# Asymptotic Security using Bayesian Defense Mechanism with Application to Cyber Deception

Hampei Sasahara, *Member, IEEE*, Henrik Sandberg, *Senior Member, IEEE*

*Abstract*— **This paper addresses the question whether model knowledge can guide a defender to appropriate decisions, or not, when an attacker intrudes into control systems. The model-based defense scheme considered in this study, namely Bayesian defense mechanism, chooses reasonable reactions through observation of the system's behavior using models of the system's stochastic dynamics, the vulnerability to be exploited, and the attacker's objective. On the other hand, rational attackers take deceptive strategies for misleading the defender into making inappropriate decisions. In this paper, their dynamic decision making is formulated as a stochastic signaling game. It is shown that the belief of the true scenario has a limit in a stochastic sense at an equilibrium based on martingale analysis. This fact implies that there are only two possible cases: the defender asymptotically detects the attack with a firm belief, or the attacker takes actions such that the system's behavior becomes nominal after a finite number of time steps. Consequently, if different scenarios result in different stochastic behaviors, the Bayesian defense mechanism guarantees the system to be secure in an asymptotic manner provided that effective countermeasures are implemented. As an application of the finding, a defensive deception utilizing asymmetric recognition of vulnerabilities exploited by the attacker is analyzed. It is shown that the attacker possibly withdraws even if the defender is unaware of the exploited vulnerabilities, as long as the defender's unawareness is concealed by the defensive deception.**

*Index Terms*— **Bayesian methods, game theory, intrusion detection, security, stochastic systems.**

## I. INTRODUCTION

SOCIETAL monetary loss from cyber crime is estimated to be about a thousand billion USD per year presently, and even worse, a rising trend can be observed [1]. Another trend is that not only information systems but also control systems, which are typically governed by physical laws, are exposed to cyber threats as demonstrated by recent incidents [2]–[5]. *Deception* is a key notion to predict the consequence of incidents. Rational attackers take deceptive strategies, i.e., the attacker tries to conceal her existence and even mislead the defender into taking inappropriate decisions. An example of

H. Sasahara is with the Department of Systems and Control Engineering, Tokyo Institute of Technology, Tokyo, 152-8552 Japan e-mail: sasahara@sc.e.titech.ac.jp.

H. Sandberg is with the Division of Decision and Control Systems, KTH Royal Institute of Technology, Stockholm, SE-100 44 Sweden e-mail: hsan@kth.se.

deception is replay attacks, which hijacks sensors of the plant, eavesdrops the nominal data transmitted when the system is operated under normal conditions, and replays the observed nominal data during the execution of another damaging attack. A replay attack was executed in the Stuxnet incident, and it was an essential factor leading to serious damage in the targeted plant [6]. The incident suggests that prevention of deception is a fundamental requirement for secure system design.

Assuming the situation where an attacker might intrude into a control system where a defense mechanism is implemented, this paper addresses the following question: *Can model knowledge guide the defender to appropriate decisions against attacker's deceptive strategies?* Specifically, we consider the case where the stochastic model of the control system, the vulnerability to be exploited, and the objective of the attacker are known. The setting naturally leads to *Bayesian defense mechanisms*, which monitor the system's behavior and form a belief on the existence of the attacker using the model. If the system's behavior is inconsistent with the nominal one, the belief increases owing to Bayes' rule. When the belief is strong enough, the Bayesian defense mechanism proactively carries out a proper reaction. On the other hand, we also suppose a powerful attacker who knows the model and the defense scheme to be implemented. The attacker aims at achieving her objective while avoiding being detected by deceiving the defender.

For mathematical analysis, we formulate the decision making as a dynamic game with incomplete information. More specifically, we refer to the game as a stochastic signaling game, because it is a stochastic game [7] in the sense that the system's dynamics is given as a Markov decision process (MDP) governed by two players and it is also a signaling game [8] in the sense that one player's type is unknown to the opponent. In this game, the attacker strategically chooses harmful actions while avoiding being detected, while the defender, namely, the Bayesian defense mechanism, chooses appropriate counteractions according to her belief.

Based on the game-theoretic formulation, we find that *model knowledge can always lead the defender to appropriate decisions in an asymptotic sense* as long as the system's dynamics admits no stealthy attacks. More specifically, there are only two possible cases: one is that the defender asymptotically forms a firm belief on the existence of an attacker and the other is that the attacker takes harmless actions after finite time such that the system converges to nominal behavior. This finding leads to the conclusion that the Bayesian defense mechanism

guarantees the system to be secure in an asymptotic manner.

The analysis means that the defender always wins in an asymptotic manner when the stochastic model of the system is available and *the vulnerability exploited for the intrusion is known and modeled.* However, in practice, it is hard to be aware of all possible vulnerabilities in advance. As an application of the finding above, we consider *defensive deception using bluffing* that utilizes asymmetric recognitions between the attacker and the defender. Specifically, we suppose that, the defender is unaware of the exploited vulnerability but the attacker is unaware of the defender's unawareness. If the state of the system does not possess any information about the defender's recognition on the vulnerability, the attacker cannot identify whether the defender is aware of the vulnerability, or not. The result obtained in the former part suggests that the attacker may possibly withdraw if the defender's reactions affect only the attacker's utility without influence to the system's behavior. The difficulty of the analysis is that standard incomplete information games, which assume common prior, cannot describe this situation. The common prior implicitly assumes that the attacker is aware of the defender's unawareness. To overcome the difficulty, we employ the Mertens-Zamir model, which can represent incomplete information games without common prior assumption, using the notion of belief hierarchy [9], [10]. Based on this setting, we show, in a formal manner, that the defensive deception effectively works when the attacker strongly believes that the defender is aware of the vulnerability.

### Related Work

Model-based security analysis helps the system designer to prioritize security investments [11]. Attack graphs [12] and attack trees [13] are basic models of vulnerabilities, attacks, and consequences. Incorporating defensive actions into the graphical representation induces defense trees [14]. For dynamic models, attack countermeasure trees, partially observable MDP, and Bayesian network model have been used [15]–[17]. Those probabilistic models naturally lead to Bayesian defense mechanisms, such as Bayesian intrusion detection [18], [19], Bayesian intrusion response [20], and Bayesian security risk management [21]. Meanwhile, the model of the dynamical system to be protected is also used for control system security [22], [23]. For example, identifying existence of stealthy attacks and removing the vulnerability require the dynamical model [24], [25], and attack detection performance can be enhanced by model knowledge [26]. Our Bayesian defense mechanisms can be interpreted as a generalization of those approaches. This work reveals a fundamental property of such commonly used model-based defense schemes.

Game theory is a standard approach to modeling the decision making in cyber security, where there inevitably arises a need to address strategic interactions between the attacker and the defender [27], [28]. In particular, games with incomplete information play a crucial role in deceptive situations [29]–[32]. The modeling in this study follows the signaling game framework in [33], [34]. Our main concern is especially on

asymptotic phenomena in the dynamic deception and effectiveness of model knowledge.

Our finding is based on analysis of an asymptotic behavior of Bayesian inference. The convergence property of Bayesian inference on the true parameter is referred to as Bayesian consistency, which has been investigated mainly in the context of statistics [35], [36]. However, those existing results are basically applicable only to independent and identically distributed (i.i.d.) samples because the discussion mostly relies on the strong law of large numbers (SLLN). Although there is an extension to Markov chains [37], the observable variable in our work is not Markov. Indeed, sophisticated attackers can choose strategies such that the states at all steps are correlated with the entire previous trajectory. Thus, existing results for Bayesian consistency cannot be applied to our problem in a straightforward manner.

Preliminary versions of this work have been presented in [38], [39], but they made the claim of Theorem 2 as an assumption rather than proving it. Moreover, they did not include rigorous proofs of the claims in Section III and analysis of the bluffing proposed in Section IV.
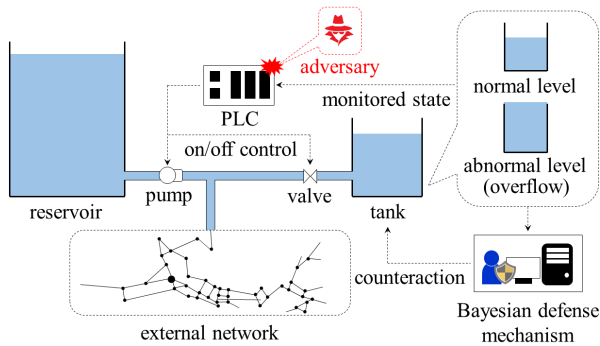
### Organization and Preliminaries

In Section II, we present a motivating example of water supply networks, and subsequently, formulate the decision making as a stochastic signaling game. Section III analyzes the consequence of the formulated game and shows that Bayesian defense mechanisms can achieve asymptotic security of the system to be protected. In Section IV, we analyze a defensive deception that utilizes asymmetric recognition as an application of the finding of Section III. The game of interest is reformulated using the Mertens-Zamir model. It is shown that the attacker possibly stops the execution even if the defender is unaware of the exploited vulnerabilities, as long as the defender's belief is concealed. Section V verifies the theoretical results through numerical simulation. Finally, Section VI concludes and summarizes the paper.

Let $\mathbb{N}$, $\mathbb{Z}_+$, and $\mathbb{R}$ be the sets of natural numbers, non-negative integers, and real numbers, respectively. The $k$-ary Cartesian power of the set $\mathcal{X}$ is denoted by $\mathcal{X}^k$. The tuple $(x_0, \ldots, x_k)$ is denoted by $x_{0:k}$. The cardinality of a set $\mathcal{X}$ is denoted by $|\mathcal{X}|$. For a set $\mathcal{X}$, the Kronecker delta denoted by $\delta : \mathcal{X} \times \mathcal{X} \to \{0, 1\}$ is defined by $\delta(x, y) = 1$ if $x = y$ and $\delta(x, y) = 0$ otherwise. The $\sigma$-algebra generated by a random variable $X$ is denoted by $\sigma(X)$. For a sequence of events $E_k$ for $k \in \mathbb{N}$, the supremum set $\cap_{N=1}^{\infty} \cup_{k=N}^{\infty} E_k$, namely, the event where $E_k$ occurs infinitely often, is denoted by $\{E_k \text{ i.o.}\}$. Jensen's inequality, which is often applied in this paper, is given as follows: For a real convex function $\varphi$ and a finite set $\mathcal{X}$, the inequality

$$\sum_{x \in \mathcal{X}} p(x)\varphi(a(x)) \geq \varphi \left( \sum_{x \in \mathcal{X}} p(x)a(x) \right) \qquad (1)$$

holds where $a : \mathcal{X} \to \mathbb{R}$ and $p : \mathcal{X} \to [0, 1]$ that satisfies the equation $\sum_{x \in \mathcal{X}} p(x) = 1$. The inequality is reversed if $\varphi$ is concave. The generalized Borel-Cantelli's second lemma is given as follows [40, Theorem 4.3.4]: Let $\mathcal{F}_k$ for $k \in \mathbb{Z}_+$ be a filtration of a probability space $(\Omega, \mathcal{F}, \mathbb{P})$ with $\mathcal{F}_0 := \{\emptyset, \Omega\}$

Fig. 1. Motivating example: water tank system connected to a reservoir within a water distribution network. The programmable logic controller (PLC) transmits on/off control signals to the pump and the valve monitoring the state, the water level of the tank. In the scenario, an adversarial software possibly intrudes into the PLC and then the infected PLC tries to cause overflow by sending inappropriate control signals without being detected. A Bayesian defense mechanism, which utilizes the data of the monitored state and forms her belief on existence of the attacker based on the system model, is also equipped to deal with the attack.

and let $E_k$ for $k \in \mathbb{Z}_+$ be a sequence of events with $E_k \in \mathcal{F}_{k+1}$. Then

$$\{E_k \text{ i.o.}\} = \{\omega \in \Omega : \sum_{k=0}^{\infty} \mathbb{P}(E_k|\mathcal{F}_k)(\omega) = \infty\}. \quad (2)$$

The appendix contains the proofs of the claims made in the paper.

## II. MODELING USING STOCHASTIC SIGNALING GAMES

### A. Motivating Example

As a motivating example, we consider water distribution networks (WDNs), which supply drinking water of suitable quality to customers. Because of their indispensability to our life, WDNs are an attractive target for adversaries and expose their architecture to cyber-physical attacks [41]. In particular, we treat the water tank system illustrated by Fig. 1, where a tank is connected to a reservoir within a WDN. The amount of the water in the tank varies due to usage for drinking and flow between the external network. Thus the tank system is needed to be properly controlled through actuation of the pump and the valve to keep the water amount within a desired range [42]. A programmable logic controller (PLC) transmits on/off control signals to the pump and the valve monitoring the state, namely, the water level of the tank. The dynamics is modeled as a MDP, where the state space and the action space are given by quantized water levels and finite control actions. Interaction to the external network is modeled as the randomness in the process.

We here suppose an attack scenario considered in [43]. The adversary succeeds to hijack the PLC and can directly manipulate its control logic. Such an intrusion can be carried out by stealthy and evasive maneuvers in advanced persistent threats [44]. The objective of the attack is to damage the system by causing water overflow through inappropriate control signals without being detected. To deal with this attack, we consider a Bayesian defense mechanism, which utilizes the data of the monitored state and forms her belief on existence of the attacker based on the system model. The Bayesian defense

mechanism chooses a proper reaction by identifying if the system is under attack through an observation of the state. If the system's behavior is highly suspicious, for example, the defense mechanism takes an aggressive reaction such as log analysis, dispatch of operators, or emergency shutdown.

The defender's belief on the existence of an attacker plays a key role to analyze the consequence of the threat. When the attacker naively executes an attack, the system's behavior becomes different from the one of the normal operation and accordingly the belief increases. On the other hand, if the attacker chooses sophisticated attacks that deceive the defender, the belief may decrease. Our main interest in this study is to investigate the defense capability achieved by the Bayesian defense mechanism.

### B. Modeling using Stochastic Signaling Game

We introduce the general description based on dynamic games with incomplete information. In particular, we refer to the game as a stochastic signaling game where the system's dynamics is given as an MDP and the type of a player is unknown to the opponent.
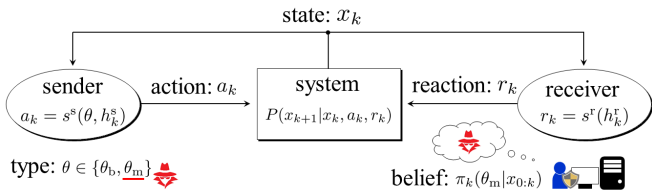
The system to be protected with a Bayesian defense mechanism is depicted in Fig. 2. The system is modeled by a finite MDP governed by two players as in standard stochastic games. Formally, the MDP considered in this paper is given by the tuple $\mathcal{M} := (\mathcal{X}, \mathcal{A}, \mathcal{R}, P, P_0)$ where $\mathcal{X}$ is a finite state space, $\mathcal{A}$ and $\mathcal{R}$ are finite action spaces, $P : \mathcal{X} \times \mathcal{X} \times \mathcal{A} \times \mathcal{R} \to [0, 1]$ is a transition probability, and $P_0 : \mathcal{X} \to [0, 1]$ is the probability distribution of the initial state. The state at the $k$th step is denoted by $x_k \in \mathcal{X}$. There is an agent who can alter the system through an *action* $a_k \in \mathcal{A}$ for $k \in \mathbb{Z}_+$. We refer to the agent as *sender* as in standard signaling games. Based on the measured output, the Bayesian defense mechanism, called a *receiver*, chooses an action $r_k \in \mathcal{R}$ at each time step. We henceforth refer to $r_k$ as a *reaction* for emphasizing that $r_k$ denotes a counteraction against potentially malicious attacks. The system dynamics is given by $P$, where the transition probability from $x$ to $x'$ with $a$ and $r$ is denoted by $P(x'|x, a, r)$. To eliminate the possibility of trivial stealthy attacks, we assume that the system's behavior varies in a stochastic sense when different actions are taken.

**Assumption 1** For any $x \in \mathcal{X}$ and $r \in \mathcal{R}$, there exists $x' \in \mathcal{X}$ such that

$$P(x'|x, a, r) \neq P(x'|x, a', r) \quad (3)$$

for different actions $a \neq a'$.

Next, we determine the class of the decision rules. Let $\theta \in \Theta$ denote the *type* of the sender. For simplicity, the type is assumed to be binary, i.e., $\Theta = \{\theta_\mathrm{b}, \theta_\mathrm{m}\}$, where $\theta_\mathrm{b}$ and $\theta_\mathrm{m}$ correspond to benign and malicious senders, respectively. The types $\theta_\mathrm{b}$ and $\theta_\mathrm{m}$ describe the situations where there does not and does exist an adversary, respectively. The true type $\theta$ is known to the sender, but unknown to the receiver. Let $\bar{s}^\mathrm{s} := (\bar{s}_k^\mathrm{s})_{k \in \mathbb{Z}_+}$ and $\bar{s}^\mathrm{r} := (\bar{s}_k^\mathrm{r})_{k \in \mathbb{Z}_+}$ denote the sender's and receiver's pure *strategy*, respectively. It is assumed that the receiver's available information about the sender type is only

Fig. 2. Block diagram of the system to be protected using the Bayesian defense mechanism. The system is governed by actions and reactions, which are decided by the sender and the receiver, respectively. The sender type $\theta_{\mathrm{b}}$ means that the system is normally operated. The other type $\theta_{\mathrm{m}}$ means that there exists an attacker who executes malicious actions. The receiver is the Bayesian defense mechanism that forms her belief on the existence of an attacker utilizing the measured data and chooses reactions based on the belief.

the state, i.e., she cannot observe her instantaneous utility, defined below, nor the sender's action. Similarly, it is assumed that the sender can observe only the state and her action. The strategies at the $k$th step with the available information are given by $\bar{s}_k^{\mathrm{s}} : \Theta \times \mathcal{H}_k^{\mathrm{s}} \to \mathcal{A}$, and $\bar{s}_k^{\mathrm{r}} : \mathcal{H}_k^{\mathrm{r}} \to \mathcal{R}$ where $h_k^{\mathrm{s}} \in \mathcal{H}_k^{\mathrm{s}}$, and $h_k^{\mathrm{r}} \in \mathcal{H}_k^{\mathrm{r}}$ are histories at the $k$th step given by $h_k^{\mathrm{s}} = (x_{0:k}, a_{0:k-1})$ and $h_k^{\mathrm{r}} = (x_{0:k}, r_{0:k-1})$. Note that the resulting state trajectory is not Markov since the strategies depend on the entire history. Because we consider pure strategies, it suffices to consider the state-history dependent strategies $s_k^{\mathrm{s}} : \Theta \times \mathcal{X}^{k+1} \to \mathcal{A}$ and $s_k^{\mathrm{r}} : \mathcal{X}^{k+1} \to \mathcal{R}$, recursively defined by

$$
\begin{aligned}
s_k^{\mathrm{s}}(\theta, x_{0:k}) &:= \bar{s}_k^{\mathrm{s}}(\theta, x_{0:k}, s_{0:k-1}^{\mathrm{s}}(x_{0:k-1})), \\
s_k^{\mathrm{r}}(x_{0:k}) &:= \bar{s}_k^{\mathrm{r}}(x_{0:k}, s_{0:k-1}^{\mathrm{r}}(x_{0:k-1})).
\end{aligned}
\tag{4}
$$

The strategy profile is denoted by $s := (s^{\mathrm{s}}, s^{\mathrm{r}})$. The sender's and receiver's admissible strategy sets are denoted by $\mathcal{S}^{\mathrm{s}}$ and $\mathcal{S}^{\mathrm{r}}$, respectively. The set of admissible strategy profiles is denoted by $\mathcal{S} := \mathcal{S}^{\mathrm{s}} \times \mathcal{S}^{\mathrm{r}}$. Note that, although we do not specify $\mathcal{S}$ here, it can be taken to be any set of state-history dependent strategies. While we consider a general strategy set in Sec. III, we impose a constraint on $\mathcal{S}$ in Sec. IV.

Once a strategy profile is fixed, the stochastic property of the system is induced. Construct the canonical measurable space $(\Omega, \mathcal{F})$ of the MDP with the sender type where $\Omega := \Theta \times \Pi_{k=0}^{\infty} (\mathcal{X} \times \mathcal{A} \times \mathcal{R})$ and $\mathcal{F}$ is its product $\sigma$-algebra [45, Chapter 2]. We denote $\omega = (\theta, (x_0, a_0, r_0), (x_1, a_1, r_1), \ldots) \in \Omega$. The random variables $\Theta, X_k, A_k$, and $R_k$ are defined on the measurable space $(\Omega, \mathcal{F})$ by the projections of $\omega$ such that $\Theta(\omega) := \theta, X_k(\omega) := x_k, A_k(\omega) := a_k, R_k(\omega) := r_k$. The probability measure on $(\Omega, \mathcal{F})$, induced by $s$, is denoted by $\mathbb{P}^s$, which satisfies

$$
\begin{cases}
\mathbb{P}^s(X_0 = x_0) = P_0(x_0), \\
\mathbb{P}^s(A_k = a_k | \Theta = \theta, X_{0:k} = x_{0:k}) = \delta(a_k, s_k^{\mathrm{s}}(\theta, x_{0:k})), \\
\mathbb{P}^s(R_k = r_k | X_{0:k} = x_{0:k}) = \delta(r_k, s_k^{\mathrm{r}}(x_{0:k})), \\
\mathbb{P}^s(X_{k+1} = x_{k+1} | X_{0:k} = x_{0:k}, A_k = a_k, R_k = r_k) \\
\quad = P(x_{k+1} | x_k, a_k, r_k), \\
\mathbb{P}^s(\Theta = \theta) = \pi_0(\theta)
\end{cases}
\tag{5}
$$

for any $k \in \mathbb{Z}_+$ with the initial distribution of the sender type $\pi_0 : \Theta \to [0,1]$. We denote the conditional probability $\mathbb{P}^s(\cdot | \Theta = \theta)$ by $\mathbb{P}_\theta^s$. To simplify the notation, we denote the

conditional probability mass function with type $\theta$ by

$$
p_\theta^s(x_{k+1} | x_{0:k}) := \mathbb{P}_\theta^s(X_{k+1} = x_{k+1} | X_0 = x_0, \ldots, X_k = x_k).
\tag{6}
$$

The expectation with respect to $\mathbb{P}^s$ is denoted by $\mathbb{E}^s$.

We introduce each player's *belief* on the uncertain variables next. The receiver's belief at the $k$th step is given by

$$
\begin{aligned}
&\pi_k^{\mathrm{r}}(\theta, a_{0:k-1} | x_{0:k}, r_{0:k-1}) \\
&:= \mathbb{P}^s(\Theta = \theta, A_{0:k-1} = a_{0:k-1} | X_{0:k} = x_{0:k}, R_{0:k-1} = r_{0:k-1})
\end{aligned}
\tag{7}
$$

for $k \in \mathbb{Z}_+$. The belief can be recursively computed by Bayes' rule

$$
\begin{aligned}
&\pi_{k+1}^{\mathrm{r}}(\theta, a_{0:k} | x_{0:k+1}, r_{0:k}) = \delta(a_k, s_k^{\mathrm{s}}(\theta, x_{0:k})) \\
&\times \frac{P(x_{k+1} | x_k, s_k^{\mathrm{s}}(\theta, x_{0:k}), r_k) \pi_k^{\mathrm{r}}(\theta, a_{0:k-1} | x_{0:k}, r_{0:k-1})}{\sum_{\phi \in \Theta} P(x_{k+1} | x_k, s_k^{\mathrm{s}}(\phi, x_{0:k}), r_k) \pi_k^{\mathrm{r}}(\phi, a_{0:k-1} | x_{0:k}, r_{0:k-1})}
\end{aligned}
\tag{8}
$$

when the denominator is nonzero. To simplify notation, we introduce the receiver's belief only of the sender type:

$$
\pi_k^{\mathrm{r}}(\theta | x_{0:k}) := \pi_k^{\mathrm{r}}(\theta, s_{0:k-1}^{\mathrm{s}}(\theta, x_{0:k-1}) | x_{0:k}, s_{0:k-1}^{\mathrm{r}}(x_{0:k-1})),
\tag{9}
$$

which follows Bayes' rule

$$
\begin{aligned}
&\pi_{k+1}^{\mathrm{r}}(\theta | x_{0:k+1}) \\
&= \frac{P(x_{k+1} | x_k, s_k^{\mathrm{s}}(\theta, x_{0:k}), s_k^{\mathrm{r}}(x_{0:k})) \pi_k^{\mathrm{r}}(\theta | x_{0:k})}{\sum_{\phi \in \Theta} P(x_{k+1} | x_k, s_k^{\mathrm{s}}(\phi, x_{0:k}), s_k^{\mathrm{r}}(x_{0:k})) \pi_k^{\mathrm{r}}(\phi | x_{0:k})}.
\end{aligned}
\tag{10}
$$

The sender's belief can similarly be defined and is denoted by $\pi_k^{\mathrm{s}}(r_{0:k-1} | \theta, x_{0:k}, a_{0:k-1})$.

In Sec. III, the initial beliefs are assumed to be known to both players, i.e., we make the common prior assumption. Since we consider pure strategies, $r_{0:k-1}$ is uniquely determined by $x_{0:k-1}$ once the strategy is fixed. Hence, the sender's belief does not appear explicitly in Sec. III. On the other hand, in Sec. IV, we consider the case where the initial belief is unknown to the sender, modeling the possibility of *bluffing*.

Let $U^{\mathrm{s}} : \Theta \times \mathcal{X} \times \mathcal{A} \times \mathcal{R} \to \mathbb{R}$ be the sender's instantaneous utility. For a given strategy profile $s \in \mathcal{S}$ and type $\theta \in \Theta$, the sender's expected average utility at the $k$th step with the horizon length $T$ is given by

$$
\begin{aligned}
&\bar{U}_{k,T}^{\mathrm{s}}(s_{k:k+T} | \theta, x_{0:k}) \\
&:= \mathbb{E}^s \left[ \frac{1}{T+1} \sum_{\tau=k}^{k+T} U^{\mathrm{s}}(\Theta, X_\tau, s_\tau^{\mathrm{s}}(\Theta, X_{0:\tau}), s_\tau^{\mathrm{r}}(X_{0:\tau})) \bigg| \theta, x_{0:k} \right].
\end{aligned}
\tag{11}
$$

Similarly, with the receiver's instantaneous utility given by $U^{\mathrm{r}} : \Theta \times \mathcal{X} \times \mathcal{A} \times \mathcal{R} \to \mathbb{R}$, the receiver's expected average utility at the $k$th step with the horizon length $T$ is given by

$$
\begin{aligned}
&\bar{U}_{k,T}^{\mathrm{r}}(s_{k:k+T} | x_{0:k}) \\
&:= \mathbb{E}^s \left[ \frac{1}{T+1} \sum_{\tau=k}^{k+T} U^{\mathrm{r}}(\Theta, X_\tau, s_\tau^{\mathrm{s}}(\Theta, X_{0:\tau}), s_\tau^{\mathrm{r}}(X_{0:\tau})) \bigg| x_{0:k} \right].
\end{aligned}
\tag{12}
$$

We denote the limits by $(\bar{U}_k^{\mathrm{s}}, \bar{U}_k^{\mathrm{r}}) := \lim_{T \to \infty} (\bar{U}_{k,T}^{\mathrm{s}}, \bar{U}_{k,T}^{\mathrm{r}})$ assuming they exist. Under this notation, the strategy profile $s = (s^{\mathrm{s}}, s^{\mathrm{r}})$ is said to be a *perfect Bayesian equilibrium* (PBE) if

$$
\begin{cases}
s_{k:\infty}^{\mathrm{s}} \in \mathrm{BR}_k^{\mathrm{s}}(s_{r:\infty}^{\mathrm{r}} | \theta, x_{0:k}), \quad \forall \theta \in \Theta, \\
s_{k:\infty}^{\mathrm{r}} \in \mathrm{BR}_k^{\mathrm{r}}(s_{k:\infty}^{\mathrm{s}} | x_{0:k})
\end{cases}
\tag{13}
$$

for any $k \in \mathbb{Z}_+$ and $x_{0:k} \in \mathcal{X}^{k+1}$ where $\mathrm{BR}_k^{\mathrm{s}}$ and $\mathrm{BR}_k^{\mathrm{r}}$ are best responses defined by

$$\mathrm{BR}_k^{\mathrm{s}}(s_{k:\infty}^{\mathrm{r}}|\theta, x_{0:k}) := \underset{\tilde{s}_{k:\infty}^{\mathrm{s}} \in \mathcal{S}_{k:\infty}^{\mathrm{s}}}{\arg\max} \ \bar{U}_k^{\mathrm{s}}((\tilde{s}_{k:\infty}^{\mathrm{s}}, s_{k:\infty}^{\mathrm{r}})|\theta, x_{0:k}),$$
$$\mathrm{BR}_k^{\mathrm{r}}(s_{k:\infty}^{\mathrm{s}}|x_{0:k}) := \underset{\tilde{s}_{k:\infty}^{\mathrm{r}} \in \mathcal{S}_{k:\infty}^{\mathrm{r}}}{\arg\max} \ \bar{U}_k^{\mathrm{r}}((s_{k:\infty}^{\mathrm{s}}, \tilde{s}_{k:\infty}^{\mathrm{r}})|x_{0:k}).$$
(14)

Note that, our analysis can be extended to the case of general objective functions rather than expected average utilities as long as the adversary with the utilities avoids being detected, which is formally stated in Definition 1 below.

We define the game formulated above by

$$\mathcal{G}_1 := (\mathcal{M}, \mathcal{S}, U, \Theta, \pi_0),$$
(15)

where the initial belief is common information. This game belongs to the class of incomplete, imperfect, and asymmetric information stochastic games. Owing to the existence of the type $\theta$, which is unknown to the receiver, the information is incomplete. Because the actions taken by each player are unobservable to the opponent, the information is imperfect and asymmetric. Although investigating existence and computing equilibria of the game are challenging, we discuss properties of equilibria on the premise that they exist and are given because our interest here lies in the consequences for the threat.

### III. ANALYSIS: ASYMPTOTIC SECURITY

In this section, we analyze asymptotic behaviors of beliefs and actions when the adversary avoids being detected. It is shown that the system is guaranteed to be secure in an asymptotic manner as long as the defender possesses an effective counteraction.

#### A. Belief's Asymptotic Behavior

The random variable of the belief on the type $\theta \in \Theta$ at the $k$th step $\pi_k^\theta : \Omega \to [0, 1]$ is given by

$$\pi_k^\theta(\omega) := \pi_k^{\mathrm{r}}(\theta|X_{0:k}(\omega)).$$
(16)

Recall that $\pi_k^\theta$ represents the defender's confidence on existence of an attacker. If the belief is low in spite of existence of malicious signals, this means that the Bayesian defense mechanism is deceived. Because we are interested in whether the Bayesian defense mechanism is permanently deceived, or not, we examine asymptotic behavior of the belief.

We first investigate increment of the belief sequence. The following lemma is key to our analysis.

**Lemma 1** Consider the game $\mathcal{G}_1$. The belief of the true type $\pi_k^\theta$ is a submartingale with respect to the probability $\mathbb{P}_\theta^s$ and the filtration $\sigma(X_{0:k})$ for any type $\theta$ and strategy profile $s$.

Lemma 1 roughly implies that the expectation of the belief on the true type is non-decreasing. As a direct conclusion of this lemma, the following theorem holds.

**Theorem 1** Consider the game $\mathcal{G}_1$. There exists an integrable random variable $\pi_\infty^\theta : \Omega \to [0, 1]$ such that

$$\lim_{k \to \infty} \pi_k^\theta = \pi_\infty^\theta \quad \mathbb{P}_\theta^s\text{-a.s.}$$
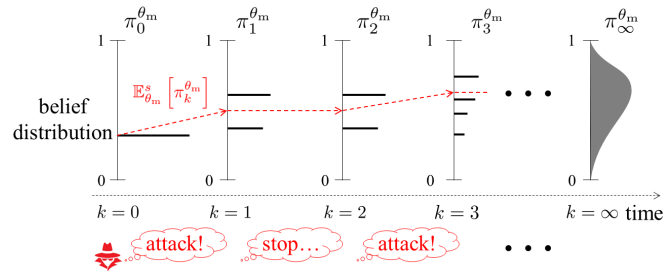(17)



Fig. 3. Distributions of the belief sequence when there exists an attacker. Lemma 1 and Theorem 1 claim that its expectation is non-decreasing over time and the belief has a limit. When the adversary stops the attack, the belief is invariant.

for any type $\theta$ and strategy profile $s$.

Theorem 1 implies that the belief has a limit even if an intermittent attack is executed. Fig. 3 depicts the distributions of the belief sequence when there exists an attacker. Owing to the model knowledge, if the adversary stops the attack at some time step then the belief is invariant, which is illustrated as the transition of the belief at $k = 1$ in Fig. 3. Moreover, the expectation of the belief is non-decreasing over time as claimed by Lemma 1. Thus, there exists a limit $\pi_\infty^{\theta_{\mathrm{m}}}$ as shown at the right of Fig. 3.

We next investigate the limit. An undesirable limit is $\pi_\infty^\theta = 0$, which means that *the defender is completely deceived.* We show that *this does not happen as long as the initial belief is nonzero.* The following lemma holds.

**Lemma 2** Consider the game $\mathcal{G}_1$. If $\pi_0^\theta > 0$ for any type $\theta$, then $\log(\pi_k^\theta)$ with any basis converges $\mathbb{P}_\theta^s$-almost surely to an integrable random variable as $k \to \infty$ for any type $\theta$ and strategy profile $s$.

Lemma 2 leads to the following theorem.

**Theorem 2** Consider the game $\mathcal{G}_1$. If $\pi_0^\theta > 0$ then

$$\pi_\infty^\theta > 0 \quad \mathbb{P}_\theta^s\text{-a.s.}$$
(18)

for any type $\theta$ and strategy profile $s$.

Theorem 2 implies that the complete deception described by $\pi_\infty^\theta = 0$ does not occur.

*Remark:* Theorems 1 and 2 can heuristically be justified from an information-theoretic perspective as follows. Suppose that the state sequence $x_{0:k}$ is observed at the $k$th step. Then the belief is given by

$$\pi_k(\theta|x_{0:k}) = \frac{\pi_0(\theta)}{\sum_{\phi \neq \theta} \frac{p_\phi^s(x_{0:k})}{p_\theta^s(x_{0:k})} \pi_0(\phi) + \pi_0(\theta)}$$
$$= \frac{\pi_0(\theta)}{\sum_{\phi \neq \theta} \exp(kS_k^\phi(x_{0:k})) \pi_0(\phi) + \pi_0(\theta)}$$
(19)

where $p_\theta^s(x_{0:k})$ and $p_\phi^s(x_{0:k})$ are the joint probability mass functions of $x_{0:k}$ with respect to $\mathbb{P}_\theta^s$ and $\mathbb{P}_\phi^s$, respectively, and

$$S_k^\phi(x_{0:k}) := \frac{1}{k} \sum_{i=1}^k \log \frac{p_\phi^s(x_i|x_{0:i-1})}{p_\theta^s(x_i|x_{0:i-1})}.$$
(20)

Assuming that $p_\theta^s(x_k|x_{0:k-1})$ approaches a stationary distribution $p_\theta^s(x)$ on $\mathcal{X}$ and SLLN can be applied, we have

$$\lim_{k\to\infty} S_k^\phi = \mathbb{E}_{x\sim p_\theta^s}\left[\log p_\phi^s(x)/p_\theta^s(x)\right] = -D_{\mathrm{KL}}(p_\theta^s||p_\phi^s) \quad (21)$$

where $D_{\mathrm{KL}}$ denotes the Kullback-Leibler divergence. Since $D_{\mathrm{KL}}$ is nonnegative for any pair of distributions, $S_k^\phi$ converges to a nonpositive number, which results in convergence of $\pi_k^\theta$. If $p_\theta^s \neq p_\phi^s$ for any $\phi \in \Theta\backslash\{\theta\}$, the limit of $S_k^\phi$ becomes negative, and hence $\lim_{k\to\infty}\exp(kS_k^\phi(x_{0:k})) = 0$, which leads to

$$\lim_{k\to\infty} \pi_k(\theta|x_{0:k}) = \frac{\pi_0(\theta)}{\displaystyle\sum_{\phi\neq\theta} \lim_{k\to\infty}\exp(kS_k^\phi(x_{0:k}))\pi_0(\phi) + \pi_0(\theta)}$$
$$= 1.$$
$$(22)$$

Thus, the belief of the true type converges to one. Such a convergence property of the Bayesian estimator on the true parameter, referred to as Bayesian consistency, has been investigated mainly in the context of statistics [35], [36]. In this sense, Theorems 1 and 2 can be regarded as another representation of Bayesian consistency. However, note again that this discussion is not a rigorous proof but a heuristic justification because the state is essentially non-i.i.d. and even non-ergodic in our game-theoretic formulation.

### B. Asymptotic Security

It has turned out that the belief has a positive limit. To clarify our interest, we define the notion of detection-averse utilities.

**Definition 1 (Detection-averse Utilities)** A pair $(U^{\mathrm{s}}, U^{\mathrm{r}})$ in the game $\mathcal{G}_1$ are detection-averse utilities when

$$\pi_\infty^{\theta_{\mathrm{m}}} < 1 \quad \mathbb{P}_{\theta_{\mathrm{m}}}^s-\text{a.s.} \quad (23)$$

for any PBE $s$.

Definition 1 characterizes utilities where the malicious sender avoids having the defender form a firm belief on the existence of an attacker. An example of detection-averse utilities is given in Appendix I. Naturally, strategies reasonable for the attacker should be detection-averse as long as the defender possesses an effective counteraction. If the utilities of interest are not detection-averse, this means that the defense mechanism cannot cope with the attack because the attacker is not afraid to reveal herself. For protecting such systems, appropriate counteractions should be implemented beforehand.

Suppose that there is an effective countermeasure, and hence the utilities are detection-averse. A simple malicious sender's strategy that satisfies (23) is to imitate the benign sender's strategy after a finite number of time steps. We give a formal definition of such strategies.

**Definition 2 (Asymptotically Benign Strategy)** A strategy profile $s$ in the game $\mathcal{G}_1$ is asymptotically benign when

$$\lim_{k\to\infty} \delta\left(A_k^{\theta_{\mathrm{m}}}, A_k^{\theta_{\mathrm{b}}}\right) = 1 \quad \mathbb{P}_{\theta_{\mathrm{m}}}^s-\text{a.s.} \quad (24)$$

where $A_k^\theta$ is the action taken by the sender with the type $\theta$ defined by $A_k^\theta := s_k^s(\theta, X_{0:k})$.

The objective of this subsection is to show that Bayesian defense mechanisms can restrict all reasonable strategies to be asymptotically benign as long as an effective countermeasure is implemented.

As a preparation for proving our main claim, we investigate the asymptotic behavior of state transition. From Theorems 1 and 2, we can expect that the state eventually loses information on the type, which is justified by the following lemma.

**Lemma 3** Consider the game $\mathcal{G}_1$ with detection-averse utilities. If $\pi_0^{\theta_{\mathrm{m}}} > 0$, then

$$\lim_{k\to\infty} \left|p_{\theta_{\mathrm{m}}}^s(X_{k+1}|X_{0:k}) - p_{\theta_{\mathrm{b}}}^s(X_{k+1}|X_{0:k})\right| = 0 \quad \mathbb{P}_{\theta_{\mathrm{m}}}^s-\text{a.s.}$$
$$(25)$$

for any PBE $s$.

Under Assumption 1, which eliminates the possibility of stealthy attacks, Lemma 3 implies that the actions themselves must be identical. This fact yields the following theorem, one of the main results in this paper.

**Theorem 3** Consider the game $\mathcal{G}_1$ with detection-averse utilities. Let Assumption 1 hold and assume $\pi_0^{\theta_{\mathrm{m}}} > 0$. Then, every PBE of $\mathcal{G}_1$ is asymptotically benign.

Theorem 3 implies that the malicious sender's action converges to the benign action. Equivalently, an attacker necessarily behaves as a benign sender after a finite number time steps. Therefore, the system is guaranteed to be secure in an asymptotic manner, i.e., Bayesian defense mechanisms *can* prevent deception in an asymptotic sense. This result indicates the powerful defense capability achieved by model knowledge.

## IV. APPLICATION: ANALYSIS OF DEFENSIVE DECEPTION UTILIZING ASYMMETRIC RECOGNITION

### A. Idea of Defensive Deception using Bluffing

The result in Section III claims that the defender, namely, the Bayesian defense mechanism, always wins in an asymptotic manner when the stochastic model of the system is available and *the vulnerability to be exploited for intrusion is known and modeled*. The latter condition is quantitatively described by the condition $\pi_0(\theta_{\mathrm{m}}) > 0$. Although the derived result proves a quite powerful defense capability, it is also true that it is almost impossible to be aware of all possible vulnerabilities in advance. Moreover, it is also challenging to implement effective countermeasures for all scenarios and to compute the equilibrium of the dynamic game.

In this section, as an application of the finding in the previous section, we consider *defensive deception* using bluffing that utilizes asymmetric recognitions between the attacker and the defender. Suppose that an attacker exploits a vulnerability of which the defender is unaware but the attacker is unaware of the defender's unawareness. Then their recognition becomes asymmetric in the sense that the attacker does not correctly recognize the defender's recognition of the vulnerability. This
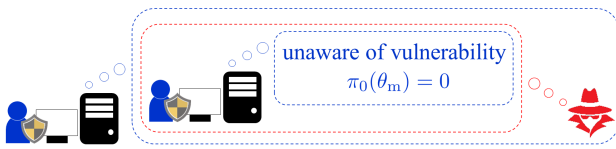
Fig. 4. The attacker's belief of the defender's belief with symmetric recognition, which is the case of the game $\mathcal{G}_1$. The attacker is aware of the fact that the defender is unaware of the vulnerability. Moreover, the defender is aware of the attacker's awareness.
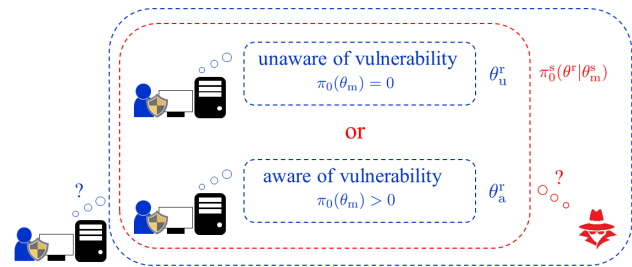


Fig. 5. The attacker's belief of the defender's belief with asymmetric recognition. Because the defender's true belief is unknown to the attacker, the attacker forms a belief on both cases that the defender is aware or unaware of the vulnerability. Moreover, the defender forms a belief on the attacker's belief. This process induces the notion of belief hierarchy.

situation naturally arises in practice because the defender's recognition is private information. By utilizing the asymmetric recognition, the defender can possibly deceive the attacker such that the attacker believes that the defender might be aware of the vulnerability and carrying out effective counteractions. Specifically, we consider the bluffing strategies where the system's state does not possess information about the defender's belief. For instance, if the defender chooses the reactions that affect only the players' utilities without influence to the system, the state is independent of the reaction. By concealing the defender's unawareness, the defender's recognition, which is quantified by her belief, is completely unknown to the attacker over time.

The defensive deception is possibly able to force the attacker to withdraw even if the defender is actually unaware of the exploited vulnerability. For instance, consider the example in Sec. II-A and suppose that emergency shutdown of the system can be carried out by the defender. Suppose also that the attacker wants to keep administrative privileges of the PLC. In this case, the attacker may rationally terminate her evasive maneuvers after a finite number of time steps due to the risk of sudden shutdown. The objective of this section is to show that the hypothesis is true in a formal manner.

### B. Reformulation using Type Structure

The situation of interest in this section is that the defender is unaware of the vulnerability to be exploited but the attacker is *not* necessarily aware of this unawareness. To address the uncertainty on defender's recognition, the attacker forms her belief on the defender's belief. Fig. 4 illustrates the attacker's belief on the defender's belief with the common prior assumption, i.e., the initial defender's belief is known to the attacker, which has been made in the previous section. In this case, the attacker has a firm belief that the defender is unaware of the vulnerability. On the other hand, Fig. 5 illustrates the attacker's belief without the common prior assumption. Then the attacker's belief is no longer firm as depicted by the figure. In addition, because of the lack of the common prior assumption, the defender also forms another belief on the attacker's belief on the defender's belief on the existence of an attacker. This procedure repeats indefinitely and induces infinitely many beliefs.

The notion of *belief hierarchy* has been proposed to handle the infinitely many beliefs [9], [10], [46]. A belief hierarchy is formed as follows. Let $\Delta(\cdot)$ denote the set of probability measures over a set. The first-order initial belief is given as $\pi_0^1 \in \Delta(\Theta)$, which describes the defender's initial belief on

existence of the attacker. The second-order initial belief is given as $\pi_0^2 \in \Delta(\Delta(\Theta))$, which describes the attacker's initial belief on the defender's first-order belief. In a similar manner, the belief at any level is given, and the tuple of beliefs at all levels is referred to as a belief hierarchy.

To handle belief hierarchies, the *Mertens-Zamir model* has been introduced [9], [10], [46]. The model considers *type structure*, in which a belief hierarchy is embedded. A type structure consists of players, sets of types, and initial beliefs. In particular, a type structure for our situation of interest can be given by

$$\mathcal{T} = ((\mathrm{s}, \mathrm{r}), (\Theta^{\mathrm{s}}, \Theta^{\mathrm{r}}), (\pi_0^{\mathrm{s}}, \pi_0^{\mathrm{r}})) \qquad (26)$$

where $(\mathrm{s}, \mathrm{r})$ represents the sender and the receiver, $\Theta^{\mathrm{s}}$ and $\Theta^{\mathrm{r}}$ represent the sets of player types, and $\pi_0^{\mathrm{s}} : \Theta^{\mathrm{r}} \times \Theta^{\mathrm{s}} \to [0, 1]$ and $\pi_0^{\mathrm{r}} : \Theta^{\mathrm{s}} \times \Theta^{\mathrm{r}} \to [0, 1]$ represent the initial beliefs. The value $\pi_0^{\mathrm{s}}(\theta^{\mathrm{r}}|\theta^{\mathrm{s}})$ denotes the sender's initial belief of the receiver type $\theta^{\mathrm{r}}$ when the sender type is $\theta^{\mathrm{s}}$, and $\pi_0^{\mathrm{r}}(\theta^{\mathrm{s}}|\theta^{\mathrm{r}})$ denotes the corresponding receiver's initial belief. The first-order initial belief is given by $\pi_0^1(\theta^{\mathrm{s}}) = \pi_0^{\mathrm{r}}(\theta^{\mathrm{s}}|\theta^{\mathrm{r}})$ for the true receiver type $\theta^{\mathrm{r}} \in \Theta^{\mathrm{r}}$, and the second-order initial belief is given by $\pi_0^2(\pi^{\mathrm{r}}(\cdot|\theta^{\mathrm{r}})|\theta^{\mathrm{s}}) = \pi_0^{\mathrm{s}}(\theta^{\mathrm{r}}|\theta^{\mathrm{s}})$ for the true sender type $\theta^{\mathrm{s}} \in \Theta^{\mathrm{s}}$. By repeating it, the belief at any level of the belief hierarchy can be derived from the type structure. Importantly, for any reasonable belief hierarchy there exists a type structure that can generate the belief hierarchy of interest. For a formal discussion, see [9], [10], [46].

We model the situation of interest by using the binary type sets:

$$\Theta^{\mathrm{s}} = \{\theta_{\mathrm{b}}^{\mathrm{s}}, \theta_{\mathrm{m}}^{\mathrm{s}}\}, \quad \Theta^{\mathrm{r}} = \{\theta_{\mathrm{u}}^{\mathrm{r}}, \theta_{\mathrm{a}}^{\mathrm{r}}\}. \qquad (27)$$

While $\theta_{\mathrm{b}}^{\mathrm{s}}$ and $\theta_{\mathrm{m}}^{\mathrm{s}}$ represent benign and malicious senders, respectively, $\theta_{\mathrm{u}}^{\mathrm{r}}$ and $\theta_{\mathrm{a}}^{\mathrm{r}}$ represent receivers being unaware and aware of the vulnerability, respectively. The receiver's initial beliefs are set to

$$\pi_0^{\mathrm{r}}(\theta_{\mathrm{b}}^{\mathrm{s}}|\theta_{\mathrm{u}}^{\mathrm{r}}) = 1, \quad \pi_0^{\mathrm{r}}(\theta_{\mathrm{m}}^{\mathrm{s}}|\theta_{\mathrm{u}}^{\mathrm{r}}) = 0 \qquad (28)$$

and

$$\pi_0^{\mathrm{r}}(\theta_{\mathrm{b}}^{\mathrm{s}}|\theta_{\mathrm{a}}^{\mathrm{r}}) = \alpha, \quad \pi_0^{\mathrm{r}}(\theta_{\mathrm{m}}^{\mathrm{s}}|\theta_{\mathrm{a}}^{\mathrm{r}}) = 1 - \alpha \qquad (29)$$

with $\alpha \in [0, 1)$. The initial beliefs mean that, the receiver $\theta_{\mathrm{u}}^{\mathrm{r}}$ is unaware of the vulnerability and firmly believes that the

This article has been accepted for publication in IEEE Transactions on Automatic Control. This is the author's version which has not been fully edited and content may change prior to final publication. Citation information: DOI 10.1109/TAC.2023.3340978

8      IEEE TRANSACTIONS ON AUTOMATIC CONTROL, VOL. XX, NO. XX, XXXX 2023

TABLE I
INITIAL BELIEFS ON OPPONENT TYPE

| $\pi_0^{\mathrm{r}}(\cdot\|\theta_{\mathrm{u}}^{\mathrm{r}})$ | $\theta_{\mathrm{b}}^{\mathrm{s}}$ | $\theta_{\mathrm{m}}^{\mathrm{s}}$ |
|---|---|---|
| $\pi_0^{\mathrm{r}}(\cdot\|\theta_{\mathrm{u}}^{\mathrm{r}})$ | 1 | 0 |
| $\pi_0^{\mathrm{r}}(\cdot\|\theta_{\mathrm{a}}^{\mathrm{r}})$ | $\alpha$ | $1-\alpha$ |

| $\pi_0^{\mathrm{s}}(\cdot\|\theta_{\mathrm{b}}^{\mathrm{s}})$ | $\theta_{\mathrm{u}}^{\mathrm{r}}$ | $\theta_{\mathrm{a}}^{\mathrm{r}}$ |
|---|---|---|
| $\pi_0^{\mathrm{s}}(\cdot\|\theta_{\mathrm{b}}^{\mathrm{s}})$ | 1 | 0 |
| $\pi_0^{\mathrm{s}}(\cdot\|\theta_{\mathrm{m}}^{\mathrm{s}})$ | $\beta$ | $1-\beta$ |

system is normally operated, while the receiver $\theta_{\mathrm{a}}^{\mathrm{r}}$ is aware of the vulnerability and suspects existence of an attacker with probability $1-\alpha$. The sender's initial beliefs are assumed to be given by

$$\pi_0^{\mathrm{s}}(\theta_{\mathrm{u}}^{\mathrm{r}}|\theta_{\mathrm{b}}^{\mathrm{s}})=1, \quad \pi_0^{\mathrm{s}}(\theta_{\mathrm{a}}^{\mathrm{r}}|\theta_{\mathrm{b}}^{\mathrm{s}})=0, \tag{30}$$

and

$$\pi_0^{\mathrm{s}}(\theta_{\mathrm{u}}^{\mathrm{r}}|\theta_{\mathrm{m}}^{\mathrm{s}})=\beta, \quad \pi_0^{\mathrm{s}}(\theta_{\mathrm{a}}^{\mathrm{r}}|\theta_{\mathrm{m}}^{\mathrm{s}})=1-\beta \tag{31}$$

with $\beta \in [0,1]$. The malicious sender does not know the true receiver type, i.e., whether the sender is aware of the vulnerability or not. The given initial beliefs are summarized in Table I.

In accordance with the introduction of the type structure, the definition of strategies and the solution concept are needed to be slightly modified. The contrasting ingredients of the game with symmetric recognition and the one with asymmetric recognition are listed in Table II, where those with asymmetric recognition can analogically be defined. The conditional probability $\mathbb{P}^s(\cdot|\Theta^{\mathrm{s}} = \theta^{\mathrm{s}}, \Theta^{\mathrm{r}} = \theta^{\mathrm{r}})$, which is the probability measure induced by $s^{\mathrm{s}}(\theta^{\mathrm{s}}, \cdot)$ and $s^{\mathrm{r}}(\theta^{\mathrm{r}}, \cdot)$, is denoted by $\mathbb{P}^s_{\theta^{\mathrm{s}}, \theta^{\mathrm{r}}}$. The sender's expected average utility at the $k$th step with the horizon length $T$ is given by

$$\bar{U}_{k,T}^{\mathrm{s}}(s_{k:k+T}|\theta^{\mathrm{s}}, x_{0:k}) := \frac{1}{T+1}$$
$$\times \mathbb{E}^s\left[ \sum_{\tau=k}^{k+T} U^{\mathrm{s}}(\Theta^{\mathrm{s}}, X_\tau, s_\tau^{\mathrm{s}}(\Theta^{\mathrm{s}}, X_{0:\tau}), s_\tau^{\mathrm{r}}(\Theta^{\mathrm{r}}, X_{0:\tau})) \middle| \theta^{\mathrm{s}}, x_{0:k} \right]. \tag{32}$$

The receiver's expected average utility at the $k$th step with the horizon length $T$ is given by

$$\bar{U}_{k,T}^{\mathrm{r}}(s_{k:k+T}|\theta^{\mathrm{r}}, x_{0:k}) := \frac{1}{T+1}$$
$$\times \mathbb{E}^s\left[ \sum_{\tau=k}^{k+T} U^{\mathrm{r}}(\Theta^{\mathrm{s}}, X_\tau, s_\tau^{\mathrm{s}}(\Theta^{\mathrm{s}}, X_{0:\tau}), s_\tau^{\mathrm{r}}(\Theta^{\mathrm{r}}, X_{0:\tau})) \middle| \theta^{\mathrm{r}}, x_{0:k} \right]. \tag{33}$$

A strategy $s$ is said to be a PBE when the limit of the utilities $(\bar{U}_k^{\mathrm{s}}, \bar{U}_k^{\mathrm{r}})$ satisfies

$$\begin{cases} s_{r:\infty}^{\mathrm{s}} \in \mathrm{BR}_k^{\mathrm{s}}(s_{r:\infty}^{\mathrm{r}}|\theta^{\mathrm{s}}, x_{0:k}), & \forall \theta^{\mathrm{s}} \in \Theta^{\mathrm{s}}, \\ s_{k:\infty}^{\mathrm{r}} \in \mathrm{BR}_k^{\mathrm{r}}(s_{k:\infty}^{\mathrm{s}}|\theta^{\mathrm{r}}, x_{0:k}), & \forall \theta^{\mathrm{r}} \in \Theta^{\mathrm{r}}, \end{cases} \tag{34}$$

for any $k \in \mathbb{Z}_+$ and $x_{0:k} \in \mathcal{X}^{k+1}$ where

$$\mathrm{BR}_k^{\mathrm{s}}(s_{k:\infty}^{\mathrm{r}}|\theta^{\mathrm{s}}, x_{0:k}) := \operatorname*{arg\,max}_{\tilde{s}_{k:\infty}^{\mathrm{s}} \in \mathcal{S}_{k:\infty}^{\mathrm{s}}} \bar{U}_k^{\mathrm{s}}((\tilde{s}_{k:\infty}^{\mathrm{s}}, s_{k:\infty}^{\mathrm{r}})|\theta^{\mathrm{s}}, x_{0:k}),$$
$$\mathrm{BR}_k^{\mathrm{r}}(s_{k:\infty}^{\mathrm{s}}|\theta^{\mathrm{r}}, x_{0:k}) := \operatorname*{arg\,max}_{\tilde{s}_{k:\infty}^{\mathrm{r}} \in \mathcal{S}_{k:\infty}^{\mathrm{r}}} \bar{U}_k^{\mathrm{r}}((s_{k:\infty}^{\mathrm{s}}, \tilde{s}_{k:\infty}^{\mathrm{r}})|\theta^{\mathrm{r}}, x_{0:k}). \tag{35}$$

We define the game formulated above by

$$\mathcal{G}_2 := (\mathcal{M}, \mathcal{S}, U, (\Theta^{\mathrm{s}}, \Theta^{\mathrm{r}}), (\pi_0^{\mathrm{s}}, \pi_0^{\mathrm{r}})), \tag{36}$$

where the defender's initial belief is *not* common information in contrast to $\mathcal{G}_1$.

TABLE II
CONTRASTING INGREDIENTS OF THE GAMES WITH SYMMETRIC AND ASYMMETRIC RECOGNITIONS

|  | symmetric recognition | asymmetric recognition |
|---|---|---|
| receiver's strategy | $s_k^{\mathrm{r}}(x_{0:k})$ | $s_k^{\mathrm{r}}(\theta^{\mathrm{r}}, x_{0:k})$ |
| sender's belief | N/A | $\pi^{\mathrm{s}}(\theta^{\mathrm{r}}|\theta^{\mathrm{s}})$ |
| receiver's belief | $\pi^{\mathrm{r}}(\theta)$ | $\pi^{\mathrm{r}}(\theta^{\mathrm{s}}|\theta^{\mathrm{r}})$ |
| sender's utility | $U^{\mathrm{s}}(s, \theta)$ | $U^{\mathrm{s}}(s, \theta^{\mathrm{s}})$ |
| receiver's utility | $U^{\mathrm{r}}(s)$ | $U^{\mathrm{r}}(s, \theta^{\mathrm{r}})$ |

In the following discussion, we analyze $\mathcal{G}_2$ through $\mathcal{G}_1$. To clarify their relationship, we describe the game $\mathcal{G}_1$ using the modified formulation. Define another game

$$\hat{\mathcal{G}}_2 := (\mathcal{M}, \mathcal{S}, U, (\Theta^{\mathrm{s}}, \Theta^{\mathrm{r}}), (\hat{\pi}_0^{\mathrm{s}}, \pi_0^{\mathrm{r}})), \tag{37}$$

where

$$\hat{\pi}_0^{\mathrm{s}}(\theta_{\mathrm{a}}^{\mathrm{r}}|\theta_{\mathrm{m}}^{\mathrm{s}})=1. \tag{38}$$

The initial belief means that the adversary believes that the defender is aware of the vulnerability. The situation of $\hat{\mathcal{G}}_2$ is the same as that of $\mathcal{G}_1$ if the defender is aware of the vulnerability. Thus, these games lead to the same consequence when the true types are $\theta_{\mathrm{m}}^{\mathrm{s}}$ and $\theta_{\mathrm{a}}^{\mathrm{r}}$. The following lemma holds.

**Lemma 4** Consider the games $\mathcal{G}_1$ and $\hat{\mathcal{G}}_2$. For a strategy profile $\hat{s}_2 = (\hat{s}_2^{\mathrm{s}}, \hat{s}_2^{\mathrm{r}})$ in $\hat{\mathcal{G}}_2$, let $s_1 = (s_1^{\mathrm{s}}, s_1^{\mathrm{r}})$ be a strategy profile in $\mathcal{G}_1$ such that

$$s_1^{\mathrm{s}} := \hat{s}_2^{\mathrm{s}}, \quad s_1^{\mathrm{r}} := \hat{s}_2^{\mathrm{r}}|_{\theta^{\mathrm{r}}=\theta_{\mathrm{a}}^{\mathrm{r}}} \tag{39}$$

where $\hat{s}_2^{\mathrm{r}}|_{\theta^{\mathrm{r}}=\theta_{\mathrm{a}}^{\mathrm{r}}}$ is the restriction of $\hat{s}_2^{\mathrm{r}}$ with $\theta^{\mathrm{r}} = \theta_{\mathrm{a}}^{\mathrm{r}}$. Then the probability measures induced by $s_1$ and $\hat{s}_2$ are equal when $\theta^{\mathrm{s}} = \theta_{\mathrm{m}}^{\mathrm{s}}$ and $\theta^{\mathrm{r}} = \theta_{\mathrm{a}}^{\mathrm{r}}$, i.e.,

$$\mathbb{P}_{\theta_{\mathrm{m}}^{\mathrm{s}}}^{s_1} = \mathbb{P}_{\theta_{\mathrm{m}}^{\mathrm{s}}, \theta_{\mathrm{a}}^{\mathrm{r}}}^{\hat{s}_2}. \tag{40}$$

Also, if $\hat{s}_{2,k:\infty}^{\mathrm{s}} \in \mathrm{BR}_k^{\mathrm{s}}(\hat{s}_{2,k:\infty}^{\mathrm{r}}|\theta_{\mathrm{m}}^{\mathrm{s}}, x_{0:k})$ then $s_{1,k:\infty}^{\mathrm{s}} \in \mathrm{BR}_k^{\mathrm{s}}(s_{1,k:\infty}^{\mathrm{r}}|\theta_{\mathrm{m}}^{\mathrm{s}}, x_{0:k})$.

We extend the notions of detection-averse utilities and asymptotically benign strategies to $\mathcal{G}_2$. Our objective is to investigate the effectiveness of the proposed defensive deception. It is possible to define detection-averse utilities directly using the game $\mathcal{G}_2$ as utilities where the resulting equilibrium leads the adversary to avoid being detected. However, this definition immediately means that the defensive deception works well, and any results from the definition cannot show its effectiveness. Instead, we say that utilities in $\mathcal{G}_2$ are detection-averse when the adversary avoids being detected if she is certain that the defender is aware of the vulnerability.

**Definition 3 (Detection-averse Utilities in $\mathcal{G}_2$)** A pair of utilities $(U^{\mathrm{s}}, U^{\mathrm{r}})$ in the game $\mathcal{G}_2$ are detection-averse utilities when

$$\lim_{k\to\infty} \pi_k^{\mathrm{r}}(\theta_{\mathrm{m}}^{\mathrm{s}}|\theta_{\mathrm{a}}^{\mathrm{r}}) < 1 \quad \mathbb{P}_{\theta_{\mathrm{m}}^{\mathrm{s}}, \theta_{\mathrm{a}}^{\mathrm{r}}}^{s}-\text{a.s.} \tag{41}$$

for any PBE $s$ of $\hat{\mathcal{G}}_2$.

Note that Definition 3 is a necessary requirement to make the game interesting, because the adversary is not afraid of being detected at all without this condition.

Next, we define desirable strategies that should be achieved by Bayesian defense mechanisms. We say a strategy in $\mathcal{G}_2$ to be asymptotically benign when it becomes benign regardless of the defender's awareness.

**Definition 4 (Asymptotically Benign Strategies in $\mathcal{G}_2$)** A strategy profile $s$ in the game $\mathcal{G}_2$ is asymptotically benign when

$$\lim_{k \to \infty} \delta \left( A_k^{\theta_{\mathrm{m}}}, A_k^{\theta_{\mathrm{b}}} \right) = 1 \quad \mathbb{P}_{\theta_{\mathrm{m}}^{\mathrm{s}}, \theta^{\mathrm{r}}}^{s}-\text{a.s.} \tag{42}$$

for any $\theta^{\mathrm{r}} \in \Theta^{\mathrm{r}}$.

Note that Definition 4 requires the strategy to be asymptotically benign for any $\theta^{\mathrm{r}} \in \Theta^{\mathrm{r}}$. In other words, the strategy is needed to be asymptotically benign even if the defender is unaware of the vulnerability.

### C. Passively Bluffing Strategies

We expect that there exists a chance of preventing attacks that exploit unnoticed vulnerabilities if the state does not possess information about the defender's recognition. To formally verify this expectation, we define passively bluffing strategies.

**Definition 5 (Passively Bluffing Strategies)** A strategy profile $s$ in $\mathcal{G}_2$ is a passively bluffing strategy profile when the sender's belief satisfies

$$\pi_k^{\mathrm{s}}(\theta^{\mathrm{r}}|X_{0:k}, \theta^{\mathrm{s}}) = \pi_0^{\mathrm{s}}(\theta^{\mathrm{r}}|\theta^{\mathrm{s}}) \quad \mathbb{P}_{\theta^{\mathrm{s}}, \theta^{\mathrm{r}}}^{s}-\text{a.s.} \tag{43}$$

for any $\theta^{\mathrm{s}} \in \Theta^{\mathrm{s}}, \theta^{\mathrm{r}} \in \Theta^{\mathrm{r}}$, and $k \in \mathbb{Z}_+$. A strategy profile set $\mathcal{S}$ in $\mathcal{G}_2$ is a passively bluffing strategy set when its all elements are passively bluffing.

Definition 5 requires the sender's belief to be invariant over time. If the strategy is passively bluffing, the adversary cannot identify whether the defender is aware of the exploited vulnerability or not even in an asymptotic sense. Note that the introduced passively bluffing strategies can be regarded as a commitment. It is well known that restricting feasible strategies, referred to as commitment, can be beneficial in a game [47], [48]. In what follows, we investigate the effectiveness of the specific commitment.

Passively bluffing strategies can relax the condition for asymptotically benign strategies. The following lemma holds.

**Lemma 5** Consider the game $\mathcal{G}_2$. If a passively bluffing strategy profile $s$ satisfies

$$\lim_{k \to \infty} \delta \left( A_k^{\theta_{\mathrm{m}}}, A_k^{\theta_{\mathrm{b}}} \right) = 1 \quad \mathbb{P}_{\theta_{\mathrm{m}}^{\mathrm{s}}, \theta_{\mathrm{a}}^{\mathrm{r}}}^{s}-\text{a.s.} \tag{44}$$

then $s$ is asymptotically benign.

The difference between (42) and (44) is the required receiver type. Lemma 5 implies that if a passively bluffing strategy profile is asymptotically benign when the receiver is aware of the vulnerability then the strategy is needed to be asymptotically benign even when the receiver is unaware of the vulnerability.

*Remark:* Although Definition 5 depends not only on the receiver's strategy but also on the sender's strategy for generality, the bluffing should be realized only by the defender in practice.

A simple defender's approach to achieving the bluffing is to choose reactions that do not influence the system's behavior. Let $\mathcal{R}_{\mathrm{pb}} \subset \mathcal{R}$ be the set of reactions such that the system's dynamics is independent of the reaction, i.e., the transition probability satisfies

$$P(x'|x, a, r) = P(x'|x, a, r') \tag{45}$$

for any $x' \in \mathcal{X}, x \in \mathcal{X}, a \in \mathcal{A}, r \in \mathcal{R}_{\mathrm{pb}}, r' \in \mathcal{R}_{\mathrm{pb}}$. If the receiver's strategy takes only reactions in $\mathcal{R}_{\mathrm{pb}}$, every strategy profile is passively bluffing. Indeed, because the transition probability is independent of $r \in \mathcal{R}_{\mathrm{pb}}$, the probability distribution of the state is independent of $\theta^{\mathrm{r}}$. Thus, from Bayes' rule, we have

$$\begin{aligned}
\pi_k^{\mathrm{s}}(\theta^{\mathrm{r}}|x_{0:k}, \theta^{\mathrm{s}}) &= \frac{p_{\theta^{\mathrm{s}}, \theta^{\mathrm{r}}}^{s}(x_{0:k})\pi_0^{\mathrm{s}}(\theta^{\mathrm{r}}|\theta^{\mathrm{s}})}{\sum_{\phi^{\mathrm{r}} \in \Theta^{\mathrm{r}}} p_{\theta^{\mathrm{s}}, \phi^{\mathrm{r}}}^{s}(x_{0:k})\pi_0^{\mathrm{s}}(\phi^{\mathrm{r}}|\theta^{\mathrm{s}})} \\
&= \frac{p_{\theta^{\mathrm{s}}}^{s}(x_{0:k})\pi_0^{\mathrm{s}}(\theta^{\mathrm{r}}|\theta^{\mathrm{s}})}{\sum_{\phi^{\mathrm{r}} \in \Theta^{\mathrm{r}}} p_{\theta^{\mathrm{s}}}^{s}(x_{0:k})\pi_0^{\mathrm{s}}(\phi^{\mathrm{r}}|\theta^{\mathrm{s}})} \\
&= \frac{\pi_0^{\mathrm{s}}(\theta^{\mathrm{r}}|\theta^{\mathrm{s}})}{\sum_{\phi^{\mathrm{r}} \in \Theta^{\mathrm{r}}} \pi_0^{\mathrm{s}}(\phi^{\mathrm{r}}|\theta^{\mathrm{s}})} \\
&= \pi_0^{\mathrm{s}}(\theta^{\mathrm{r}}|\theta^{\mathrm{s}})
\end{aligned} \tag{46}$$

when $p_{\theta^{\mathrm{s}}, \theta^{\mathrm{r}}}^{s}(x_{0:k}) \neq 0$. An example of such reactions is just analyzing the network log and raising an alarm inside the operation room without applying control on the system itself. Note that the reaction still affects the players' decision making through their utility functions, even if (45) holds.

### D. Analysis

Our expectation can be described in a quantitative form based on the definition of passively bluffing strategies, which lead to a simple representation of the sender's utility. If $s$ is passively bluffing, the sender's belief is invariant over time. Hence, the sender's utility with infinite horizon is given by

$$\bar{U}_k^{\mathrm{s}}(s_{k:\infty}|\theta_{\mathrm{m}}^{\mathrm{s}}, x_{0:k}) = \sum_{\theta^{\mathrm{r}} \in \Theta^{\mathrm{r}}} \bar{U}_{k, \theta^{\mathrm{r}}}^{\mathrm{s}}(s_{k:\infty}|\theta_{\mathrm{m}}^{\mathrm{s}}, x_{0:k})\pi_0^{\mathrm{s}}(\theta^{\mathrm{r}}|\theta_{\mathrm{m}}^{\mathrm{s}}) \tag{47}$$

where

$$\begin{aligned}
\bar{U}_{k, \theta^{\mathrm{r}}}^{\mathrm{s}}(s_{k:\infty}|\theta_{\mathrm{m}}^{\mathrm{s}}, x_{0:k}) &:= \lim_{T \to \infty} \frac{1}{T+1} \\
&\times \mathbb{E}_{\theta^{\mathrm{r}}}^{s}\left[ \sum_{\tau=k}^{k+T} U^{\mathrm{s}}(\theta_{\mathrm{m}}^{\mathrm{s}}, X_{\tau}, s_{\tau}^{\mathrm{s}}(\theta_{\mathrm{m}}^{\mathrm{s}}, X_{0:\tau}), s_{\tau}^{\mathrm{r}}(\theta^{\mathrm{r}}, X_{0:\tau})) \middle| x_{0:k} \right].
\end{aligned} \tag{48}$$

Note that $\bar{U}_{k, \theta_{\mathrm{a}}^{\mathrm{r}}}^{\mathrm{s}}$ and $\bar{U}_{k, \theta_{\mathrm{u}}^{\mathrm{r}}}^{\mathrm{s}}$ denote the sender's utilities of the two cases where the defender is aware and unaware of the vulnerability, respectively. Thus (47) implies that the sender's utility is simply given as a sum weighted by her initial beliefs when the strategy is passively bluffing. Therefore, we can expect that the sender possibly stops the execution in the middle of the attack if $\pi_0^{\mathrm{s}}(\theta_{\mathrm{a}}^{\mathrm{r}}|\theta_{\mathrm{m}}^{\mathrm{s}})$ is sufficiently large. We show the existence of such sender's initial belief. Note that $\pi_0^{\mathrm{s}}(\theta_{\mathrm{a}}^{\mathrm{r}}|\theta_{\mathrm{m}}^{\mathrm{s}}) = 1$ is the trivial case, and thus we assume that sender's initial beliefs that are strictly less than one.

First, we rephrase the result in Sec. III. Let $\mathcal{S}_{\mathrm{nab}}$ denote the set of non-asymptotically-benign strategies in $\mathcal{G}_2$. Our aim here is to show that the set of PBE of $\mathcal{G}_2$ does not overlap with $\mathcal{S}_{\mathrm{nab}}$ when the attacker strongly believes that the defender is

aware of the vulnerability. It suffices to show that there is no overlap between the set of PBE of $\mathcal{G}_2$ and $\mathcal{S}^*_{\mathrm{nab}} := \mathcal{S}_{\mathrm{nab}} \cap \mathcal{S}^*$ where

$$
\mathcal{S}^* := \{(s^{\mathrm{s}}, s^{\mathrm{r}}) \in \mathcal{S} : s^{\mathrm{s}}_{k:\infty} \in \mathrm{BR}^{\mathrm{r}}_k(s^{\mathrm{r}}_{k:\infty}|\theta^{\mathrm{s}}_{\mathrm{b}}, x_{0:k}),
$$
$$
s^{\mathrm{r}}_{k:\infty} \in \mathrm{BR}^{\mathrm{r}}_k(s^{\mathrm{s}}_{k:\infty}|\theta^{\mathrm{r}}, x_{0:k}), \forall \theta^{\mathrm{r}} \in \Theta^{\mathrm{r}}, k \in \mathbb{Z}_+, x_{0:k} \in \mathcal{X}^{k+1}\}, \quad (49)
$$

where the benign sender and the receiver with any type take their best response strategies. Note that, $\mathcal{S}_{\mathrm{nab}}$ and $\mathcal{S}^*$ of the games $\mathcal{G}_2$ and $\hat{\mathcal{G}}_2$ are identical because the sets are independent of the malicious sender's belief. The following lemma is another description of the claim of Theorem 3 with respect to $\bar{U}^{\mathrm{s}}_{\theta^{\mathrm{r}}_{\mathrm{a}}}$ and $\mathcal{S}^*_{\mathrm{nab}}$.

**Lemma 6** Consider the game $\hat{\mathcal{G}}_2$ with detection-averse utilities. Let Assumption 1 hold. For any strategy profile $s = (s^{\mathrm{s}}, s^{\mathrm{r}})$ in $\mathcal{S}^*_{\mathrm{nab}}$, there exists $\tilde{s}^{\mathrm{s}} \in \mathcal{S}^{\mathrm{s}}$ such that

$$
D_{\theta^{\mathrm{r}}_{\mathrm{a}}}(s, \tilde{s}^{\mathrm{s}}) > 0 \quad (50)
$$

holds where

$$
D_{\theta^{\mathrm{r}}}(s, \tilde{s}^{\mathrm{s}}) := \bar{U}^{\mathrm{s}}_{\theta^{\mathrm{r}}}((\tilde{s}^{\mathrm{s}}, s^{\mathrm{r}}), \theta^{\mathrm{s}}_{\mathrm{m}}) - \bar{U}^{\mathrm{s}}_{\theta^{\mathrm{r}}}(s, \theta^{\mathrm{s}}_{\mathrm{m}}). \quad (51)
$$

Lemma 6 implies the existence of a function

$$
g : \mathcal{S}^*_{\mathrm{nab}} \to \mathcal{S}^{\mathrm{s}} \quad \text{s.t.} \quad D_{\theta^{\mathrm{r}}_{\mathrm{a}}}(s, g(s)) > 0 \quad (52)
$$

for any $s \in \mathcal{S}^*_{\mathrm{nab}}$. Thus we have $\gamma \geq 0$ where

$$
\gamma := \inf_{s \in \mathcal{S}^*_{\mathrm{nab}}} D_{\theta^{\mathrm{r}}_{\mathrm{a}}}(s, g(s)). \quad (53)
$$

We here make an assumption that $D_{\theta^{\mathrm{r}}_{\mathrm{a}}}(s, g(s))$ is uniformly lower bounded by a positive value.

**Assumption 2** For the game $\hat{\mathcal{G}}_2$, there exists $g$ in (52) such that the infimum (53) is positive, i.e., $\gamma > 0$.

Assumption 2 eliminates the case where the difference between the sender's utilities achievable by asymptotically benign strategies and non-asymptotically-benign strategies is infinitesimally small.

The following theorem, the main result of this section, holds.

**Theorem 4** Consider the game $\mathcal{G}_2$ with detection-averse utilities and a passively bluffing strategy set. Let Assumptions 1 and 2 hold. Then, there exists a sender's initial belief $\pi^{\mathrm{s}}_0(\theta^{\mathrm{r}}_{\mathrm{a}}|\theta^{\mathrm{s}}_{\mathrm{m}}) < 1$ such that every PBE of $\mathcal{G}_2$ is asymptotically benign.

Theorem 4 implies that the system can possibly be protected by passively bluffing strategies if the attacker strongly believes that the defender is aware of the vulnerability. The result suggests the importance of concealing the defender's recognition and the effectiveness of defensive deception.

## V. SIMULATION

In this section, we confirm the theoretical results through numerical simulation.

TABLE III
TRANSITION PROBABILITIES FROM THE ABNORMAL STATE TO ABNORMAL STATE.

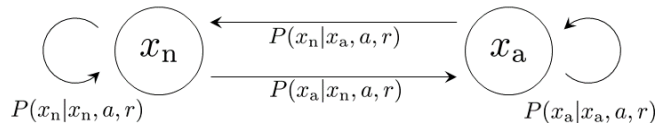| $P(x_{\mathrm{a}}|x_{\mathrm{a}}, a, r)$ | $r_{\mathrm{b}}$ | $r_{\mathrm{m}}$ | $r^{\mathrm{b}}_{\mathrm{m}}$ |
|---|---|---|---|
| $a_{\mathrm{b}}$ | 0.5 | 0.3 | 0.5 |
| $a_{\mathrm{m}}$ | 0.6 | 0.4 | 0.6 |



Fig. 6. State transition diagram of the numerical example.

### A. Fundamental Setup

We assume the state space and the action space to be binary, i.e., $\mathcal{X} = \{x_{\mathrm{n}}, x_{\mathrm{a}}\}$ and $\mathcal{A} = \{a_{\mathrm{b}}, a_{\mathrm{m}}\}$. The states $x_{\mathrm{n}}$ and $x_{\mathrm{a}}$ represent the normal and abnormal states, respectively, and $a_{\mathrm{b}}$ and $a_{\mathrm{m}}$ represent benign and malicious actions, respectively. The benign and malicious actions correspond to nominal and malicious control signals, respectively. The reaction set is given by $\mathcal{R} = \{r_{\mathrm{b}}, r_{\mathrm{m}}, r^{\mathrm{b}}_{\mathrm{m}}\}$. The state transition diagram is depicted in Fig. 6. The initial state is set to $x_{\mathrm{n}}$. The transition probability is given as follows. Set the transition probability from $x_{\mathrm{n}}$ to be given by

$$
P(x_{\mathrm{a}}|x_{\mathrm{n}}, a, r) = \begin{cases} 0.2 & \text{if } a = a_{\mathrm{b}}, \\ 0.3 & \text{if } a = a_{\mathrm{m}} \end{cases} \quad (54)
$$

for any $r \in \mathcal{R}$, which means that the probability from the normal state to the abnormal state is increased by the malicious action and it is independent of the reaction. The transition probability from $x_{\mathrm{a}}$ to $x_{\mathrm{a}}$ is given by Table III. The probability from the abnormal state to the abnormal state is increased by the malicious action and it is decreased by the reaction $r_{\mathrm{m}}$. The reaction $r^{\mathrm{b}}_{\mathrm{m}}$ corresponds to bluffing since it induces the same transition probability as $r_{\mathrm{b}}$.

The utilities are given as follows. The benign sender's utility is

$$
U^{\mathrm{s}}(\theta_{\mathrm{b}}, x, a, r) = \begin{cases} 1 & \text{if } x = x_n, \\ 0 & \text{otherwise} \end{cases} \quad (55)
$$

for any $a \in \mathcal{A}$ and $r \in \mathcal{R}$, which means that the benign sender prefers the nominal state regardless of other variables. The malicious sender's utility is given by Table IV. The benign action $a_{\mathrm{b}}$ is a risk-free action, which always induces zero utility, while the malicious action $a_{\mathrm{m}}$ is a risky action. If the reaction is $r_{\mathrm{b}}$, the malicious sender obtains positive utility, where the abnormal state $x_{\mathrm{a}}$ is more preferred than $x_{\mathrm{n}}$. On the other hand, if the reaction is $r_{\mathrm{m}}$, the malicious sender incurs loss. The receiver's utility is set to be independent on $a \in \mathcal{A}$ and given by Table V, where $a$ is omitted. The receiver obtains utility only when she takes an appropriate reaction depending on the sender type. When an appropriate reaction is chosen, the normal state is more preferred than the abnormal state. Note that $r^{\mathrm{b}}_{\mathrm{m}}$ induces the same utilities as those with $r_{\mathrm{m}}$ but it increases the probability of the abnormal

TABLE IV
MALICIOUS SENDER'S UTILITY

| $U^{\mathrm{s}}(\theta_{\mathrm{m}}, x, a_{\mathrm{b}}, r)$ | $r_{\mathrm{b}}$ | $r_{\mathrm{m}}$ | $r_{\mathrm{m}}^{\mathrm{b}}$ |
|---|---|---|---|
| $x_{\mathrm{n}}$ | 0 | 0 | 0 |
| $x_{\mathrm{a}}$ | 0 | 0 | 0 |
| $U^{\mathrm{s}}(\theta_{\mathrm{m}}, x, a_{\mathrm{m}}, r)$ | $r_{\mathrm{b}}$ | $r_{\mathrm{m}}$ | $r_{\mathrm{m}}^{\mathrm{b}}$ |
| $x_{\mathrm{n}}$ | 1 | -3 | -3 |
| $x_{\mathrm{a}}$ | 2 | -3 | -3 |

TABLE V
RECEIVER'S UTILITY

| $U^{\mathrm{r}}(\theta_{\mathrm{b}}, x, r)$ | $r_{\mathrm{b}}$ | $r_{\mathrm{m}}$ | $r_{\mathrm{m}}^{\mathrm{b}}$ | $U^{\mathrm{r}}(\theta_{\mathrm{m}}, x, r)$ | $r_{\mathrm{b}}$ | $r_{\mathrm{m}}$ | $r_{\mathrm{m}}^{\mathrm{b}}$ |
|---|---|---|---|---|---|---|---|
| $x_{\mathrm{n}}$ | 5 | 0 | 0 | $x_{\mathrm{n}}$ | 0 | 5 | 5 |
| $x_{\mathrm{a}}$ | 1 | 0 | 0 | $x_{\mathrm{a}}$ | 0 | 1 | 1 |

state. Therefore, there is no motivation to choose $r_{\mathrm{m}}^{\mathrm{b}}$ when the defender's recognition is known to the attacker.

Since it is difficult to compute an exact equilibrium for the infinite time horizon problem, we treat a sequence of equilibria for a finite time horizon problem as a tractable approximation [49]. Letting $(s_k^{\mathrm{s}}, s_k^{\mathrm{r}}, \ldots, s_{k+T-1}^{\mathrm{s}}, s_{k+T-1}^{\mathrm{r}})$ be the resulting equilibrium of the finite time horizon game, we use $s_k^{\mathrm{s}}$ and $s_k^{\mathrm{r}}$ as the $k$th strategies as with receding horizon control. The equilibrium is obtained through brute-force search. For the game $\mathcal{G}_2$, the strategies in the simulation are given in a similar manner. The horizon length is set to $T = 2$. In the numerical examples, the equilibrium is uniquely determined.

### B. Simulation: Asymptotic Security

In the first scenario, we consider the case where the vulnerability is known, and thus this situation corresponds to the game $\mathcal{G}_1$ in (15). The initial belief is given by $\pi_0^{\mathrm{r}}(\theta_{\mathrm{m}}) = 0.01$, which is known to the sender. The true sender type is given by $\theta^{\mathrm{s}} = \theta_{\mathrm{m}}^{\mathrm{s}}$.

Under the setting, sample paths of the belief on the malicious sender, the action, and the reaction with $\theta = \theta_{\mathrm{m}}$ are depicted in Fig. 7. The belief converges to a nonzero value over time as claimed by Theorems 1 and 2. The action converges to the benign action as claimed by Theorem 3. The graphs evidence asymptotic security achieved by the Bayesian defense mechanism. In more detail, it can be observed that the malicious sender takes the malicious action $a_{\mathrm{m}}$ while the receiver takes the reaction $r_{\mathrm{b}}$ until about the time step $k = 50$. This is because the receiver's belief on the malicious sender is low during the beginning of the game. On the other hand, between the time steps $k = 50$ and $k = 100$, $a_{\mathrm{b}}$ and $r_{\mathrm{m}}$ sporadically appear because the belief is increased. Finally, after the time step $k = 100$, the belief exceeds a threshold, which results in the fixed actions $a = a_{\mathrm{b}}$ and $r = r_{\mathrm{m}}$. It is notable that $r_{\mathrm{m}}^{\mathrm{b}}$ is not chosen at all since there is no reason for it, as explained above.

### C. Simulation: Defensive Deception using Bluffing

In the second scenario, we consider the case where the defender is unaware of the vulnerability and the attacker is unaware of the defender's unawareness. Then this situation
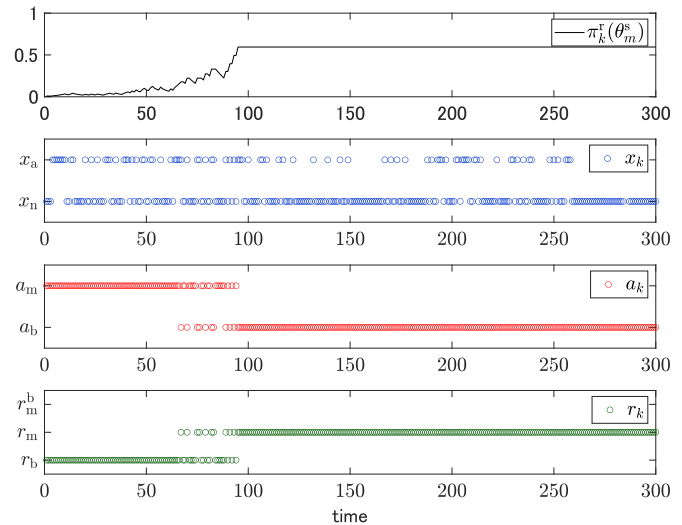


Fig. 7. Sample paths of the belief on the malicious sender, state, action, and reactions with $\theta = \theta_{\mathrm{m}}$. The belief converges to a nonzero value over time as claimed by Theorems 1 and 2. The action converges to the benign action as claimed by Theorem 3. The results evidence asymptotic security achieved by the Bayesian defense mechanism.
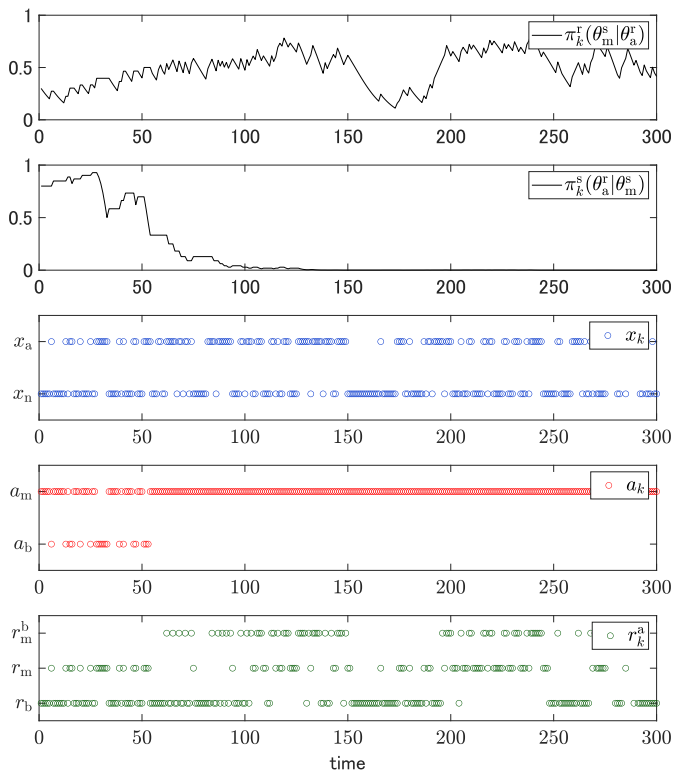
corresponds to the game $\mathcal{G}_2$ in (36). The initial beliefs are given by $\pi_0^{\mathrm{r}}(\theta_{\mathrm{m}}^{\mathrm{s}}|\theta_{\mathrm{a}}^{\mathrm{r}}) = 0.3$ and $\pi_0^{\mathrm{s}}(\theta_{\mathrm{a}}^{\mathrm{r}}|\theta_{\mathrm{m}}^{\mathrm{s}}) = 0.8$. The true types are given by $\theta^{\mathrm{s}} = \theta_{\mathrm{m}}^{\mathrm{s}}$ and $\theta^{\mathrm{r}} = \theta_{\mathrm{u}}^{\mathrm{r}}$. Note that $\pi_0^{\mathrm{r}}(\theta_{\mathrm{m}}^{\mathrm{s}}|\theta_{\mathrm{u}}^{\mathrm{r}}) = 0$ and hence the defender is completely unaware of the attack while the game is proceeding.

We first consider the case where the strategy is *not* passively bluffing. The same transition probability as that used in the previous simulation, where it depends on the receiver's reaction. As a result, the state possesses information about the receiver type.

Fig. 8 depicts sample paths of the receiver's belief on the malicious sender if the receiver were aware of the vulnerability, the sender's belief on the receiver being aware, the actual state, the actual action, and reactions that would be taken by the receiver being aware. It can be observed that the sender's belief converges to zero, i.e., the sender notices that the receiver is unaware of the vulnerability. As a result, malicious actions are constantly taken after a sufficiently large number of time steps. The result indicates that the defense mechanism fails to defend the system in this case.

We next consider the bluffing case. As a commitment for passively bluffing strategy, we restrict the reaction set to $\mathcal{R} = \{r_{\mathrm{b}}, r_{\mathrm{m}}^{\mathrm{b}}\}$. Then the transition probability is independent of the reaction, and hence any strategy becomes passively bluffing.

Fig. 9 depicts sample paths of those depicted in Fig. 8 under the bluffing setting. The sender's belief is invariant over time because the state does not possess information about the receiver type. Thus, the malicious sender remains cautious about being detected. As a result, the benign action is continuously taken after a sufficiently large number of time steps in contrast to Fig. 8. The result indicates that asymptotic security is achieved by the bluffing even if the defender is unaware of the vulnerability. The simulation suggests importance of concealing the defender's belief even if it degrades control performance.
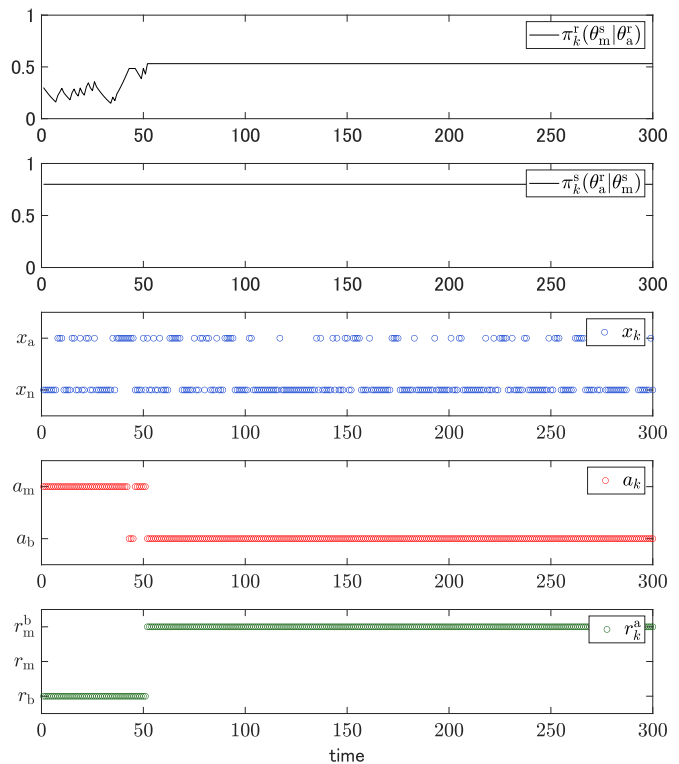
Fig. 8. Sample paths of the receiver's belief on the malicious sender when the receiver is aware, the sender's belief on the receiver being aware, state, action, and reactions that would be taken when the receiver were aware, where the strategy is *not* passively bluffing. The sender's belief converges to zero, i.e., the attacker notices that the defender is unaware of the vulnerability. As a result, the malicious action is continuously taken after a sufficiently large number of time steps.



Fig. 9. Sample paths of the receiver's belief on the malicious sender when the receiver is aware, the sender's belief on the receiver being aware, state, action, and reactions that would be taken when the receiver were aware, where the strategy is passively bluffing. The sender's belief is invariant over time because the state does not possess information about the receiver type. Thus, the attacker keeps to be cautious about being detected. As a result, the benign action is continuously taken after a sufficiently large number of time steps in contrast to Fig. 8.

## VI. CONCLUSION

This study has analyzed defense capability achieved by Bayesian defense mechanisms. It has been shown that the system to be protected can be guaranteed to be secure by Bayesian defense mechanisms provided that effective countermeasures are implemented. This fact implies that model knowledge can prevent the defender from being deceived in an asymptotic sense. As a defensive deception utilizing the derived asymptotic security, bluffing utilizing asymmetric recognition has been considered. It has also been shown that the attacker possibly stops the execution in the middle of the attack in a rational manner when she strongly believes the defender to be aware of the vulnerability, even if the vulnerability is unnoticed.

Important future work includes an extension to infinite state spaces because the state space in control systems typically is a subset of the Euclidean space. For this purpose, existing Bayesian consistency analysis for general sample space should be useful [36]. Moreover, although it is assumed that the state is observable in this framework, a generalization to partially observable setting is a more practical setting. We expect that the key properties such as Lemma 1 still holds if the system requirement, such as Assumption 1, can be appropriately modified. Another direction is to extend the results to non-binary types. Finally, finding a general condition for detection-averse utilities is an important issue. For this purpose, the example in Appendix I should be helpful.

## APPENDIX I
## EXAMPLE OF DETECTION-AVERSE UTILITIES

Consider an MDP with binary spaces $\mathcal{X} = \{x_\mathrm{n}, x_\mathrm{a}\}, \mathcal{A} = \{a_\mathrm{b}, a_\mathrm{m}\}$, and $\mathcal{R} = \{r_\mathrm{b}, r_\mathrm{m}\}$. The initial state is $x_\mathrm{n}$ and it goes to $x_\mathrm{a}$ with probability $p > 0$ when $a = a_\mathrm{m}$ and stays at $x_\mathrm{m}$ otherwise. The objective of the receiver is to detect the true state, which is modeled by

$$U^\mathrm{r}(\theta, x, a, r) = \begin{cases} 1 & \text{if } (\theta, r) = (\theta_\mathrm{b}, r_\mathrm{b}) \text{ or } (\theta_\mathrm{m}, r_\mathrm{m}), \\ 0 & \text{otherwise.} \end{cases} \tag{56}$$

The malicious sender's utility is given by

$$U^\mathrm{s}(\theta_\mathrm{m}, x, a, r_\mathrm{b}) = \begin{cases} 1 & \text{if } a = a_\mathrm{m}, \\ 0 & \text{otherwise,} \end{cases} \quad U^\mathrm{s}(\theta_\mathrm{m}, x, a, r_\mathrm{m}) = -1, \tag{57}$$

which means that the adversary wants to avoid being detected. The benign sender is assumed to choose $r_\mathrm{b}$ anytime. The initial belief satisfies $\pi_0^\mathrm{r}(\theta_\mathrm{m}) < 1/2$. Let $(s^\mathrm{s}, s^\mathrm{r})$ be a strategy such that $\pi_\infty^{\theta_\mathrm{m}} = 1$ with probability $q > 0$. Take such $x_{0:\infty}$ and then there exists $N \in \mathbb{Z}_+$ such that $\pi_k^\mathrm{r}(\theta_\mathrm{m}|x_{0:k}) = 1$ for any $k > N$ since

$$\pi_k^\mathrm{r}(\theta_\mathrm{m}|x_{0:k}) = \begin{cases} \pi_0^\mathrm{r}(\theta_\mathrm{m}) & \text{if } x_\tau = x_\mathrm{n} \ \forall \tau \in \{0, \dots, k\}, \\ 1 & \text{otherwise.} \end{cases} \tag{58}$$

The receiver's best response at $x_{0:k}$ is to take $r_{\mathrm{m}}$, which leads to the sender's average utility at $x_{0:k}$ equal $-1$. Thus, the sender's average utility at the initial state is $-q < 0$. However, if the sender takes the strategy such that $a_{\mathrm{k}} = a_{\mathrm{b}}$ for any $k \in \mathbb{Z}_+$, then the sender's average utility at the initial state is $0$. Thus, $s^{\mathrm{s}}$ is not a best response to $s^{\mathrm{r}}$, which means that the utilities are detection-averse.

## APPENDIX II
## PROOFS OF PROPOSITIONS

In the proofs, we omit the symbol $s$ in the notation for simplicity when no confusion arises.

*Proof of Lemma 1:* It is clear that $\pi_k^\theta$ is adapted to the filtration $\sigma(X_{0:k})$. It is also clear that $\pi_k^\theta$ is integrable with respect to $\mathbb{P}_\theta$ since it is bounded. Thus, it suffices to show

$$\mathbb{E}\left[\pi_{k+1}^\theta | \sigma(X_{0:k})\right] \geq \pi_k^\theta \quad \mathbb{P}_\theta-\text{a.s.} \quad (59)$$

for the claim. For a fixed outcome $\omega \in \Omega$ with which $X_{0:k}(\omega) = x_{0:k}$, the inequality is equivalent to

$$\sum_{x_{k+1} \in \mathcal{X}} p_\theta(x_{k+1}|x_{0:k})\pi_{k+1}(\theta|x_{0:k+1}) \geq \pi_k^{\mathrm{r}}(\theta|x_{0:k}). \quad (60)$$

Thus it suffices to show (60) for any $k \in \mathbb{N}$ and $x_{0:k} \in \mathcal{X}^{k+1}$.

First, we reduce the index of the summation in (60). When $\pi_k^{\mathrm{r}}(\theta|x_{0:k}) = 0$, the inequality (60) always holds. Thus we assume $\pi_k^{\mathrm{r}}(\theta|x_{0:k}) > 0$ in the following. Define

$$\mathcal{X}_k^0 := \left\{ x_{k+1} \in \mathcal{X} : \sum_{\phi \in \Theta} p_\phi(x_{k+1}|x_{0:k})\pi_k^{\mathrm{r}}(\phi|x_{0:k}) = 0 \right\}. \quad (61)$$

Because $\pi_k^{\mathrm{r}}(\phi|x_{0:k})$ is positive for any $\phi \in \Theta$, if $x_{k+1}$ belongs to $\mathcal{X}_k^0$ then $p_\theta(x_{k+1}|x_{0:k}) = 0$. Hence (60) is equivalent to

$$\sum_{x_{k+1} \in \mathcal{X}_k^+} p_\theta(x_{k+1}|x_{0:k})\pi_{k+1}^{\mathrm{r}}(\theta|x_{0:k+1}) \geq \pi_k^{\mathrm{r}}(\theta|x_{0:k}) \quad (62)$$

where $\mathcal{X}_k^+ := \mathcal{X} \setminus \mathcal{X}_k^0$.

To simplify notation, we define

$$\pi(\theta) := \pi_k^{\mathrm{r}}(\theta|x_{0:k}), \ p_\phi(x) := p_\phi(x|x_{0:k}), \ \mathcal{X}^+ := \mathcal{X}_k^+ \quad (63)$$

for fixed $k$ and $x_{0:k}$. Then the inequality (62) is equivalent to

$$\sum_{x \in \mathcal{X}^+} p_\theta(x)\frac{p_\theta(x)\pi(\theta)}{\sum_{\phi \in \Theta} p_\phi(x)\pi(\phi)} \geq \pi(\theta). \quad (64)$$

Because $\pi(\theta) > 0$, this inequality is equivalent to

$$\underbrace{\sum_{x \in \mathcal{X}^+} p_\theta(x)\frac{p_\theta(x)}{\sum_{\phi \in \Theta} p_\phi(x)\pi(\phi)}}_{=:G(\theta)} \geq 1. \quad (65)$$

By rewriting $G(\theta)$, we have

$$G(\theta) = \sum_{x \in \mathcal{X}^+} p_\theta(x)\frac{1}{\sum_{\phi \in \Theta} \frac{p_\phi(x)}{p_\theta(x)}\pi(\phi)}. \quad (66)$$

By applying Jensen's inequality in (1) with the functions $p(x) := p_\theta(x), a(x) := \sum_{\phi \in \Theta} \frac{p_\phi(x)}{p_\theta(x)}\pi(\phi)$, and $\varphi(\xi) := 1/\xi$, we have

$$
\begin{aligned}
G(\theta) &\geq \varphi\left(\sum_{x \in \mathcal{X}^+} p_\theta(x) \sum_{\phi \in \Theta} \frac{p_\phi(x)}{p_\theta(x)}\pi(\phi)\right) \\
&= \varphi\left(\sum_{x \in \mathcal{X}^+} \sum_{\phi \in \Theta} p_\phi(x)\pi(\phi)\right) \\
&= \varphi\left(\sum_{\phi \in \Theta} \left(\sum_{x \in \mathcal{X}^+} p_\phi(x)\right)\pi(\phi)\right) \\
&= \varphi\left(\sum_{\phi \in \Theta} \pi(\phi)\right) \\
&= \varphi(1) \\
&= 1,
\end{aligned}
\quad (67)
$$

which leads to the claim. $\qquad \square$

*Proof of Theorem 1:* Because the belief is uniformly bounded over time, we have $\sup_{k \in \mathbb{N}} \mathbb{E}_\theta\left[\pi_k^\theta\right] < \infty$. From Lemma 1 and Doob's convergence theorem [50, Theorem 4.4.1], the claim holds. $\qquad \square$

*Proof of Lemma 2:* Since $\pi_0^\theta > 0$, we have $\pi_k^\theta > 0$ $\mathbb{P}_\theta$-almost surely for any $k \in \mathbb{N}$. Thus $\log(\pi_k^\theta)$ is well-defined. We first show that $\log(\pi_k^\theta)$ is a submartingale with respect to the probability measure $\mathbb{P}_\theta$ and the filtration $\sigma(X_{0:k})$. It is clear that $\log(\pi_k^\theta)$ is adapted to the filtration $\sigma(X_{0:k})$. Because the number of elements in the support of $\log(\pi_k^\theta)$ is finite, $\log(\pi_k^\theta)$ is integrable for any $k \in \mathbb{N}$. Thus it suffices to show that

$$\mathbb{E}_\theta\left[\log(\pi_{k+1}^\theta)|\sigma(X_{0:k})\right] \geq \log(\pi_k^\theta) \quad \mathbb{P}_\theta-\text{a.s.} \quad (68)$$

As in the proof of Lemma 1, this inequality is equivalent to

$$\sum_{x_{k+1} \in \mathcal{X}_k^+} p_\theta(x_{k+1}|x_{0:k})\log(\pi_{k+1}^{\mathrm{r}}(\theta|x_{0:k+1})) \geq \log(\pi_k^{\mathrm{r}}(\theta|x_{0:k})) \quad (69)$$

for any $x_{0:k} \in \mathcal{X}^{k+1}$. With the notation (63), this is equivalent to

$$\sum_{x \in \mathcal{X}^+} p_\theta(x)\log\left(\frac{p_\theta(x)\pi(\theta)}{\sum_{\phi \in \Theta} p_\phi(x)\pi(\phi)}\right) \geq \log(\pi(\theta)). \quad (70)$$

Because the left-hand side can be rewritten by

$$
\begin{aligned}
&\sum_{x \in \mathcal{X}^+} p_\theta(x)\log\left(\frac{p_\theta(x)\pi(\theta)}{\sum_{\phi \in \Theta} p_\phi(x)\pi(\phi)}\right) \\
&= \sum_{x \in \mathcal{X}^+} p_\theta(x)\left\{\log\left(\frac{p_\theta(x)}{\sum_{\phi \in \Theta} p_\phi(x)\pi(\phi)}\right) + \log(\pi(\theta))\right\} \\
&= \sum_{x \in \mathcal{X}^+} p_\theta(x)\log\left(\frac{p_\theta(x)}{\sum_{\phi \in \Theta} p_\phi(x)\pi(\phi)}\right) + \log(\pi(\theta)),
\end{aligned}
\quad (71)
$$

the inequality (70) is equivalent to

$$\sum_{x \in \mathcal{X}^+} p_\theta(x)\log\left(\frac{p_\theta(x)}{\sum_{\phi \in \Theta} p_\phi(x)\pi(\phi)}\right) \geq 0, \quad (72)$$

which is also equivalent to

$$\underbrace{\sum_{x \in \mathcal{X}^+} p_\theta(x)\log\left(\frac{\sum_{\phi \in \Theta} p_\phi(x)\pi(\phi)}{p_\theta(x)}\right)}_{H(\theta)} \leq 0. \quad (73)$$

By applying Jensen's inequality for a concave function with $p(x) := p_\theta(x), a(x) := \frac{\sum_{\phi \in \Theta} p_\phi(x) \pi(\phi)}{p_\theta(x)}$, and $\varphi(\xi) := \log(\xi)$, we have

$$
\begin{aligned}
H(\theta) &\leq \log\left(\sum_{x \in \mathcal{X}^+} p_\theta(x) \frac{\sum_{\phi \in \Theta} p_\phi(x) \pi(\phi)}{p_\theta(x)}\right) \\
&= \log\left(\sum_{x \in \mathcal{X}^+} \sum_{\phi \in \Theta} p_\phi(x) \pi(\phi)\right) \\
&= \log(1) \\
&= 0,
\end{aligned}
\tag{74}
$$

which implies that $\log(\pi_k^\theta)$ is a submartingale.

From Doob's convergence theorem [50, Theorem 4.4.1], it suffices to show that the expectation of the nonnegative part of $\log(\pi_k^\theta)$ is uniformly bounded. Because $\pi_k^\theta \in (0, 1]$, $\log(\pi_k^\theta)$ is nonpositive for any $k \in \mathbb{N}$, and hence the uniform boundedness holds. $\square$

*Proof of Theorem 2:* We prove the claim by contradiction. Define $E$ as the inverse image of $\{0\}$ for $\pi_\infty^\theta$. Assume that $\mathbb{P}_\theta(E) > 0$. For any $\omega \in E$, $\pi_k^\theta(\omega) \to 0$ as $k \to \infty$. Hence, from the continuity of logarithm functions, it turns out that $\log(\pi_k^\theta(\omega)) \to -\infty$ as $k \to \infty$. This means that $\log(\pi_k^\theta(\omega))$ diverges for $\omega \in E$. However, Lemma 2 states that $\log(\pi_k^\theta)$ converges $\mathbb{P}_\theta$-almost surely. This is a contradiction. $\square$

*Proof of Lemma 3:* We first show that the coefficient of Bayes' rule converges to one, i.e.,

$$
\lim_{k \to \infty} f_k(\theta_m, X_{0:k}) = 1 \quad \mathbb{P}_{\theta_m}-\text{a.s.}
\tag{75}
$$

where $f_{k+1}(\theta, x_{0:k}) := \pi_{k+1}^\theta / \pi_k^\theta$. Because $\pi_\infty^{\theta_m}$ is nonzero $\mathbb{P}_{\theta_m}$-almost surely from Theorem 2, we have

$$
\begin{aligned}
\lim_{k \to \infty} f_{k+1}(\theta_m, X_{0:k}) &= \lim_{k \to \infty}\left(\pi_{k+1}^{\theta_m} / \pi_k^{\theta_m}\right) \\
&= \pi_\infty^{\theta_m} / \pi_\infty^{\theta_m} \\
&= 1
\end{aligned}
\tag{76}
$$

$\mathbb{P}_{\theta_m}$-almost surely. Thus (75) holds.

It is observed that $f_k$ can be calculated by

$$
f_k^s(\theta, x_{0:k}) = \frac{p_\theta^s(x_k|x_{0:k-1})}{\sum_{\phi \in \Theta} p_\phi^s(x_k|x_{0:k-1}) \pi_{k-1}^r(\phi|x_{0:k-1})}
\tag{77}
$$

from (10). Define

$$
\begin{aligned}
f_{N,k+1}(\omega) &:= p_{\theta_m}(X_{k+1}(\omega)|X_{0:k}(\omega)), \\
f_{D,k+1}(\omega) &:= \sum_{\phi \in \Theta} p_\phi(X_{k+1}(\omega)|X_{0:k}(\omega)) \pi_k^\phi(\omega),
\end{aligned}
\tag{78}
$$

which denote the numerator and the denominator of the coefficient $f_{k+1}(\theta_m, X_{0:k+1}(\omega))$, respectively. Since $\pi_0^{\theta_m} > 0$, we have $0 < f_{D,k+1} \leq 1$ for any $\omega \in \Omega, k \in \mathbb{Z}_+$. Thus

$$
0 \leq f_{D,k+1}|f_{k+1} - 1| \leq |f_{k+1} - 1| \quad \mathbb{P}_{\theta_m}-\text{a.s.}
\tag{79}
$$

for any $k \in \mathbb{Z}_+$. Because $\lim_{k \to \infty} |f_{k+1} - 1| = 0$ $\mathbb{P}_{\theta_m}$-a.s. from (75), the squeeze theorem for (79) yields

$$
\lim_{k \to \infty} f_{D,k+1}|f_{k+1}^{\theta_m,s} - 1| = 0 \quad \mathbb{P}_{\theta_m}-\text{a.s.}
\tag{80}
$$

Now we have

$$
\begin{aligned}
&f_{D,k+1}|f_{k+1}^{\theta_m,s} - 1| \\
&= |f_{N,k+1} - f_{D,k+1}| \\
&= |p_{\theta_m}(X_{k+1}|X_{0:k}) - \sum_{\phi \in \Theta} p_\theta(X_{k+1}|X_{0:k}) \pi_k^\theta| \\
&= |p_{\theta_m}(X_{k+1}|X_{0:k})(1 - \pi_k^{\theta_m}) - p_{\theta_b}(X_{k+1}|X_{0:k}) \pi_k^{\theta_b}| \\
&= |p_{\theta_m}(X_{k+1}|X_{0:k}) - p_{\theta_b}(X_{k+1}|X_{0:k})|(1 - \pi_k^{\theta_m}).
\end{aligned}
\tag{81}
$$

Since $s$ is a PBE with detection-averse utilities, $\lim_{k \to \infty}(1 - \pi_k^{\theta_m}) \neq 0$ $\mathbb{P}_{\theta_m}$-a.s. Therefore, (80) leads to the claim. $\square$

*Proof of Theorem 3:* From the definition of the conditional probability mass function, we have

$$
p_\theta(X_{k+1}|X_{0:k}) = P(X_{k+1}|X_k, s_k^s(\theta, X_{0:k}), s_k^r(X_{0:k})).
\tag{82}
$$

Thus the claim of Lemma 3 can be rewritten by

$$
\begin{aligned}
&|P(X_{k+1}|X_k, A_k^{\theta_m}, s_k^r(X_{0:k})) - P(X_{k+1}|X_k, A_k^{\theta_b}, s_k^r(X_{0:k}))| \\
&\to 0
\end{aligned}
\tag{83}
$$

$\mathbb{P}_{\theta_m}$-almost surely as $k \to \infty$. From finiteness of the MDP, the condition (83) is equivalent to

$$
\mathbb{P}_{\theta_m}(\{E_k \text{ i.o.}\}) = 0
\tag{84}
$$

where

$$
E_k := \left\{P(X_{k+1}|X_k, A_k^{\theta_m}, R_k) \neq P(X_{k+1}|X_k, A_k^{\theta_b}, R_k)\right\}.
\tag{85}
$$

By applying the generalized Borel-Cantelli's second lemma in (2) with $\mathcal{F}_k := \sigma(X_{0:k})$, we have

$$
\mathbb{P}_{\theta_m}(E) = 0
\tag{86}
$$

where

$$
E := \left\{\omega \in \Omega : \sum_{k=0}^\infty \mathbb{P}_{\theta_m}(E_k|\sigma(X_{0:k}))(\omega) = \infty\right\}.
\tag{87}
$$

We derive a simpler description of the event $E$. For any $\omega \in \Omega$, the set of nonnegative integers $\mathbb{Z}_+$ can be divided into two disjoint subsets $\hat{\mathbb{Z}}_+(\omega)$ and $\mathbb{Z}_+ \setminus \hat{\mathbb{Z}}_+(\omega)$ such that

$$
\begin{cases}
A_k^{\theta_m}(\omega) \neq A_k^{\theta_b}(\omega), & \forall k \in \hat{\mathbb{Z}}_+(\omega), \\
A_k^{\theta_m}(\omega) = A_k^{\theta_b}(\omega), & \forall k \in \mathbb{Z}_+ \setminus \hat{\mathbb{Z}}_+(\omega).
\end{cases}
\tag{88}
$$

For a fixed $\omega \in \Omega$, $\mathbb{P}_{\theta_m}(E_k|\sigma(X_{0:k}))(\omega) = 0$ for $k \in \mathbb{Z}_+ \setminus \hat{\mathbb{Z}}_+(\omega)$. Thus we have

$$
\sum_{k=0}^\infty \mathbb{P}_{\theta_m}(E_k|\sigma(X_{0:k}))(\omega) = \sum_{k \in \hat{\mathbb{Z}}_+(\omega)} \mathbb{P}_{\theta_m}(E_k|\sigma(X_{0:k}))(\omega).
\tag{89}
$$

Moreover, for any $k \in \mathbb{Z}_+$ and $\omega \in \Omega$, we have

$$
\mathbb{P}_{\theta_m}(E_k|\sigma(X_{0:k}))(\omega) = \sum_{x_{k+1} \in \mathcal{X}_{k+1}(\omega)} P(x_{k+1}|x_k, a_k^{\theta_m}, r_k)
\tag{90}
$$

where

$$
\mathcal{X}_{k+1}(\omega) := \{x \in \mathcal{X} : P(x|x_k, a_k^{\theta_m}, r_k) \neq P(x|x_k, a_k^{\theta_b}, r_k)\}
\tag{91}
$$

with $x_{0:k} := X_{0:k}(\omega)$, $a_k^\theta := A_k^\theta(\omega)$, and $r_k := R_k(\omega)$. Thus, the condition in the definition of $E$ can be rewritten by

$$
\sum_{k \in \hat{\mathbb{Z}}_+(\omega)} \sum_{x_{k+1} \in \mathcal{X}_{k+1}(\omega)} P(x_{k+1}|x_k, a_k^{\theta_m}, r_k) = \infty.
\tag{92}
$$

Now we define $F := \{A_k^{\theta_m} \neq A_k^{\theta_b} \text{ i.o.}\}$. We show the claim by contradiction. Assume $\mathbb{P}_{\theta_m}(F) > 0$. Then $\mathbb{P}_{\theta_m}(E|F)$ is well-defined. From (86), we have $\mathbb{P}_{\theta_m}(E \cap F) = \mathbb{P}_{\theta_m}(E|F)\mathbb{P}_{\theta_m}(F) = 0$. Because $\mathbb{P}_{\theta_m}(F)$ is assumed to be

This article has been accepted for publication in IEEE Transactions on Automatic Control. This is the author's version which has not been fully edited and content may change prior to final publication. Citation information: DOI 10.1109/TAC.2023.3340978

SASAHARA *et al.*: ASYMPTOTIC SECURITY USING BAYESIAN DEFENSE MECHANISM WITH APPLICATION TO CYBER DECEPTION                                        15

nonzero, this equation implies $\mathbb{P}_{\theta_{\mathrm{m}}}(E|F) = 0$. We now calculate $\mathbb{P}_{\theta_{\mathrm{m}}}(E|F)$ from its definition. Let

$$\gamma(\omega) := \inf_{k \in \hat{\mathbb{Z}}_+(\omega)} \sum_{x_{k+1} \in \mathcal{X}_{k+1}(\omega)} P(x_{k+1}|x_k, a_k^{\theta_{\mathrm{m}}}, r_k). \quad (93)$$

For any $\omega \in \Omega$, the set $\mathcal{X}_{k+1}(\omega)$ is nonempty for $k \in \hat{\mathbb{Z}}_+(\omega)$ from Assumption 1. This fact and the finiteness of the MDP lead to that $\gamma(\omega) > 0$. The infimum leads to the inequality

$$\sum_{k \in \hat{\mathbb{Z}}_+(\omega)} \sum_{x_{k+1} \in \mathcal{X}_{k+1}(\omega)} P(x_{k+1}|x_k, a_k^{\theta_{\mathrm{m}}}, r_k) \geq \left|\hat{\mathbb{Z}}_+(\omega)\right| \gamma(\omega). \quad (94)$$

If $\omega \in F$, then $\hat{\mathbb{Z}}_+(\omega)$ has infinite elements and hence $|\hat{\mathbb{Z}}_+(\omega)|\gamma(\omega) = \infty$. Thus, for any $\omega \in F$, from the inequality (94), the condition (92) holds. Therefore, $\mathbb{P}_{\theta_{\mathrm{m}}}(E|F) = 1$, which is a contradiction. Hence $\mathbb{P}_{\theta_{\mathrm{m}}}(F) = 0$ holds. $\square$

*Proof of Lemma 4:* The former claim is obvious since the probability measures are independent of the strategies with $\theta_{\mathrm{b}}^{\mathrm{s}}$ and $\theta_{\mathrm{u}}^{\mathrm{r}}$. Assume $\hat{s}_{2,k:\infty}^{\mathrm{s}} \in \mathrm{BR}_k^{\mathrm{s}}(\hat{s}_{2,k:\infty}^{\mathrm{r}}|\theta_{\mathrm{m}}^{\mathrm{s}}, x_{0:k})$ for any $k \in \mathbb{Z}_+$ and $x_{0:k} \in \mathcal{X}^{k+1}$. Because $\hat{\pi}_k^{\mathrm{s}}(\theta_{\mathrm{a}}^{\mathrm{r}}|\theta_{\mathrm{m}}^{\mathrm{s}}) = 1$ and $\hat{\pi}_k^{\mathrm{s}}(\theta_{\mathrm{u}}^{\mathrm{r}}|\theta_{\mathrm{m}}^{\mathrm{s}}) = 0$ for $k \in \mathbb{Z}_+$, the malicious sender's expected average utility in $\hat{\mathcal{G}}_2$ is given by

$$\bar{U}_{k,T}^{\mathrm{s}}(\hat{s}_{2,k:k+T}|\theta_{\mathrm{m}}^{\mathrm{s}}, x_{0:k}) := \frac{1}{T+1}$$
$$\times \mathbb{E}^s\left[\sum_{\tau=k}^{k+T} U^{\mathrm{s}}(\theta_{\mathrm{m}}^{\mathrm{s}}, X_k, \hat{s}_{2,\tau}^{\mathrm{s}}(\theta_{\mathrm{m}}^{\mathrm{s}}, X_{0:\tau}), \hat{s}_{2,\tau}^{\mathrm{r}}(\theta_{\mathrm{a}}^{\mathrm{r}}, X_{0:\tau})) \middle| x_{0:k}\right]. \quad (95)$$

which is the malicious sender's utility in $\mathcal{G}_1$. Thus, $\mathrm{BR}_k^{\mathrm{s}}(\hat{s}_{2,k:\infty}^{\mathrm{r}}|\theta_{\mathrm{m}}^{\mathrm{s}}, x_{0:k})$ in $\hat{\mathcal{G}}_2$ is equal to $\mathrm{BR}_k^{\mathrm{s}}(s_{1,k:\infty}^{\mathrm{r}}|\theta_{\mathrm{m}}, x_{0:k})$ in $\mathcal{G}_1$. Hence, $s_{1,k:\infty}^{\mathrm{s}} \in \mathrm{BR}_k^{\mathrm{s}}(s_{1,k:\infty}^{\mathrm{r}}|\theta_{\mathrm{m}}, x_{0:k})$. $\square$

*Proof of Lemma 5:* If $s$ is a passively bluffing strategy, then the distribution of $H_k^{\mathrm{s}}$ is independent of the receiver type. Thus the distribution of $\delta\left(A_k^{\theta_{\mathrm{m}}}, A_k^{\theta_{\mathrm{b}}}\right)$ is also independent of the receiver type. Hence, if (44) holds, the same condition holds for $\theta_{\mathrm{u}}^{\mathrm{r}}$ as well. $\square$

*Proof of Lemma 6:* Take $s = (s^{\mathrm{s}}, s^{\mathrm{r}}) \in \mathcal{S}_{\mathrm{nab}}^*$. Consider $\mathcal{G}_1$ corresponding to $\hat{\mathcal{G}}_2$. Let $s_1 = (s_1^{\mathrm{s}}, s_1^{\mathrm{r}})$ be a strategy profile in $\mathcal{G}_1$ given by (39). From Lemma 4, we have

$$\mathbb{P}_{\theta_{\mathrm{m}}^{\mathrm{s}}}^{s_1}\left(\delta(A_k^{\theta_{\mathrm{m}}}, A_k^{\theta_{\mathrm{b}}}) = 0\right) = \mathbb{P}_{\theta_{\mathrm{m}}^{\mathrm{s}}, \theta_{\mathrm{a}}^{\mathrm{r}}}^{\hat{s}_2}\left(\delta(A_k^{\theta_{\mathrm{m}}}, A_k^{\theta_{\mathrm{b}}}) = 0\right). \quad (96)$$

From the contraposition of Lemma 5, this equation implies that $s_1$ is not asymptotically benign in the sense of the game $\mathcal{G}_1$ since $s \in \mathcal{S}_{\mathrm{nab}}$. Thus, $s_1$ is not a PBE of $\mathcal{G}_1$ from Theorem 3. This means that $s_1^{\mathrm{s}}$ contains a strategy that is not a best response. Because $s \in \mathcal{S}^*$, this means $s_{1,k:\infty}^{\mathrm{s}} \notin \mathrm{BR}_k^{\mathrm{s}}(s_{1,k:\infty}^{\mathrm{r}}|\theta_{\mathrm{m}}, x_{0:k})$ for some $k$. From the contraposition of Lemma 4, $s_{k:\infty}^{\mathrm{s}} \notin \mathrm{BR}_k^{\mathrm{s}}(s_{k:\infty}^{\mathrm{r}}|\theta_{\mathrm{m}}^{\mathrm{s}}, x_{0:k})$, which is equivalent to (50). $\square$

*Proof of Theorem 4:* We prove the existence of $\pi_0^{\mathrm{s}}(\theta_{\mathrm{a}}^{\mathrm{r}}|\theta_{\mathrm{m}}^{\mathrm{s}}) < 1$ such that the contraposition of the condition holds, i.e., if $s$ is not asymptotically benign then $s$ is not a PBE. Let $s \in \mathcal{S}_{\mathrm{nab}}$. If $s \notin \mathcal{S}^*$, $s$ is not a PBE. Thus we suppose $s \in \mathcal{S}_{\mathrm{nab}}^*$. It suffices to show that there exists $\pi_0^{\mathrm{s}}(\theta_{\mathrm{a}}^{\mathrm{r}}|\theta_{\mathrm{m}}^{\mathrm{s}}) < 1$ such that

$$\inf_{s \in \mathcal{S}_{\mathrm{nab}}^*} D(s, g(s)) > 0 \quad (97)$$

where $D(s, \tilde{s}^{\mathrm{s}}) := \bar{U}^{\mathrm{s}}((\tilde{s}^{\mathrm{s}}, s^{\mathrm{r}}), \theta_{\mathrm{m}}^{\mathrm{s}}, \pi^{\mathrm{s}}) - \bar{U}^{\mathrm{s}}(s, \theta_{\mathrm{m}}^{\mathrm{s}}, \pi^{\mathrm{s}})$ and $g$ is given in (52).

From (47), we have $D(s, \tilde{s}^{\mathrm{s}}) = \sum_{\theta^{\mathrm{r}} \in \Theta^{\mathrm{r}}} \pi_0^{\mathrm{s}}(\theta^{\mathrm{r}}|\theta_{\mathrm{m}}^{\mathrm{s}})$. From the definition of $\gamma$ in (53), we have

$$D(s, g(s)) \geq D_{\theta_{\mathrm{u}}^{\mathrm{r}}}(s, g(s))\pi_0^{\mathrm{s}}(\theta_{\mathrm{u}}^{\mathrm{r}}|\theta_{\mathrm{m}}^{\mathrm{s}}) + \gamma\pi_0^{\mathrm{s}}(\theta_{\mathrm{a}}^{\mathrm{r}}|\theta_{\mathrm{m}}^{\mathrm{s}})$$
$$= \gamma + (D_{\theta_{\mathrm{u}}^{\mathrm{r}}}(s, g(s)) - \gamma)\pi_0^{\mathrm{s}}(\theta_{\mathrm{u}}^{\mathrm{r}}|\theta_{\mathrm{m}}^{\mathrm{s}}) \quad (98)$$

Consider the case where $D_{\theta_{\mathrm{u}}^{\mathrm{r}}}(s, g(s)) - \gamma \geq 0$ for any $s \in \mathcal{S}_{\mathrm{nab}}^*$. Then (98) implies that $D(s, g(s)) \geq \gamma$ for any $\pi_0^{\mathrm{s}}(\theta_{\mathrm{a}}^{\mathrm{r}}|\theta_{\mathrm{m}}^{\mathrm{s}}) < 1$. From Assumption 2, this inequality leads to $\inf_{s \in \mathcal{S}_{\mathrm{nab}}^*} D(s, g(s)) \geq \gamma > 0$, which implies (97).

Next, consider the case where $D_{\theta_{\mathrm{u}}^{\mathrm{r}}}(s, g(s)) - \gamma < 0$ for some $s \in \mathcal{S}_{\mathrm{nab}}^*$. By taking an initial belief $\pi_0^{\mathrm{s}}(\theta_{\mathrm{u}}^{\mathrm{r}}|\theta_{\mathrm{m}}^{\mathrm{s}}) > 0$ such that

$$\pi_0^{\mathrm{s}}(\theta_{\mathrm{u}}^{\mathrm{r}}|\theta_{\mathrm{m}}^{\mathrm{s}}) < \inf_{s \in \mathcal{T}} \frac{-\gamma}{D_{\theta_{\mathrm{u}}^{\mathrm{r}}}(s, g(s)) - \gamma} \quad (99)$$

where $\mathcal{T} := \{s \in \mathcal{S}_{\mathrm{nab}}^* : D_{\theta_{\mathrm{u}}^{\mathrm{r}}}(s, g(s)) - \gamma < 0\}$, we have (97) from (98). From the definition of $D_{\theta_{\mathrm{u}}^{\mathrm{r}}}$ in (51), $D_{\theta_{\mathrm{u}}^{\mathrm{r}}}(s, g(s))$ is bounded in $\mathcal{T}$ because $\bar{U}_{\theta^{\mathrm{r}}}^{\mathrm{s}}$ is bounded for any strategy. Thus we have

$$\inf_{s \in \mathcal{T}} \frac{-\gamma}{D_{\theta_{\mathrm{u}}^{\mathrm{r}}}(s, g(s)) - \gamma} > 0. \quad (100)$$

Thus a nonzero initial belief that satisfies (99) exists.

$\square$

## REFERENCES

[1] McAfee, "The hidden costs of cybercrime," Tech. Rep., 2020, [Online]. Available: https://www.mcafee.com/enterprise/en-us/assets/reports/rp-hidden-costs-of-cybercrime.pdf.

[2] Cybersecurity & Infrastructure Security Agency, "Stuxnet malware mitigation," Tech. Rep. ICSA-10-238-01B, 2014, [Online]. Available: https://www.us-cert.gov/ics/advisories/ICSA-10-238-01B.

[3] ——, "HatMan - safety system targeted malware," Tech. Rep. MAR-17-352-01, 2017, [Online]. Available: https://www.us-cert.gov/ics/MAR-17-352-01-HatMan-Safety-System-Targeted-Malware-Update-B.

[4] ——, "Cyber-attack against Ukrainian critical infrastructure," Tech. Rep. IR-ALERT-H-16-056-01, 2018, [Online]. Available: https://www.us-cert.gov/ics/alerts/IR-ALERT-H-16-056-01.

[5] ——, "DarkSide ransomware: Best practices for preventing business disruption from ransomware attacks," Tech. Rep. AA21-131A, 2021, [Online]. Available: https://us-cert.cisa.gov/ncas/alerts/aa21-131a.

[6] N. Falliere, L. O. Murchu, and E. Chien, "W32. Stuxnet Dossier," Symantec, Tech. Rep., 2011.

[7] A. Neyman and S. Sorin, Eds., *Stochastic Games and Applications*. Springer, 2003.

[8] I.-K. Cho and D. M. Kreps, "Signaling Games and Stable Equilibria," *The Quarterly Journal of Economics*, vol. 102, no. 2, pp. 179–221, 1987.

[9] J. Mertens and S. Zamir, "Formulation of Bayesian analysis for games with incomplete information," *International Journal of Game Theory*, vol. 14, pp. 1–29, 1985.

[10] E. Dekel and M. Siniscalchi, "Epistemic game theory," in *Handbook of Game Theory*. Elsevier, 2015, ch. 12, pp. 619–702.

[11] M. S. Lund, B. Solhaug, and K. Stolen, *Model-Driven Risk Analysis*. Springer, 2011.

[12] C. Phillips and L. P. Swiler, "A graph-based system for network-vulnerability analysis," in *Proc. 1998 Workshop on New Security Paradigms*, 1998, pp. 71–79.

[13] B. Schneier, "Attack trees," *Dr. Dobb's Journal*, vol. 24, no. 12, pp. 21–29, 1999.

[14] S. Bistarelli, F. Fioravanti, and P. Peretti, "Defense trees for economic evaluation of security investments," in *Proc. First International Conference on Availability, Reliability and Security*, 2006.

[15] A. Roy, D. S. Kim, and K. S. Trivedi, "Attack countermeasure trees (ACT): towards unifying the constructs of attack and defense trees," *Security and Communication Networks*, vol. 5, no. 8, pp. 929–943, 2012.

[16] E. Miehling, M. Rasouli, and D. Teneketzis, "A POMDP approach to the dynamic defense of large-scale cyber networks," *IEEE Trans. Inf. Forensics Security*, vol. 13, no. 10, pp. 2490–2505, 2018.

[17] S. Chockalingam, W. Pieters, A. Teixeira, and P. Gelder, "Bayesian network models in cyber security: A systematic review," in *Secure IT Systems*, ser. Lecture Notes in Computer Science.   Springer, 2017.

[18] C. Kruegel, D. Mutz, W. Robertson, and F. Valeur, "Bayesian event classification for intrusion detection," in *Proc. 19th Annual Computer Security Applications Conference*, 2003, pp. 14–23.

[19] W. Alhakami, A. ALharbi, S. Bourouis, R. Alroobaea, and N. Bouguila, "Network anomaly intrusion detection using a nonparametric Bayesian approach and feature selection," *IEEE Access*, vol. 7, pp. 52 181–52 190, 2019.

[20] S. A. Zonouz, H. Khurana, W. H. Sanders, and T. M. Yardley, "RRE: A game-theoretic intrusion response and recovery engine," *IEEE Trans. Parallel Distrib. Syst.*, vol. 25, no. 2, pp. 395–406, 2014.

[21] N. Poolsappasit, R. Dewri, and I. Ray, "Dynamic security risk management using Bayesian attack graphs," *IEEE Trans. Dependable Secure Comput.*, vol. 9, no. 1, pp. 61–74, 2012.

[22] H. Sandberg, S. Amin, and K. H. Johansson, "Cyberphysical security in networked control systems: An introduction to the issue," *IEEE Control Systems Magazine*, vol. 35, no. 1, pp. 20–23, 2015.

[23] S. M. Dibaji, M. Pirani, D. B. Flamholz, A. M. Annaswamy, K. H. Johansson, and A. Chakrabortty, "A systems and control perspective of CPS security," *Annual Reviews in Control*, vol. 47, pp. 394–411, 2019.

[24] A. Teixeira, I. Shames, H. Sandberg, and K. H. Johansson, "Revealing stealthy attacks in control systems," in *Proc. 50th Annual Allerton Conference on Communication, Control, and Computing*, 2012, pp. 1806–1813.

[25] J. Milošević, A. Teixeira, T. Tanaka, K. H. Johansson, and H. Sandberg, "Security measure allocation for industrial control systems: Exploiting systematic search techniques and submodularity," *International Journal of Robust and Nonlinear Control*, vol. 30, no. 11, pp. 4278–4302, 2020.

[26] J. Giraldo *et al.*, "A survey of physics-based attack detection in cyber-physical systems," *ACM Comput. Surv.*, vol. 51, no. 4, 2018.

[27] M. Tambe, *Security and Game Theory: Algorithms, Deployed Systems, Lessons Learned*.   Cambridge University Press, 2012.

[28] T. Alpcan and T. Başar, *Network Security: A Decision and Game-Theoretic Approach*.   Cambridge University Press, 2010.

[29] J. Pawlick, E. Colbert, and Q. Zhu, "A game-theoretic taxonomy and survey of defensive deception for cybersecurity and privacy," *ACM Computing Surveys*, vol. 52, no. 4, 2019.

[30] M. O. Sayin and T. Başar, "Deception-as-defense framework for cyber-physical systems," in *Safety, Security and Privacy for Cyber-Physical Systems*, R. M. Ferrari and A. M. H. Teixeira, Eds.   Springer International Publishing, 2021, pp. 287–317.

[31] H. Sasahara and H. Sandberg, "Epistemic signaling games for cyber deception with asymmetric recognition," *IEEE Contr. Syst. Lett.*, vol. 6, pp. 854–859, 2022.

[32] J. Pawlick and Q. Zhu, *Game Theory for Cyber Deception: From Theory to Applications*, ser. Static & Dynamic Game Theory: Foundations & Applications.   Springer, 2021.

[33] J. Pawlick, E. Colbert, and Q. Zhu, "Modeling and analysis of leaky deception using signaling games with evidence," *IEEE Trans. Inf. Forensics Security*, vol. 14, no. 7, pp. 1871–1886, July 2019.

[34] Q. Zhu and Z. Xu, "Secure estimation of CPS with a digital twin," in *Cross-Layer Design for Secure and Resilient Cyber-Physical Systems*.   Springer, 2020, pp. 115–138.

[35] P. Diaconis and D. Freedman, "On the consistency of Bayes estimation," *The Annals of Statistics*, vol. 14, no. 1, pp. 1–26, 1986.

[36] S. Walker, "New approaches to Bayesian consistency," *The Annals of Statistics*, vol. 32, no. 5, pp. 2028–2043, 2004.

[37] P. Eichelsbacher and A. Ganesh, "Bayesian inference for Markov chains," *Journal of Applied Probability*, vol. 39, no. 1, pp. 91–99, 2002.

[38] H. Sasahara, S. Sarıtaş, and H. Sandberg, "Asymptotic security of control systems by covert reaction: Repeated signaling game with undisclosed belief," in *Proc. 59th IEEE Conference on Decision and Control*, 2020.

[39] H. Sasahara and H. Sandberg, "Asymptotic security by model-based incident handlers for Markov decision processes," in *Proc. 60th IEEE Conference on Decision and Control*, 2021.

[40] R. Durrett, *Probability: Theory and Examples*, ser. Cambridge Series in Statistical and Probabilistic Mathematics.   Cambridge University Press, 2019.

[41] A. Rasekh, A. Hassanzadeh, S. Mulchandani, S. Modi, and M. K. Banks, "Smart water networks and cyber security," *Journal of Water Resources Planning and Management*, vol. 142, no. 7, 2016.

[42] E. Creaco, A. Campisano, N. Fontana, G. Marini, P. R. Page, and T. Walski, "Real time control of water distribution newtorks: A state-of-the-art review," *Water Research*, vol. 161, pp. 517–530, 2019.

[43] R. Taormina, S. Galelli, N. O. Tippenhauer, E. Salomons, and A. Ostfeld, "Characterizing cyber-physical attacks on water distribution systems," *Journal of Water Resources Planning and Management*, vol. 143, no. 5, 2017.

[44] P. Chen, L. Desmet, and C. Huygens, "A study on advanced persistent threats," in *Proc. International Conference on Communications and Multimedia Security*, 2014, pp. 63–72.

[45] M. L. Puterman, *Markov Decision Processes: Discrete Stochastic Dynamic Programming*.   John Wiley & Sons, Inc., 1994.
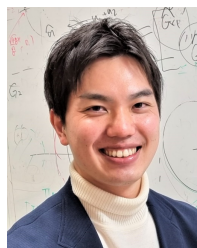
[46] S. Zamir, "Bayesian games: Games with incomplete information," in *Encyclopedia of Complexity and Systems Science*.   Springer, 2009, pp. 426–441.

[47] A. Heifetz, "Commitment," in *Game Theory: Interactive Strategies in Economics and Management*.   Cambridge University Press, 2012, ch. 20.

[48] J. Letchford, D. Korzhyk, and V. Conitzer, "On the value of commitment," *Autonomous Agents and Multi-Agent Systems*, vol. 28, no. 6, pp. 986–1016, 2014.

[49] H. S. Chang and S. I. Marcus, "Two-person zero-sum Markov games: Receding horizon approach," *IEEE Trans. Autom. Control*, vol. 48, no. 11, pp. 1951–1961, 2003.

[50] E. Çinlar, *Probability and Statistics*, ser. Graduate Texts in Mathematics.   Springer, 2011.

**Hampei Sasahara** (M'15) received the Ph.D. degree in engineering from Tokyo Institute of Technology in 2019. He is currently an Assistant Professor with Tokyo Institute of Technology, Tokyo, Japan. From 2019 to 2021, he was a Postdoctoral Scholar with KTH Royal Institute of Technology, Stockholm, Sweden. His main interests include secure control system design and control of large-scale systems.

**Henrik Sandberg** (F'23) received the M.Sc. degree in engineering physics and the Ph.D. degree in automatic control from Lund University, Lund, Sweden, in 1999 and 2004, respectively. He is a Professor with the Division of Decision and Control Systems, KTH Royal Institute of Technology, Stockholm, Sweden. From 2005 to 2007, he was a Postdoctoral Scholar with the California Institute of Technology, Pasadena, CA, USA. In 2013, he was a Visiting Scholar with the Laboratory for Information and Decision Systems, Massachusetts Institute of Technology, Cambridge, MA, USA. He has also held visiting appointments with the Australian National University, Canberra, ACT, USA, and the University of Melbourne, Parkville, VIC, Australia. His current research interests include security of cyber-physical systems, power systems, model reduction, and fundamental limitations in control. Dr. Sandberg received the Best Student Paper Award from the IEEE Conference on Decision and Control in 2004, an Ingvar Carlsson Award from the Swedish Foundation for Strategic Research in 2007, and Consolidator Grant from the Swedish Research Council in 2016. He has served on the editorial boards of IEEE Transactions on Automatic Control and the IFAC Journal Automatica.