning rate decreases so slow that new samples contribute significantly to the current $Q$-value. It causes a large variation and a slower convergence rate of $\max_{a' \in \mathcal{A}(s_{12})} Q^k(s_{12}, a')$.

We show the convergence of the policy and value under $k_c = 1, \epsilon = 0.2$, in the video demo (See URL: https://bit.ly/2QUz3Ok). In the video, the color of each node $n^k$ distinguishes the defender's action $a^k$ at state $s^k$ and the size of the node is proportional to $\max_{a' \in \mathcal{A}(s^k)} Q^k(s^k, a')$ at stage $k$. To show the convergence, we decrease the value of $\epsilon$ gradually to 0 after 5000 steps. Since the convergence trajectory is stochastic, we run the simulation for 100 times and plot the mean and the variance of $Q^k(s_{12}, a_P)$ of state $s_{12}$ under the optimal policy $\pi(s_{12}) = a_P$ in Fig. 11. The mean in red converges to the theoretical value in about 400 steps and the variance in blue reduces dramatically as step $k$ increases.



(a) The convergence rate under different values of $k_c$.

(b) The evolution of the mean and the variance of $Q^k(s_{12}, a_P)$.

Fig. 11: Convergence results of $Q$-learning over SMDP.

## 5 Conclusion and Discussion

This chapter has introduced three defense schemes, i.e., defensive deception to detect and counter adversarial deception, feedback-driven Moving Target Defense (MTD) to increase the attacker's probing and reconnaissance costs, and adaptive honeypot engagement to gather fundamental threat information. These schemes satisfy the Principle of 3A Defense as they actively protect the system prior to the attack damages, provide strategic defenses autonomously, and apply learning to adapt to uncertainty and changes. These schemes possess three progressive levels

of information restrictions, which lead to different strategic learning schemes to estimate the parameter, the payoff, and the environment. All these learning schemes, however, have a feedback loop to sense samples, estimate the unknowns, and take actions according to the estimate. Our work lays a solid foundation for strategic learning in active, adaptive, autonomous defenses under incomplete information and leads to the following challenges and future directions.

First, multi-agent learning in non-cooperative environments is challenging due to the coupling and interaction between these heterogeneous agents. The learning results depend on all involving agents yet other players' behaviors, levels of rationality, and learning schemes are not controllable and may change abruptly. Moreover, as attackers become aware of the active defense techniques and the learning scheme under incomplete information, the savvy attacker can attempt to interrupt the learning process. For example, attackers may sacrifice their immediate rewards and take incomprehensible actions instead so that the defender learns incorrect attack characteristics. The above challenges motivate robust learning methods under non-cooperative and even adversarial environments.

Second, since the learning process is based on samples from real interactions, the defender needs to concern the system safety and security during the learning period, while in the same time, attempts to achieve more accurate learning results of the attack's characteristics. Moreover, since the learning under non-cooperative and adversarial environments may terminate unpredictably at any time, the asymptotic convergence would not be critical for security. The defender needs to care more about the time efficiency of the learning, i.e., how to achieve a sufficiently good estimate in a finite number of steps.

Third, instead of learning from scratch, the defender can attempt to reuse the past experience with attackers of similar behaviors to expedite the learning process, which motivates the investigation of transfer learning in reinforcement learning [69]. Some side-channel information may also contribute to the learning to allow agents to learn faster.

# References

1. "Verizon 2019 data breach investigations report," 2019.
2. D. Shackleford, "Combatting cyber risks in the supply chain," *SANS. org*, 2015.
3. L. Huang and Q. Zhu, "Adaptive strategic cyber defense for advanced persistent threats in critical infrastructure networks," *ACM SIGMETRICS Performance Evaluation Review*, vol. 46, no. 2, pp. 52–56, 2019.
4. ——, "Analysis and computation of adaptive defense strategies against advanced persistent threats for cyber-physical systems," in *International Conference on Decision and Game Theory for Security*. Springer, 2018, pp. 205–226.
5. L. Huang and Q. Zhu, "A Dynamic Games Approach to Proactive Defense Strategies against Advanced Persistent Threats in Cyber-Physical Systems," *arXiv e-prints*, p. arXiv:1906.09687, Jun 2019.
6. Q. Zhu and T. Başar, "Game-theoretic approach to feedback-driven multi-stage moving target defense," in *International Conference on Decision and Game Theory for Security*. Springer, 2013, pp. 246–263.

7. L. Huang and Q. Zhu, "Adaptive Honeypot Engagement through Reinforcement Learning of Semi-Markov Decision Processes," *arXiv e-prints*, p. arXiv:1906.12182, Jun 2019.

8. J. Pawlick, E. Colbert, and Q. Zhu, "A game-theoretic taxonomy and survey of defensive deception for cybersecurity and privacy," *arXiv preprint arXiv:1712.05441*, 2017.

9. F. J. Stech, K. E. Heckman, and B. E. Strom, "Integrating cyber-d&d into adversary modeling for active cyber defense," in *Cyber deception*. Springer, 2016, pp. 1–22.

10. K. E. Heckman, M. J. Walsh, F. J. Stech, T. A. O'boyle, S. R. DiCato, and A. F. Herber, "Active cyber defense with denial and deception: A cyber-wargame experiment," *computers & security*, vol. 37, pp. 72–77, 2013.

11. J. Gómez-Hernández, L. Álvarez-González, and P. García-Teodoro, "R-locker: Thwarting ransomware action through a honeyfile-based approach," *Computers & Security*, vol. 73, pp. 389–398, 2018.

12. N. Virvilis, B. Vanautgaerden, and O. S. Serrano, "Changing the game: The art of deceiving sophisticated attackers," in *2014 6th International Conference On Cyber Conflict (CyCon 2014)*. IEEE, 2014, pp. 87–97.

13. J. Pawlick, E. Colbert, and Q. Zhu, "Modeling and analysis of leaky deception using signaling games with evidence," *IEEE Transactions on Information Forensics and Security*, 2018.

14. S. Jajodia, A. K. Ghosh, V. Swarup, C. Wang, and X. S. Wang, *Moving target defense: creating asymmetric uncertainty for cyber threats*. Springer Science & Business Media, 2011, vol. 54.

15. G. S. Kc, A. D. Keromytis, and V. Prevelakis, "Countering code-injection attacks with instruction-set randomization," in *Proceedings of the 10th ACM conference on Computer and communications security*. ACM, 2003, pp. 272–280.

16. A. Clark, Q. Zhu, R. Poovendran, and T. Başar, "Deceptive routing in relay networks," in *International Conference on Decision and Game Theory for Security*. Springer, 2012, pp. 171–185.

17. H. Maleki, S. Valizadeh, W. Koch, A. Bestavros, and M. van Dijk, "Markov modeling of moving target defense games," in *Proceedings of the 2016 ACM Workshop on Moving Target Defense*. ACM, 2016, pp. 81–92.

18. C. R. Hecker, "A methodology for intelligent honeypot deployment and active engagement of attackers," Ph.D. dissertation, 2012.

19. Q. D. La, T. Q. Quek, J. Lee, S. Jin, and H. Zhu, "Deceptive attack and defense game in honeypot-enabled networks for the internet of things," *IEEE Internet of Things Journal*, vol. 3, no. 6, pp. 1025–1035, 2016.

20. J. Pawlick, T. T. H. Nguyen, and Q. Zhu, "Optimal timing in dynamic and robust attacker engagement during advanced persistent threats," *CoRR*, vol. abs/1707.08031, 2017. [Online]. Available: http://arxiv.org/abs/1707.08031

21. J. Pawlick and Q. Zhu, "A Stackelberg game perspective on the conflict between machine learning and data obfuscation," in *Information Forensics and Security (WIFS), 2016 IEEE International Workshop on*. IEEE, 2016, pp. 1–6. [Online]. Available: http://ieeexplore.ieee.org/abstract/document/7823893/

22. Q. Zhu, A. Clark, R. Poovendran, and T. Basar, "Deployment and exploitation of deceptive honeybots in social networks," in *Decision and Control (CDC), 2013 IEEE 52nd Annual Conference on*. IEEE, 2013, pp. 212–219.

23. Q. Zhu, H. Tembine, and T. Basar, "Hybrid learning in stochastic games and its applications in network security," *Reinforcement Learning and Approximate Dynamic Programming for Feedback Control*, pp. 305–329, 2013.

24. Q. Zhu, Z. Yuan, J. B. Song, Z. Han, and T. Başar, "Interference aware routing game for cognitive radio multi-hop networks," *Selected Areas in Communications, IEEE Journal on*, vol. 30, no. 10, pp. 2006–2015, 2012.

25. Q. Zhu, L. Bushnell, and T. Basar, "Game-theoretic analysis of node capture and cloning attack with multiple attackers in wireless sensor networks," in *Decision and Control (CDC), 2012 IEEE 51st Annual Conference on*. IEEE, 2012, pp. 3404–3411.

26. Q. Zhu, A. Clark, R. Poovendran, and T. Başar, "Deceptive routing games," in *Decision and Control (CDC), 2012 IEEE 51st Annual Conference on*. IEEE, 2012, pp. 2704–2711.

27. Q. Zhu, H. Li, Z. Han, and T. Basar, "A stochastic game model for jamming in multi-channel cognitive radio systems." in *ICC*, 2010, pp. 1–6.
28. Z. Xu and Q. Zhu, "Secure and practical output feedback control for cloud-enabled cyber-physical systems," in *Communications and Network Security (CNS), 2017 IEEE Conference on*. IEEE, 2017, pp. 416–420.
29. ——, "A Game-Theoretic Approach to Secure Control of Communication-Based Train Control Systems Under Jamming Attacks," in *Proceedings of the 1st International Workshop on Safe Control of Connected and Autonomous Vehicles*. ACM, 2017, pp. 27–34. [Online]. Available: http://dl.acm.org/citation.cfm?id=3055381
30. ——, "Cross-layer secure cyber-physical control system design for networked 3d printers," in *American Control Conference (ACC), 2016*. IEEE, 2016, pp. 1191–1196. [Online]. Available: http://ieeexplore.ieee.org/abstract/document/7525079/
31. M. J. Farooq and Q. Zhu, "Modeling, analysis, and mitigation of dynamic botnet formation in wireless iot networks," *IEEE Transactions on Information Forensics and Security*, 2019.
32. Z. Xu and Q. Zhu, "A cyber-physical game framework for secure and resilient multi-agent autonomous systems," in *Decision and Control (CDC), 2015 IEEE 54th Annual Conference on*. IEEE, 2015, pp. 5156–5161.
33. L. Huang, J. Chen, and Q. Zhu, "A large-scale markov game approach to dynamic protection of interdependent infrastructure networks," in *International Conference on Decision and Game Theory for Security*. Springer, 2017, pp. 357–376.
34. J. Chen, C. Touati, and Q. Zhu, "A dynamic game analysis and design of infrastructure network protection and recovery," *ACM SIGMETRICS Performance Evaluation Review*, vol. 45, no. 2, p. 128, 2017.
35. F. Miao, Q. Zhu, M. Pajic, and G. J. Pappas, "A hybrid stochastic game for secure control of cyber-physical systems," *Automatica*, vol. 93, pp. 55–63, 2018.
36. Y. Yuan, Q. Zhu, F. Sun, Q. Wang, and T. Basar, "Resilient control of cyber-physical systems against denial-of-service attacks," in *Resilient Control Systems (ISRCS), 2013 6th International Symposium on*. IEEE, 2013, pp. 54–59.
37. S. Rass and Q. Zhu, "GADAPT: A Sequential Game-Theoretic Framework for Designing Defense-in-Depth Strategies Against Advanced Persistent Threats," in *Decision and Game Theory for Security*, ser. Lecture Notes in Computer Science, Q. Zhu, T. Alpcan, E. Panaousis, M. Tambe, and W. Casey, Eds. Cham: Springer International Publishing, 2016, vol. 9996, pp. 314–326.
38. Q. Zhu, Z. Yuan, J. B. Song, Z. Han, and T. Basar, "Dynamic interference minimization routing game for on-demand cognitive pilot channel," in *Global Telecommunications Conference (GLOBECOM 2010), 2010 IEEE*. IEEE, 2010, pp. 1–6.
39. T. Zhang and Q. Zhu, "Strategic defense against deceptive civilian gps spoofing of unmanned aerial vehicles," in *International Conference on Decision and Game Theory for Security*. Springer, 2017, pp. 213–233.
40. L. Huang and Q. Zhu, "Analysis and computation of adaptive defense strategies against advanced persistent threats for cyber-physical systems," in *International Conference on Decision and Game Theory for Security*, 2018.
41. ——, "Adaptive strategic cyber defense for advanced persistent threats in critical infrastructure networks," in *ACM SIGMETRICS Performance Evaluation Review*, 2018.
42. J. Pawlick, S. Farhang, and Q. Zhu, "Flip the cloud: Cyber-physical signaling games in the presence of advanced persistent threats," in *Decision and Game Theory for Security*. Springer, 2015, pp. 289–308.
43. S. Farhang, M. H. Manshaei, M. N. Esfahani, and Q. Zhu, "A dynamic bayesian security game framework for strategic defense mechanism design," in *Decision and Game Theory for Security*. Springer, 2014, pp. 319–328.
44. Q. Zhu and T. Başar, "Dynamic policy-based ids configuration," in *Decision and Control, 2009 held jointly with the 2009 28th Chinese Control Conference. CDC/CCC 2009. Proceedings of the 48th IEEE Conference on*. IEEE, 2009, pp. 8600–8605.
45. Q. Zhu, H. Tembine, and T. Basar, "Network security configurations: A nonzero-sum stochastic game approach," in *American Control Conference (ACC), 2010*. IEEE, 2010, pp. 1059–1064.

46. Q. Zhu, H. Tembine, and T. Başar, "Heterogeneous learning in zero-sum stochastic games with incomplete information," in *49th IEEE conference on decision and control (CDC)*.    IEEE, 2010, pp. 219–224.

47. J. Chen and Q. Zhu, "Security as a Service for Cloud-Enabled Internet of Controlled Things under Advanced Persistent Threats: A Contract Design Approach," *IEEE Transactions on Information Forensics and Security*, 2017. [Online]. Available: http://ieeexplore.ieee.org/abstract/document/7954676/

48. R. Zhang, Q. Zhu, and Y. Hayel, "A Bi-Level Game Approach to Attack-Aware Cyber Insurance of Computer Networks," *IEEE Journal on Selected Areas in Communications*, vol. 35, no. 3, pp. 779–794, 2017. [Online]. Available: http://ieeexplore.ieee.org/abstract/document/7859343/

49. R. Zhang and Q. Zhu, "Attack-aware cyber insurance of interdependent computer networks," 2016.

50. W. A. Casey, Q. Zhu, J. A. Morales, and B. Mishra, "Compliance control: Managed vulnerability surface in social-technological systems via signaling games," in *Proceedings of the 7th ACM CCS International Workshop on Managing Insider Security Threats*.    ACM, 2015, pp. 53–62.

51. Y. Hayel and Q. Zhu, "Attack-aware cyber insurance for risk sharing in computer networks," in *Decision and Game Theory for Security*.    Springer, 2015, pp. 22–34.

52. ——, "Epidemic protection over heterogeneous networks using evolutionary poisson games," *IEEE Transactions on Information Forensics and Security*, vol. 12, no. 8, pp. 1786–1800, 2017.

53. Q. Zhu, C. Fung, R. Boutaba, and T. Başar, "Guidex: A game-theoretic incentive-based mechanism for intrusion detection networks," *Selected Areas in Communications, IEEE Journal on*, vol. 30, no. 11, pp. 2220–2230, 2012.

54. Q. Zhu, C. A. Gunter, and T. Basar, "Tragedy of anticommons in digital right management of medical records." in *HealthSec*, 2012.

55. Q. Zhu, C. Fung, R. Boutaba, and T. Başar, "A game-theoretical approach to incentive design in collaborative intrusion detection networks," in *Game Theory for Networks, 2009. GameNets' 09. International Conference on*.    IEEE, 2009, pp. 384–392.

56. T. E. Carroll and D. Grosu, "A game theoretic investigation of deception in network security," *Security and Commun. Nets.*, vol. 4, no. 10, pp. 1162–1172, 2011.

57. J. Pawlick and Q. Zhu, "A Stackelberg game perspective on the conflict between machine learning and data obfuscation," *IEEE Intl. Workshop on Inform. Forensics and Security*, 2016.

58. T. Zhang and Q. Zhu, "Dynamic differential privacy for ADMM-based distributed classification learning," *IEEE Transactions on Information Forensics and Security*, vol. 12, no. 1, pp. 172–187, 2017. [Online]. Available: http://ieeexplore.ieee.org/abstract/document/7563366/

59. S. Farhang, Y. Hayel, and Q. Zhu, "Phy-layer location privacy-preserving access point selection mechanism in next-generation wireless networks," in *Communications and Network Security (CNS), 2015 IEEE Conference on*.    IEEE, 2015, pp. 263–271.

60. T. Zhang and Q. Zhu, "Distributed privacy-preserving collaborative intrusion detection systems for vanets," *IEEE Transactions on Signal and Information Processing over Networks*, vol. 4, no. 1, pp. 148–161, 2018.

61. N. Zhang, W. Yu, X. Fu, and S. K. Das, "gPath: A game-theoretic path selection algorithm to protect tor's anonymity," in *Decision and Game Theory for Security*.    Springer, 2010, pp. 58–71.

62. A. Garnaev, M. Baykal-Gursoy, and H. V. Poor, "Security games with unknown adversarial strategies," *IEEE transactions on cybernetics*, vol. 46, no. 10, pp. 2291–2299, 2015.

63. Q. Zhu, H. Tembine, and T. Başar, "Distributed strategic learning with application to network security," in *Proceedings of the 2011 American Control Conference*.    IEEE, 2011, pp. 4057–4062.

64. A. Servin and D. Kudenko, "Multi-agent reinforcement learning for intrusion detection: A case study and evaluation," in *German Conference on Multiagent System Technologies*.    Springer, 2008, pp. 159–170.

65. P. M. Djurić and Y. Wang, "Distributed bayesian learning in multiagent systems: Improving our understanding of its capabilities and limitations," *IEEE Signal Processing Magazine*, vol. 29, no. 2, pp. 65–76, 2012.

66. G. Chalkiadakis and C. Boutilier, "Coordination in multiagent reinforcement learning: a bayesian approach," in *Proceedings of the second international joint conference on Autonomous agents and multiagent systems*. ACM, 2003, pp. 709–716.
67. Z. Chen and D. Marculescu, "Distributed reinforcement learning for power limited many-core system performance optimization," in *Proceedings of the 2015 Design, Automation & Test in Europe Conference & Exhibition*. EDA Consortium, 2015, pp. 1521–1526.
68. J. C. Harsanyi, "Games with incomplete information played by "bayesian" players, i–iii part i. the basic model," *Management science*, vol. 14, no. 3, pp. 159–182, 1967.
69. M. E. Taylor and P. Stone, "Transfer learning for reinforcement learning domains: A survey," *Journal of Machine Learning Research*, vol. 10, no. Jul, pp. 1633–1685, 2009.