

A Game Theoretic Analysis of LQG Control under Adversarial Attack

Zuxing Li, György Dán, and Dong Liu

Abstract—Motivated by recent works addressing adversarial attacks on deep reinforcement learning, a deception attack on linear quadratic Gaussian control is studied in this paper. In the considered attack model, the adversary can manipulate the observation of the agent subject to a mutual information constraint. The adversarial problem is formulated as a novel dynamic cheap talk game to capture the strategic interaction between the adversary and the agent, the asymmetry of information availability, and the system dynamics. Necessary and sufficient conditions are provided for subgame perfect equilibria to exist in pure strategies and in behavioral strategies; and characteristics of the equilibria and the resulting control rewards are given. The results show that pure strategy equilibria are informative, while only babbling equilibria exist in behavioral strategies. Numerical results are shown to illustrate the impact of strategic adversarial interaction.

I. INTRODUCTION

Deep reinforcement learning (DRL) has recently emerged as a promising solution for solving large Markov decision processes (MDPs) and partially observable MDPs (POMDPs), thanks to deep neural networks used as policy approximators [1]. DRL has, however, been shown to be vulnerable to small perturbations of the state observation, called adversarial examples, which were found to mislead the control agent to take suboptimal control actions [2]. While there has been a significant recent interest in the design of adversarial examples against DRL [2], [3], [4], [5], [8], there has been little work on characterizing the ability of agents to adapt to those.

Recent work proposed to use adversarial examples for making DRL agents more robust to perturbations, by letting the adversary and the agent play against each other, and formulating the interaction as a stochastic game (SG) [7]. Nonetheless, in the case of adversarial examples the agent cannot observe the system state directly, neither can the adversary affect the state transition probabilities directly, only through the actions taken by the agent. Hence, the SG model does not capture the information structure of the problem. Effectively, in the presence of adversarial examples the agent has to solve a POMDP, where the observations are subverted by the adversary so as to mislead the agent.

As a model of this interaction, in this paper we propose a game theoretical model to study the strategic interaction between an agent that has to solve a linear quadratic Gaussian (LQG) control problem, and an adversary that can manipulate the agent's observations by a randomly

chosen affine transformation subject to a mutual information constraint, and aims at minimizing the control reward. The resulting problem is formulated as a dynamic cheap talk game, which captures information asymmetry, the beliefs of the adversary and the agent, and the undetectability constraint imposed on the adversarial attacks.

Our paper contributes to the solution of the formulated game theoretical problem in two ways. First, we address necessary and sufficient conditions for the existence of subgame perfect equilibria (SPEs) in pure strategies and in behavioral strategies, and we characterize the equilibrium strategies. Second, we characterize the rewards achievable in equilibria, and relate them to the rewards of a naive agent and an alert agent under attack. The key novelty of our contribution is that we characterize the strategies to be followed by the agent and by the adversary under strategic interaction, which has not been addressed by the existing literature.

The rest of the paper is organized as follows. In Section II we review related work. In Section III we present the system model and problem formulation. In Section IV we provide analytical results. In Section V we provide numerical results. Section VI concludes the paper.

Notation: Unless otherwise specified, we denote a random variable by a capital letter and its realization by the corresponding lower-case letter. We denote by $\mathcal{N}(\cdot, \cdot)$ the Gaussian distribution, by $\mathbb{S}(\cdot)$ the support set, by $I(\cdot; \cdot)$ the mutual information, and by $|\cdot|$ the cardinality of a set.

II. RELATED WORK

Related to our work are previous researches on robust POMDP under uncertainty of the system dynamics [6]. In [6] the control action was optimized under the worst-case assumption of the system dynamics in each stage, i.e., the agent plays as the leader and the dynamic system plays as the follower in a Stackelberg game. SG and partially observable SG (POSG) were used to model the strategic interaction of players in a dynamic system, and have been employed in robust and adversarial problems [7], [8], [9]. But unlike in the case of learning under adversarial attacks, in SG and in POSG the players interact with each other through the impact of their actions on the state transitions, not on the state observations. Our work is related to the cheap talk game [18], where a sender with private information sends a message to a receiver and the receiver takes an action based on the received message and based on its belief on the inaccessible private information. Closest to our model are [10], [11], [12]. In [10] a dynamic cheap talk game was proposed to study a deception attack on a Markovian system, where the actions

This work was partly supported by the Swedish Foundation for Strategic Research (SSF) through the CLAS project and by MSB through the CERCES project.

Z. Li, G. Dán, and D. Liu are with the School of Electrical Engineering and Computer Science, KTH Royal Institute of Technology, Stockholm, Sweden {zuxing;gyuri;doli}@kth.se

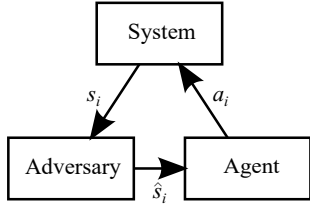


Fig. 1. Considered adversarial attack on dynamic system control. In the i -th stage, the adversary observes the true system state s_i and presents the manipulated state \hat{s}_i to the agent. The agent does not have access to the true state s_i but takes an action a_i upon observing \hat{s}_i .

do not affect the state transitions. In [11] authors developed a dynamic game model of the attacker-defender interaction, and characterized the optimal attack strategy as a function of the defense strategy, allowing for a static optimal defense strategy. In our preliminary work [12] we proposed a dynamic cheap talk framework to model deception attacks on a general MDP, and addressed computational issues.

Adversarial variants of LQG control were considered in a number of recent works. A Stackelberg game was formulated in [13], where the dynamic system is the leader, while the agent is the follower and may be an adversary. The authors formulated a finite horizon hierarchical signaling game between the sender and the receiver in a dynamic environment and showed that linear sender and receiver strategies can yield the equilibrium [14]. In [15], the opposite problem was studied but without considering the complete strategic interaction, where the adversary optimally manipulates the control actions instead of the system states. The optimal attack on both the system state and the control action in LQG control was studied in [16]. In [17], a targeted attack strategy was studied to mislead the LQG system to a particular state while evading detection.

III. ADVERSARIAL LQG CONTROL PROBLEM

We consider an N -stage LQG control problem under adversarial attack, as illustrated in Fig. 1. The system states $\{s_i\}_{i=1}^N$, the manipulated states $\{\hat{s}_i\}_{i=1}^N$, the actions $\{a_i\}_{i=1}^N$, and the instantaneous rewards $\{r_i\}_{i=1}^N$ for $1 \leq i \leq N$ are described by

$$s_{i+1} = \alpha_i s_i + \beta_i a_i + z_i, \text{ given } \alpha_i \neq 0, \beta_i \neq 0; \quad (1)$$

$$\hat{s}_i = \pi_i s_i + c_i; \quad (2)$$

$$a_i = \kappa_i \hat{s}_i + \rho_i; \quad (3)$$

$$r_i = R_i(s_i, a_i) = -\theta_i s_i^2 - \phi_i a_i^2, \text{ given } \theta_i > 0, \phi_i > 0; \quad (4)$$

$$S_1 \sim b_1 \triangleq \mathcal{N}(\mu_1, \sigma_1^2), \text{ given } \mu_1, \sigma_1^2 > 0; \quad (5)$$

$$Z_i \sim \mathcal{N}(0, \omega_i^2), \text{ given } \omega_i^2 > 0; \quad (6)$$

$$C_i \sim \mathcal{N}(0, \delta_i^2). \quad (7)$$

A. LQG Recapitulation

If $\{\pi_i\}_{i=1}^N$ and $\{\delta_i^2\}_{i=1}^N$ are known by the agent, the above problem is a standard LQG control. In the i -th stage, the agent observes \hat{s}_i but not s_i , and determines the action a_i with the aim to maximize its expected accumulated reward.

Note that it is sufficient to consider an affine function of \hat{s}_i for a_i in the standard LQG problem as the optimal action a_i^* is a linear function of the mean of the Gaussian posterior distribution of S_i for the agent after observing $\{\hat{s}_k\}_{k=1}^i$ and $\{a_k\}_{k=1}^{i-1}$ [20]. To compute the optimal coefficients κ_i^* and ρ_i^* , we first define $\tilde{\theta}_{N+1} = 0$, and for $1 \leq i \leq N$ define $\tilde{\theta}_i$ as

$$\tilde{\theta}_i = \theta_i + \tilde{\theta}_{i+1} \alpha_i^2 - \frac{\tilde{\theta}_{i+1}^2 \alpha_i^2 \beta_i^2}{\phi_i + \tilde{\theta}_{i+1} \beta_i^2}. \quad (8)$$

Furthermore, we denote by $b_i \triangleq \mathcal{N}(\mu_i, \sigma_i^2)$ the belief of the agent about S_i , which is the Gaussian posterior distribution of S_i for the agent after observing $\{\hat{s}_k\}_{k=1}^{i-1}$ and $\{a_k\}_{k=1}^{i-1}$. Then, given the manipulated state \hat{s}_i , the optimal action can be expressed as

$$\kappa_i^* = -\frac{\tilde{\theta}_{i+1} \alpha_i \beta_i \pi_i \sigma_i^2}{(\phi_i + \tilde{\theta}_{i+1} \beta_i^2)(\pi_i^2 \sigma_i^2 + \delta_i^2)}, \quad (9)$$

$$\rho_i^* = -\frac{\tilde{\theta}_{i+1} \alpha_i \beta_i \mu_i \delta_i^2}{(\phi_i + \tilde{\theta}_{i+1} \beta_i^2)(\pi_i^2 \sigma_i^2 + \delta_i^2)}, \quad (10)$$

$$a_i^* = \kappa_i^* \hat{s}_i + \rho_i^* = -\frac{\tilde{\theta}_{i+1} \alpha_i \beta_i}{\phi_i + \tilde{\theta}_{i+1} \beta_i^2} \frac{\pi_i \sigma_i^2 \hat{s}_i + \mu_i \delta_i^2}{\pi_i^2 \sigma_i^2 + \delta_i^2}, \quad (11)$$

where the coefficients κ_i^* and ρ_i^* depend on b_i ; $\frac{\pi_i \sigma_i^2 \hat{s}_i + \mu_i \delta_i^2}{\pi_i^2 \sigma_i^2 + \delta_i^2}$ is the mean of the Gaussian posterior distribution of S_i for the agent after observing $\{\hat{s}_k\}_{k=1}^i$ and $\{a_k\}_{k=1}^{i-1}$. Note that when $\pi_i \equiv 1$ and $\delta_i^2 \equiv 0$ the LQG strategy reduces to the linear quadratic regulator (LQR) strategy

$$(\kappa_i^*, \rho_i^*) = \left(-\frac{\tilde{\theta}_{i+1} \alpha_i \beta_i}{\phi_i + \tilde{\theta}_{i+1} \beta_i^2}, 0 \right). \quad (12)$$

B. Adversarial Model

The adversary can manipulate the observation of the agent and its objective is to minimize the agent's expected accumulated reward, similar to [2], [5]. In the i -th stage, the adversary chooses the manipulation parameters π_i , δ_i^2 , manipulates the state s_i to \hat{s}_i , and then reports the manipulated state \hat{s}_i to the agent. We consider that the adversarial manipulation is "small", since a large manipulation may be easily detected and may also involve a high manipulation cost. Given the agent's belief $b_i \triangleq \mathcal{N}(\mu_i, \sigma_i^2)$, we impose the following constraints on the manipulation:

$$-\infty < \varepsilon' \leq \pi_i \leq \varepsilon < \infty; \quad (13)$$

$$I(\hat{S}_i; S_i) = \frac{1}{2} \log \frac{\pi_i^2 \sigma_i^2 + \delta_i^2}{\delta_i^2} \geq \frac{1}{2} \log \lambda > 0, \forall \pi_i, \delta_i^2, \quad (14)$$

i.e., $\frac{\pi_i^2 \sigma_i^2 + \delta_i^2}{\delta_i^2} \geq \lambda > 1$. The mutual information constraint (14) implies that the manipulated state conveys at least a certain amount of information about the system state to the agent. A larger value of λ means a weaker adversary, and vice versa. Note that in order to satisfy the mutual information constraint, the adversary *cannot* use $\pi_i = 0$. We denote by $\mathcal{A}_i(b_i, \varepsilon', \varepsilon, \lambda)$ the set of feasible adversarial actions (π_i, δ_i^2) in the i -th stage subject to (13)-(14). Finally, in the end of

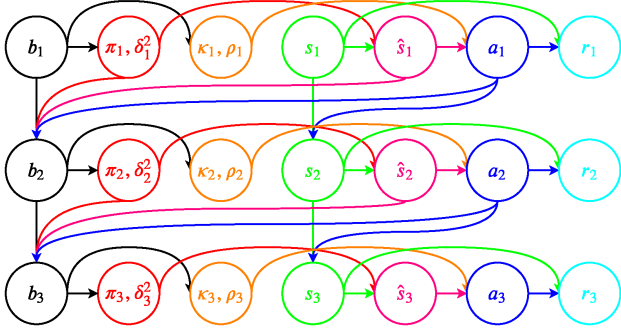


Fig. 2. Illustration of information structure for a three-stage ALQG game.

this stage, the adversary reveals π_i, δ_i^2 to the agent¹, so as to keep the adversarial model consistent with the standard LQG control.

C. Adversarial LQG Control Game

In every stage of the adversarial problem, there is a cheap talk interaction, where the adversary acts as the sender and the agent as the receiver. Different from the dynamic cheap talk game with an action-independent Markovian system [10], we propose a novel dynamic cheap talk game to model the strategic interaction of the adversary and the agent with asymmetric information in the adversarial LQG control problem. Unlike in recent works on adversarial reinforcement learning [2], [3], [4], [5], in our model of strategic interaction the agent is aware of and can adapt to the adversary.

The game is played between the adversary and the agent over N stages. In the i -th stage, the belief of the agent b_i is known to the adversary and determines the action set $\mathcal{A}_i(b_i, \epsilon', \epsilon, \lambda)$. The adversary uses a behavioral strategy $g_i(\pi_i, \delta_i^2 | b_i)$ over $\mathcal{A}_i(b_i, \epsilon', \epsilon, \lambda)$ for choosing (π_i, δ_i^2) . Then, given the observed system state s_i it generates the manipulated state \hat{s}_i with the probability measure $\mathcal{N}(\hat{s}_i | \pi_i s_i, \delta_i^2)$. The agent uses a pure strategy² $f_i(b_i)$ for choosing κ_i and ρ_i based on the belief b_i , and takes the action $a_i = \kappa_i \hat{s}_i + \rho_i$ once it receives the manipulated state \hat{s}_i . Finally, the players compute the belief $b_{i+1} \triangleq \mathcal{N}(\mu_{i+1}, \sigma_{i+1}^2)$ based on the current belief b_i , the coefficient π_i , the variance δ_i^2 , the manipulated state \hat{s}_i , and the action a_i as $\mu_{i+1} = \Lambda_\mu(b_i, \pi_i, \delta_i^2, \hat{s}_i, a_i)$ and $\sigma_{i+1}^2 = \Lambda_\nu(b_i, \pi_i, \delta_i^2)$.

We can thus express the expected accumulated agent reward using the adversarial strategies $g^N \triangleq (g_1, \dots, g_N)$ and the agent's strategies $f^N \triangleq (f_1, \dots, f_N)$ over N stages as

$$V(b_1, g^N, f^N) = E_{b_1, g^N, f^N} \left(\sum_{j=1}^N R_j(S_j, A_j) \right). \quad (15)$$

Consequently, the objective of the adversary is to minimize (15), while the agent aims at maximizing it. We refer to

¹This assumption is strong but may hold in some cases. For instance, the player in the shell game reveals the cup in which the pellet is after each round.

²The following analysis will show that it is sufficient to consider a pure agent strategy with an affine form.

this particular dynamic cheap talk game as adversarial LQG (ALQG) game. Fig. 2 illustrates a three-stage ALQG game. Our objective is to characterize SPEs in the ALQG game: the existence conditions and solution structures.

Remark 1: The adversarial LQG problem cannot be modeled as an SG or a POSG, since the adversary directly manipulates the observation of the agent. Furthermore, different from zero-sum SG, an SPE of the ALQG game does not necessarily exist. On the other hand, for $N = 1$ the game is a cheap talk game [18], where the strategies of both players depend on the belief of the agent and on the constraints on the adversarial manipulation. Nonetheless, in the ALQG game the reward function is different from that in [18], which gives rise to different equilibria, as we will show later.

IV. EQUILIBRIUM ANALYSIS

In the following, we first formulate the value function and the belief update rule for the adversarial LQG problem; we then characterize SPEs in pure strategies and in behavioral strategies, respectively.

A. Value Function and Belief Update

Assume that there is an SPE consisting of strategies (g^{N*}, f^{N*}) . Induced by this SPE, we can define the value function of a subgame starting from the i -th stage as

$$V_i^N(b_i) = V(b_i, g_i^{N*}, f_i^{N*}) = E_{b_i, g_i^{N*}, f_i^{N*}} \left(\sum_{j=i}^N R_j(S_j, A_j) \right), \quad (16)$$

i.e., the value function $V_i^N(b_i)$ is the expected accumulated agent reward in the subgame starting from the i -th stage when the belief in the i -th stage is b_i and the SPE strategies (g_i^{N*}, f_i^{N*}) are used. The evaluation of $V_i^N(b_i)$ needs the beliefs in the subgame. In the following, we specify the belief update rule.

Given the current belief $b_i \triangleq \mathcal{N}(\mu_i, \sigma_i^2)$, the coefficient π_i , the variance δ_i^2 , the manipulated state \hat{s}_i , and the action a_i , it follows from the adversarial LQG model and Bayes rule that $b_{i+1} \triangleq \mathcal{N}(\mu_{i+1}, \sigma_{i+1}^2)$ with

$$\mu_{i+1} = \Lambda_\mu(b_i, \pi_i, \delta_i^2, \hat{s}_i, a_i) = \alpha_i \frac{\pi_i \sigma_i^2 \hat{s}_i + \mu_i \delta_i^2}{\pi_i^2 \sigma_i^2 + \delta_i^2} + \beta_i a_i, \quad (17)$$

$$\sigma_{i+1}^2 = \Lambda_\nu(b_i, \pi_i, \delta_i^2) = \frac{\alpha_i^2 \sigma_i^2 \delta_i^2}{\pi_i^2 \sigma_i^2 + \delta_i^2} + \omega_i^2. \quad (18)$$

An immediate consequence of the belief update rule is the following.

Property 1: It follows from $\sigma_1^2 > 0$, the variance update rule (18), and $\pi_i \neq 0$ that $\sigma_i^2 > 0$ for all $1 \leq i \leq N$.

Observe that the value functions $\{V_i^N\}_{i=1}^{N-1}$ have to satisfy the backward dynamic programming equation

$$\begin{aligned} V_i^N(b_i) &= \min_{g_i} E_{b_i, g_i, f_i^*} \{R_i(S_i, A_i) \\ &\quad + V_{i+1}^N(\mathcal{N}(\Lambda_\mu(b_i, \Pi_i, \Delta_i^2, \hat{S}_i, A_i), \Lambda_\nu(b_i, \Pi_i, \Delta_i^2)))\} \\ &= \max_{f_i} E_{b_i, g_i^*, f_i} \{R_i(S_i, A_i) \\ &\quad + V_{i+1}^N(\mathcal{N}(\Lambda_\mu(b_i, \Pi_i, \Delta_i^2, \hat{S}_i, A_i), \Lambda_\nu(b_i, \Pi_i, \Delta_i^2)))\}. \end{aligned} \quad (19)$$

Thus the SPE has to satisfy (19), which is the basis for the analysis we present in the following.

B. Pure Strategy Equilibria

We start the analysis considering pure strategy equilibria. With slight abuse of notation, we denote by $(\pi_i, \delta_i^2) = g_i(b_i)$ a pure strategy of the adversary as a function of the belief b_i .

We first consider the case $N = 1$.

Proposition 1: Let $N = 1$. An SPE consists of (f_1^*, g_1^*) , where $(\kappa_1^*, \rho_1^*) = f_1^*(b_1) = (0, 0)$ for any belief b_1 ; and g_1^* can be any adversarial strategy defined on $\mathcal{A}_1(b_1, \varepsilon', \varepsilon, \lambda)$.

Proof: Observe that $(\kappa_1^*, \rho_1^*) = f_1^*(b_1) = (0, 0)$ is a dominant strategy for the agent for any belief b_1 . Under this strategy, the adversarial strategy has no impact on the agent's reward. This proves the result. ■

The existence of a pure strategy SPE for $N = 1$ is encouraging, even if the equilibrium is degenerate. Unfortunately, for $N \geq 2$ an SPE may not exist as shown in the following theorem.

Theorem 1: Let $N \geq 2$. If $\varepsilon' \neq \varepsilon$ or if $\varepsilon' = \varepsilon = 0$, then there is no pure strategy SPE for the ALQG game. If $\varepsilon' = \varepsilon \neq 0$, then there is a unique pure strategy SPE. The SPE strategies for $1 \leq i \leq N$ are given by

$$\tilde{\theta}_{N+1} = \hat{\theta}_{N+1} = 0; \quad (20)$$

$$\tilde{\theta}_i = \theta_i + \tilde{\theta}_{i+1} \alpha_i^2 - \frac{\tilde{\theta}_{i+1}^2 \alpha_i^2 \beta_i^2}{\phi_i + \tilde{\theta}_{i+1} \beta_i^2}; \quad (21)$$

$$\hat{\theta}_i = \theta_i + \hat{\theta}_{i+1} \alpha_i^2 - \left(\frac{\tilde{\theta}_{i+1}^2 \alpha_i^2 \beta_i^2}{\phi_i + \tilde{\theta}_{i+1} \beta_i^2} + (\hat{\theta}_{i+1} - \tilde{\theta}_{i+1}) \alpha_i^2 \right) \frac{\lambda - 1}{\lambda}; \quad (22)$$

$$\pi_i^* = g_i^*(b_i) = \varepsilon' = \varepsilon; \quad (23)$$

$$\delta_i^{2*} = g_i^*(b_i) = \frac{\varepsilon^2 \sigma_i^2}{\lambda - 1}; \quad (24)$$

$$\kappa_i^* = f_i^*(b_i) = -\frac{\tilde{\theta}_{i+1} \alpha_i \beta_i (\lambda - 1)}{(\phi_i + \tilde{\theta}_{i+1} \beta_i^2) \lambda \varepsilon}; \quad (25)$$

$$\rho_i^* = f_i^*(b_i) = -\frac{\tilde{\theta}_{i+1} \alpha_i \beta_i \mu_i}{(\phi_i + \tilde{\theta}_{i+1} \beta_i^2) \lambda}. \quad (26)$$

Corollary 1: If $\varepsilon' = \varepsilon \neq 0$, the value function induced by the unique pure strategy SPE is

$$V_i^N(b_i) = -\tilde{\theta}_i \mu_i^2 - \hat{\theta}_i \sigma_i^2 - \sum_{j=i+1}^N \hat{\theta}_j \omega_{j-1}^2. \quad (27)$$

The proofs of Theorem 1 and Corollary 1 are provided in the appendix.

Property 2: It follows from (20)-(22) that for $1 \leq i \leq N$,

$$\hat{\theta}_i \geq \tilde{\theta}_i > 0. \quad (28)$$

Remark 2: We can make the following observations based on Theorem 1 and Corollary 1.

- Since the LQG control can be seen as the best response of any given (adversarial manipulated) observation model, it is sufficient to consider a pure agent strategy in the form of an affine function for an SPE of the ALQG game with a pure adversarial strategy.

- It follows from (24) that a rational adversary will always apply a manipulation with the largest variance.
- The value function V_i^N consists of a constant term and two separable terms depending on the mean μ_i and the variance σ_i^2 of the belief, which allows a closed form solution for arbitrary N .

Time-invariant system: We now turn to the asymptotic analysis of a time-invariant system, i.e., $\alpha_i = \alpha \neq 0$, $\beta_i = \beta \neq 0$, $\omega_i^2 = \omega^2 > 0$, $\theta_i = \theta > 0$, and $\phi_i = \phi > 0$ for $i \geq 1$, and we let $N \rightarrow \infty$. Let us define the mapping $L: \mathbb{R}_{\geq 0}^2 \rightarrow \mathbb{R}_{\geq 0}^2$ as

$$L(x, y) = \left(\theta + \frac{\phi \alpha^2 x}{\phi + \beta^2 x}, \theta + \frac{\phi \alpha^2 x}{\phi + \beta^2 x} \frac{\lambda - 1}{\lambda} + \alpha^2 y \frac{1}{\lambda} \right). \quad (29)$$

Observe that L is effectively the coefficient update (20)-(22) for the time-invariant model. In what follows, we first characterize L and furthermore the pure strategy SPE of the ALQG game in the asymptotic regime.

Proposition 2: Let $\lambda > \alpha^2$. Then the mapping L admits a least fixed point $(\tilde{\theta}, \hat{\theta}) \in \mathbb{R}_{\geq 0}^2$, for which

$$\lim_{n \rightarrow \infty} L^n(0, 0) = L(\tilde{\theta}, \hat{\theta}) = (\tilde{\theta}, \hat{\theta}), \quad (30)$$

with $L^n(0, 0) \triangleq \underbrace{L(L(\dots(L(L(0, 0))))}_{n \text{ L-mappings}} \dots$.

Proof: We start the proof by observing that the mapping L is order-preserving. That is, for all $(x, y), (x', y') \in \mathbb{R}_{\geq 0}^2$ satisfying $(x, y) \preceq (x', y')$, i.e., $x \leq x'$ and $y \leq y'$, we have $L(x, y) \preceq L(x', y')$. This can be easily shown by analyzing (29).

Furthermore, the fixed point equation $L(\tilde{\theta}, \hat{\theta}) = (\tilde{\theta}, \hat{\theta})$ has a unique solution on $\mathbb{R}_{\geq 0}^2$ if and only if $\lambda > \alpha^2$. Since L is order-preserving, the convergence result $\lim_{n \rightarrow \infty} L^n(0, 0) = (\tilde{\theta}, \hat{\theta})$ follows from Kleene's fixed point theorem [21]. ■

Analytical expressions for $\tilde{\theta}$ and $\hat{\theta}$ can be obtained by solving the fixed point equation $L(\tilde{\theta}, \hat{\theta}) = (\tilde{\theta}, \hat{\theta})$, and can be used for characterizing the SPE in pure strategies, using Theorem 1 and Proposition 2, as follows.

Theorem 2: Let $\lambda > \alpha^2$, $\varepsilon' = \varepsilon \neq 0$, and $N \rightarrow \infty$. Then the ALQG game of the time-invariant model has a stationary SPE in pure strategies as: For $i \geq 1$,

$$\pi_i^* = g_i^*(b_i) = \varepsilon; \quad (31)$$

$$\delta_i^{2*} = g_i^*(b_i) = \frac{\varepsilon^2 \sigma_i^2}{\lambda - 1}; \quad (32)$$

$$\kappa_i^* = f_i^*(b_i) = -\frac{\tilde{\theta} \alpha \beta (\lambda - 1)}{(\phi + \tilde{\theta} \beta^2) \lambda \varepsilon}; \quad (33)$$

$$\rho_i^* = f_i^*(b_i) = -\frac{\tilde{\theta} \alpha \beta \mu_i}{(\phi + \tilde{\theta} \beta^2) \lambda}. \quad (34)$$

Proof: Since $\varepsilon' = \varepsilon \neq 0$, a unique SPE in pure strategies exists and the SPE strategies are given in Theorem 1. Since $\lambda > \alpha^2$ and $N \rightarrow \infty$, it follows from Proposition 2 that $\tilde{\theta}_i, \hat{\theta}_i$ converge to $\tilde{\theta}, \hat{\theta}$, respectively. This leads to the stationary SPE in Theorem 2. ■

Interestingly, for this stationary SPE in pure strategies we can obtain the expected average reward per stage in steady state in closed form.

Corollary 2: Let $b_1 \triangleq \mathcal{N}(\mu_1, \sigma_1^2)$ with bounded mean and variance. For the stationary SPE in pure strategies in Theorem 2, the expected average reward per stage in steady state is independent of the initial belief and is given by

$$\lim_{N \rightarrow \infty} \frac{V_1^N(b_1)}{N} = -\hat{\theta} \omega^2. \quad (35)$$

Proof: This result follows from Corollary 1, Proposition 2, and Theorem 2. ■

C. Equilibria in Behavioral Strategies

The previous results show that a pure strategy SPE does not exist if there are multiple choices of the coefficient π_i for the adversary. We thus turn to the analysis of SPE in behavioral strategies.

Theorem 3: Let $N \geq 2$ and $\varepsilon' < 0 < \varepsilon$. Then for $1 \leq i \leq N$, $\tilde{\theta}_i$ and $\check{\theta}_i$ are as given by (21), and

$$\tilde{\theta}_{N+1} = \check{\theta}_{N+1} = 0, \quad (36)$$

$$\check{\theta}_i = \theta_i + \check{\theta}_{i+1} \alpha_i^2 - (\check{\theta}_{i+1} - \tilde{\theta}_{i+1}) \alpha_i^2 \frac{\lambda - 1}{\lambda}. \quad (37)$$

Furthermore, there is a continuum of SPEs in behavioral strategies. Each SPE in the i -th stage consists of a behavioral strategy of the adversary and a pure strategy of the agent that satisfies

$$\mathbb{S}(g_i^* | b_i) \triangleq \left\{ (\pi_i, \delta_i^2) : \delta_i^2 = \frac{\pi_i^2 \sigma_i^2}{\lambda - 1} \right\} \subseteq \mathcal{A}_i(b_i, \varepsilon', \varepsilon, \lambda); \quad (38)$$

$$\|\mathbb{S}(g_i^* | b_i)\| \geq 2; \quad (39)$$

$$E_{g_i^*}(\Pi_i) = 0; \quad (40)$$

$$(\kappa_i^*, \rho_i^*) = f_i^*(b_i) = \left(0, -\frac{\tilde{\theta}_{i+1} \alpha_i \beta_i \mu_i}{\phi_i + \tilde{\theta}_{i+1} \beta_i^2} \right). \quad (41)$$

Interestingly, this SPE is a babbling equilibrium in which the agent's action is based on its belief, not on the manipulated state. Nonetheless, the value of the game depends on the adversarial manipulation, as shown in the following corollary.

Corollary 3: Let $\varepsilon' < 0 < \varepsilon$. For any SPE in behavioral strategies, we have

$$V_i^N(b_i) = -\tilde{\theta}_i \mu_i^2 - \check{\theta}_i \sigma_i^2 - \sum_{j=i+1}^N \check{\theta}_j \omega_{j-1}^2. \quad (42)$$

The proofs of Theorem 3 and Corollary 3 are given in the appendix.

Property 3: It follows from the update rules (20)-(22) and (36)-(37) that for all $1 \leq i \leq N$,

$$\check{\theta}_i \geq \hat{\theta}_i \geq \tilde{\theta}_i > 0. \quad (43)$$

Remark 3: We can make the following observations based on Theorem 3 and Corollary 3.

- It is sufficient for the agent to use a pure affine strategy against a behavioral strategy of the adversary.
- Although the adversary cannot use $\pi_i = 0$, the behavioral strategy g_i^* needs to achieve zero-mean of the random coefficient Π_i .
- A rational adversary will always use a manipulation with the largest variance.
- From Property 3, the value (42) of an SPE in behavioral strategies is always less than or equal to the value (27) of a pure strategy SPE.

Time-invariant system: We again turn to the time-invariant system for $N \rightarrow \infty$. Let us define the mapping $J: \mathbb{R}_{\geq 0}^2 \rightarrow \mathbb{R}_{\geq 0}^2$ as

$$J(x, y) = \left(\theta + \frac{\phi \alpha^2 x}{\phi + \beta^2 x}, \theta + \alpha^2 x \frac{\lambda - 1}{\lambda} + \alpha^2 y \frac{1}{\lambda} \right). \quad (44)$$

Observe that J is effectively the coefficient update rule (21), (36), (37) for the time-invariant system. In what follows, we characterize J and the stationary SPEs in behavioral strategies for the ALQG game.

Proposition 3: Let $\lambda > \alpha^2$. Then the mapping J admits a least fixed point $(\tilde{\theta}, \check{\theta}) \in \mathbb{R}_{\geq 0}^2$, for which

$$\lim_{n \rightarrow \infty} J^n(0, 0) = J(\tilde{\theta}, \check{\theta}) = (\tilde{\theta}, \check{\theta}). \quad (45)$$

The proof of Proposition 3 follows using the arguments in the proof of Proposition 2.

Theorem 4: Let $\lambda > \alpha^2$, $\varepsilon' < 0 < \varepsilon$, and $N \rightarrow \infty$. Then the ALQG game of the time-invariant model has a stationary SPE in behavioral strategies as: For $i \geq 1$,

$$g_i^* \left(\pi_i = \varepsilon', \delta_i^2 = \frac{\varepsilon'^2 \sigma_i^2}{\lambda - 1} \middle| b_i \right) = \frac{\varepsilon}{\varepsilon - \varepsilon'}; \quad (46)$$

$$g_i^* \left(\pi_i = \varepsilon, \delta_i^2 = \frac{\varepsilon^2 \sigma_i^2}{\lambda - 1} \middle| b_i \right) = -\frac{\varepsilon'}{\varepsilon - \varepsilon'}; \quad (47)$$

$$(\kappa_i^*, \rho_i^*) = f_i^*(b_i) = \left(0, -\frac{\tilde{\theta} \alpha \beta \mu_i}{\phi + \tilde{\theta} \beta^2} \right). \quad (48)$$

Corollary 4: Let $b_1 \triangleq \mathcal{N}(\mu_1, \sigma_1^2)$ with bounded mean and variance. For the stationary SPE in behavioral strategies in Theorem 4, the expected average reward per stage in steady state is independent of the initial belief and is given by

$$\lim_{N \rightarrow \infty} \frac{V_1^N(b_1)}{N} = -\check{\theta} \omega^2. \quad (49)$$

The proofs of Theorem 4 and Corollary 4 are based on Theorem 3, Corollary 3, and Proposition 3, and follow from similar arguments as used in the proofs of Theorem 2 and Corollary 2. Observe that Corollary 2, Corollary 4, and Property 3 jointly imply that the expected average agent reward per stage in steady state is higher when considering pure strategies, as behavioral strategies allow for more uncertainty about the attack and thus make the adversary stronger.

The behavioral strategy of the adversary in Theorem 3 needs to achieve zero-mean of the random coefficient Π_i ,

which *cannot* be satisfied if $0 \leq \varepsilon' < \varepsilon$ or $\varepsilon' < \varepsilon \leq 0$. In the following, we study SPE under these conditions.

Theorem 5: Let $N \geq 2$. If $0 = \varepsilon' < \varepsilon$ or if $\varepsilon' < \varepsilon = 0$, there is no SPE for the ALQG game.

Theorem 6: Let $N = 2$. If $0 < \varepsilon' < \varepsilon$ or if $\varepsilon' < \varepsilon < 0$, there is a unique SPE in behavioral strategies for the ALQG game: For any belief $b_1 \triangleq \mathcal{N}(\mu_1, \sigma_1^2)$,

$$g_1^* \left(\pi_1 = \varepsilon', \delta_1^2 = \frac{\varepsilon'^2 \sigma_1^2}{\lambda - 1} \middle| b_1 \right) = \frac{\varepsilon}{\varepsilon' + \varepsilon}; \quad (50)$$

$$g_1^* \left(\pi_1 = \varepsilon, \delta_1^2 = \frac{\varepsilon^2 \sigma_1^2}{\lambda - 1} \middle| b_1 \right) = \frac{\varepsilon'}{\varepsilon' + \varepsilon}; \quad (51)$$

$$\kappa_1^* = f_1^*(b_1) \quad (52)$$

$$= \frac{-\theta_2 \alpha_1 \beta_1 E_{g_1^*}(\Pi_1) \sigma_1^2}{(\phi_1 + \theta_2 \beta_1^2) \left(E_{g_1^*}(\Pi_1^2) \left(\mu_1^2 + \frac{\lambda}{\lambda-1} \sigma_1^2 \right) - E_{g_1^*}^2(\Pi_1) \mu_1^2 \right)}; \quad (53)$$

$$\rho_1^* = f_1^*(b_1) = -E_{g_1^*}(\Pi_1) \mu_1 \kappa_1^* - \frac{\theta_2 \alpha_1 \beta_1 \mu_1}{\phi_1 + \theta_2 \beta_1^2}; \quad (54)$$

for any belief $b_2 \triangleq \mathcal{N}(\mu_2, \sigma_2^2)$, $g_2^* = g_1^*$; $(\kappa_2^*, \rho_2^*) = f_2^*(b_2) = (0, 0)$; and

$$\begin{aligned} V_1^2(b_1) &= - \left(\theta_1 + \theta_2 \alpha_1^2 - \frac{\theta_2^2 \alpha_1^2 \beta_1^2}{\phi_1 + \theta_2 \beta_1^2} \right) \mu_1^2 - \theta_2 \omega_1^2 - (\theta_1 + \theta_2 \alpha_1^2) \sigma_1^2 \\ &+ \frac{\theta_2^2 \alpha_1^2 \beta_1^2 E_{g_1^*}^2(\Pi_1) \sigma_1^4}{(\phi_1 + \theta_2 \beta_1^2) \left(E_{g_1^*}(\Pi_1^2) \left(\mu_1^2 + \frac{\lambda}{\lambda-1} \sigma_1^2 \right) - E_{g_1^*}^2(\Pi_1) \mu_1^2 \right)}. \end{aligned} \quad (55)$$

The proofs of Theorems 5 & 6 are given in the appendix.

Remark 4: Different from the value functions (27) and (42), the value function (55) *cannot* be decomposed into separable terms of mean and variance, which makes the extension of Theorem 6 to $N \geq 3$ difficult.

V. NUMERICAL RESULTS

We illustrate the impact of strategic interaction on a time-invariant adversarial LQG problem. The parameters used for the evaluation are shown in Table I.

We start with illustrating the convergence of the mappings in Propositions 2 and 3. Fig. 3 shows the mappings $L^n(0, 0)$ and $J^n(0, 0)$ as functions of the iteration n for $\lambda = 1.5$ and $\lambda = 2$. Both mappings increase monotonically starting from $(0, 0)$ and converge to their least fixed points $(\tilde{\theta}, \tilde{\theta})$ and $(\hat{\theta}, \hat{\theta})$, respectively, which confirms these propositions.

We continue with the evaluation of stationary SPEs for the time-invariant system. Fig. 4 shows the expected average

TABLE I
LQG MODEL DEFAULT PARAMETERS

Parameter	μ_1	σ_1^2	α	β	ω^2	θ	ϕ
Value	0	1	-0.5	-1.5	1	2	1

reward per stage for a stationary SPE in pure strategies and that for a stationary SPE in behavioral strategies. Observe that the adversarial manipulation capability decreases as the mutual information lower bound λ increases, and hence the expected average rewards increase. The results also confirm Property 3, i.e., the expected average reward per stage for a stationary SPE in behavioral strategies *cannot* be higher than that for a stationary SPE in pure strategies.

Next, we assess the importance of strategic interaction on the agent's performance, as compared to an agent that is unaware of the attack [2], [5]. Fig. 5 shows the expected average reward per stage for a stationary SPE in behavioral strategies, as per Corollary 4, that for a naive agent under an optimal adversarial attack, and that for an alert agent under SPE adversarial behavioral strategy g_i^* in Corollary 4, for $-\varepsilon' = \varepsilon > 0$. The naive agent is unaware of the adversarial manipulation, i.e., it uses the optimal LQR strategy (12). The corresponding optimal stationary adversarial strategy can be obtained through dynamic programming, and is the pure strategy:

$$(\pi_i, \delta_i^2) = g_i^A(b_i) = \left(-\varepsilon, \frac{\varepsilon^2 \sigma_i^2}{\lambda - 1} \right).$$

The alert agent suspects an adversary but does not act strategically despite the presence of an adversary. The alert agent assumes $\hat{\pi}_i = \hat{\pi} \neq 0$, $\hat{\delta}_i^2 = \frac{\hat{\pi}^2 \sigma_i^2}{\lambda - 1}$, and uses the corresponding best response strategy

$$(\kappa_i, \rho_i) = f_i^C(b_i) = \left(-\frac{\tilde{\theta} \alpha \beta (\lambda - 1)}{(\phi + \tilde{\theta} \beta^2) \lambda \hat{\pi}}, -\frac{\tilde{\theta} \alpha \beta \mu_i}{(\phi + \tilde{\theta} \beta^2) \lambda} \right).$$

The figure shows that the expected average reward per stage of the naive agent is always lower than that for the SPE. Clearly, as the mutual information lower bound λ increases, the adversarial manipulation capability becomes weaker and the expected average rewards per stage increase. At the same time, we can observe that if the bound ε of the manipulation coefficient is higher then the adversarial manipulation capability becomes stronger and therefore the expected average reward per stage of the naive agent decreases. Note that the limit of the expected average reward per stage of the naive agent does not exist when ε is larger than a threshold due to the resulting instability of the control system. The poor performance of the naive agent is consistent with recent works on adversarial DRL [2], [5], where the naive DRL agents were found to perform poorly against strategic adversaries. The results for the SPE show, however, that an agent that is aware of the adversary can adjust its strategy to be resilient to adversarial attack. The figure also shows that the expected average reward per stage of the alert agent is always lower than that for the SPE since the alert agent does not adjust its best response strategically to the SPE adversarial strategy. Different from the SPE and the naive agent, it is interesting to observe that the alert agent's performance deteriorates as the adversarial constraint λ increases. This is due to that the alert agent's strategy f_i^C deviates more from the SPE strategy f_i^* , as shown in Corollary 4, as λ increases.

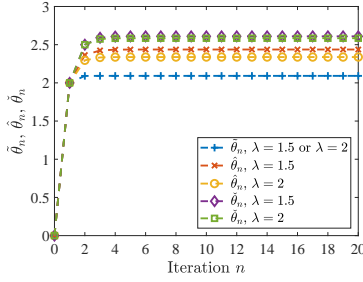


Fig. 3. $\tilde{\theta}_n$, $\hat{\theta}_n$, θ_n computed as $L^n(0,0)$ and $J^n(0,0)$ v.s. the number of iterations n , for $\lambda = 1.5$ and $\lambda = 2$ ($\lambda > \alpha^2$).

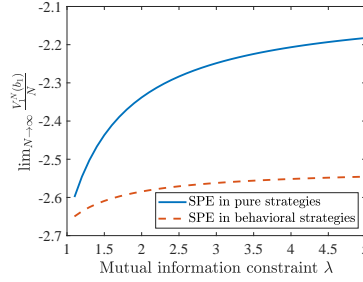


Fig. 4. Expected average reward per stage v.s. mutual information constraint λ , for stationary SPEs in pure strategies and behavioral strategies.

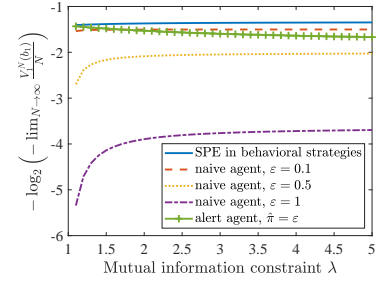


Fig. 5. Expected average reward per stage for stationary SPE in behavioral strategies, that for a naive agent, and that for an alert agent v.s. mutual information constraint λ .

VI. CONCLUSION

We proposed a game theoretic model to capture the strategic interaction, information asymmetry and system dynamics for LQG control under adversarial input subject to a mutual information constraint. We characterized the subgame perfect equilibria in pure strategies and in behavioral strategies, including stationary equilibria for time-invariant systems. Our results show that if an equilibrium exists then the agent can use an affine, pure strategy, but randomization enables the adversary to construct more powerful attacks, under a wider range of parameters, and forces the agent into a babbling equilibrium. Our numerical results show the importance of strategic interaction for LQG control, and highlight that an agent that is aware of an adversarial attack can be designed resilient. Our work could be extended in a number of interesting directions, including considering a non-scalar state dynamic system, and relaxing the assumption that the adversarial strategy is revealed to the agent after each stage.

APPENDIX

A. Proofs of Theorem 1 and Corollary 1

Proof: We prove the result using backward dynamic programming. Recall that there is no feasible adversarial strategy for $\varepsilon' = \varepsilon = 0$. Thus, it is sufficient to consider the cases $\varepsilon' \neq \varepsilon$ or $\varepsilon' = \varepsilon \neq 0$. In stage N the value function for b_N can be expressed as

$$V_N^N(b_N) = -\tilde{\theta}_N \mu_N^2 - \hat{\theta}_N \sigma_N^2, \quad (56)$$

where $\tilde{\theta}_N = \hat{\theta}_N = \theta_N$ from the update rules (20)-(22). The pure strategies, which form an SPE and achieve the value function, consist of any pure strategy g_N^* satisfying the adversarial constraints (13)-(14), and the dominant pure strategy

$$(\kappa_N^*, \rho_N^*) = f_N^*(b_N) = (0, 0). \quad (57)$$

In stage $N-1$, the Q-function of using $\pi_{N-1} \neq 0$, δ_{N-1}^2 ,

κ_{N-1} , and ρ_{N-1} given a belief b_{N-1} is

$$\begin{aligned} Q_{N-1}^N(b_{N-1}, \pi_{N-1}, \delta_{N-1}^2, \kappa_{N-1}, \rho_{N-1}) &= -\theta_{N-1} E(S_{N-1}^2) - \phi_{N-1} E(A_{N-1}^2) \\ &\quad - \tilde{\theta}_N E(\Lambda_\mu(b_{N-1}, \pi_{N-1}, \delta_{N-1}^2, \hat{S}_{N-1}, A_{N-1}))^2 \\ &\quad - \hat{\theta}_N \Lambda_V(b_{N-1}, \pi_{N-1}, \delta_{N-1}^2) \\ &= -(\theta_{N-1} + \tilde{\theta}_N \alpha_{N-1}^2) \mu_{N-1}^2 - \hat{\theta}_N \omega_{N-1}^2 \\ &\quad - (\theta_{N-1} + \hat{\theta}_N \alpha_{N-1}^2) \sigma_{N-1}^2 \\ &\quad - (\phi_{N-1} + \tilde{\theta}_N \beta_{N-1}^2) (\pi_{N-1} \kappa_{N-1} \mu_{N-1} + \rho_{N-1})^2 \\ &\quad - 2\tilde{\theta}_N \alpha_{N-1} \beta_{N-1} \mu_{N-1} (\pi_{N-1} \kappa_{N-1} \mu_{N-1} + \rho_{N-1}) \\ &\quad + (\hat{\theta}_N - \tilde{\theta}_N) \alpha_{N-1}^2 \frac{\pi_{N-1}^2 \sigma_{N-1}^2}{\pi_{N-1}^2 \sigma_{N-1}^2 + \delta_{N-1}^2} \sigma_{N-1}^2 \\ &\quad - (\phi_{N-1} + \tilde{\theta}_N \beta_{N-1}^2) \kappa_{N-1}^2 (\pi_{N-1}^2 \sigma_{N-1}^2 + \delta_{N-1}^2) \\ &\quad - 2\tilde{\theta}_N \alpha_{N-1} \beta_{N-1} \pi_{N-1} \kappa_{N-1} \sigma_{N-1}^2, \end{aligned} \quad (58)$$

where the expectations are induced by the given b_{N-1} , π_{N-1} , δ_{N-1}^2 , κ_{N-1} , and ρ_{N-1} .

Given b_{N-1} , $\pi_{N-1} \neq 0$, δ_{N-1}^2 , and κ_{N-1} , the Q-function Q_{N-1}^N is a concave quadratic function of ρ_{N-1} . As the best response to maximize the agent reward, we can substitute ρ_{N-1} in terms of b_{N-1} , π_{N-1} , and κ_{N-1} as

$$\rho_{N-1} = -\pi_{N-1} \kappa_{N-1} \mu_{N-1} - \frac{\tilde{\theta}_N \alpha_{N-1} \beta_{N-1}}{\phi_{N-1} + \tilde{\theta}_N \beta_{N-1}^2} \mu_{N-1}. \quad (59)$$

Thus, it is sufficient to consider the Q-function

$$\begin{aligned} Q_{N-1}^N(b_{N-1}, \pi_{N-1}, \delta_{N-1}^2, \kappa_{N-1}) &= - \left(\theta_{N-1} + \tilde{\theta}_N \alpha_{N-1}^2 - \frac{\tilde{\theta}_N^2 \alpha_{N-1}^2 \beta_{N-1}^2}{\phi_{N-1} + \tilde{\theta}_N \beta_{N-1}^2} \right) \mu_{N-1}^2 \\ &\quad - (\theta_{N-1} + \hat{\theta}_N \alpha_{N-1}^2) \sigma_{N-1}^2 - \hat{\theta}_N \omega_{N-1}^2 \\ &\quad + (\hat{\theta}_N - \tilde{\theta}_N) \alpha_{N-1}^2 \frac{\pi_{N-1}^2 \sigma_{N-1}^2}{\pi_{N-1}^2 \sigma_{N-1}^2 + \delta_{N-1}^2} \sigma_{N-1}^2 \\ &\quad - (\phi_{N-1} + \tilde{\theta}_N \beta_{N-1}^2) \kappa_{N-1}^2 (\pi_{N-1}^2 \sigma_{N-1}^2 + \delta_{N-1}^2) \\ &\quad - 2\tilde{\theta}_N \alpha_{N-1} \beta_{N-1} \pi_{N-1} \kappa_{N-1} \sigma_{N-1}^2. \end{aligned} \quad (60)$$

As shown in Property 1, the belief variance is always positive. Therefore, given b_{N-1} , $\pi_{N-1} \neq 0$, and δ_{N-1}^2 , the Q-function Q_{N-1}^N is a concave quadratic function of κ_{N-1} ,

and the best response of the agent in terms of κ_{N-1} for b_{N-1} , π_{N-1} , and δ_{N-1}^2 can be expressed as

$$\kappa_{N-1} = -\frac{\tilde{\theta}_N \alpha_{N-1} \beta_{N-1} \pi_{N-1} \sigma_{N-1}^2}{(\phi_{N-1} + \tilde{\theta}_N \beta_{N-1}^2)(\pi_{N-1}^2 \sigma_{N-1}^2 + \delta_{N-1}^2)} \neq 0. \quad (61)$$

Given b_{N-1} , $\pi_{N-1} \neq 0$, and $\kappa_{N-1} \neq 0$, the Q-function Q_{N-1}^N is a decreasing function of δ_{N-1}^2 . The best response of the adversary in terms of δ_{N-1}^2 for b_{N-1} and π_{N-1} is thus

$$\delta_{N-1}^2 = \frac{\pi_{N-1}^2 \sigma_{N-1}^2}{\lambda - 1}. \quad (62)$$

Consequently, it is sufficient to consider the Q-function

$$\begin{aligned} & Q_{N-1}^N(b_{N-1}, \pi_{N-1}, \kappa_{N-1}) \\ &= -\left(\theta_{N-1} + \tilde{\theta}_N \alpha_{N-1}^2 - \frac{\tilde{\theta}_N^2 \alpha_{N-1}^2 \beta_{N-1}^2}{\phi_{N-1} + \tilde{\theta}_N \beta_{N-1}^2} \right) \mu_{N-1}^2 \\ &\quad - (\theta_{N-1} + \hat{\theta}_N \alpha_{N-1}^2) \sigma_{N-1}^2 - \hat{\theta}_N \omega_{N-1}^2 \\ &\quad + (\hat{\theta}_N - \tilde{\theta}_N) \alpha_{N-1}^2 \frac{\lambda - 1}{\lambda} \sigma_{N-1}^2 \\ &\quad - (\phi_{N-1} + \tilde{\theta}_N \beta_{N-1}^2) \kappa_{N-1}^2 \frac{\lambda}{\lambda - 1} \pi_{N-1}^2 \sigma_{N-1}^2 \\ &\quad - 2\tilde{\theta}_N \alpha_{N-1} \beta_{N-1} \pi_{N-1} \kappa_{N-1} \sigma_{N-1}^2. \end{aligned} \quad (63)$$

The pure strategies (g_{N-1}^*, f_{N-1}^*) form an SPE if $\pi_{N-1}^* = g_{N-1}^*(b_{N-1})$ and $\kappa_{N-1}^* = f_{N-1}^*(b_{N-1})$ satisfy

$$\pi_{N-1}^* = \arg \min_{\varepsilon' \leq \pi_{N-1} \leq \varepsilon, \pi_{N-1} \neq 0} Q_{N-1}^N(b_{N-1}, \pi_{N-1}, \kappa_{N-1}^*), \quad (64)$$

$$\kappa_{N-1}^* = \arg \max_{\kappa_{N-1} \in \mathbb{R}} Q_{N-1}^N(b_{N-1}, \pi_{N-1}^*, \kappa_{N-1}). \quad (65)$$

If $\varepsilon' = \varepsilon \neq 0$, $\pi_{N-1}^* = \varepsilon = \varepsilon' = g_{N-1}^*(b_{N-1})$ is a dominant adversarial strategy. Therefore, the SPE must exist. The pure strategies (g_{N-1}^*, f_{N-1}^*) can be obtained by substituting $\pi_{N-1}^* = \varepsilon$ into (59), (61), and (62) as

$$\begin{aligned} \delta_{N-1}^{2*} &= g_{N-1}^*(b_{N-1}) = \frac{\varepsilon^2 \sigma_{N-1}^2}{\lambda - 1}; \\ \kappa_{N-1}^* &= f_{N-1}^*(b_{N-1}) = -\frac{\tilde{\theta}_N \alpha_{N-1} \beta_{N-1} (\lambda - 1)}{(\phi_{N-1} + \tilde{\theta}_N \beta_{N-1}^2) \lambda \varepsilon}; \\ \rho_{N-1}^* &= f_{N-1}^*(b_{N-1}) = -\frac{\tilde{\theta}_N \alpha_{N-1} \beta_{N-1} \mu_{N-1}}{(\phi_{N-1} + \tilde{\theta}_N \beta_{N-1}^2) \lambda}. \end{aligned}$$

The value function $V_{N-1}^N(b_{N-1})$ can be obtained by substituting π_{N-1}^* and κ_{N-1}^* into the Q-function (63) as

$$V_{N-1}^N(b_{N-1}) = -\tilde{\theta}_N \mu_{N-1}^2 - \hat{\theta}_N \sigma_{N-1}^2 - \hat{\theta}_N \omega_{N-1}^2.$$

Let us now consider the case $\varepsilon' \neq \varepsilon$. Assume that there exists an SPE with $\pi_{N-1}^* = g_{N-1}^*(b_{N-1}) \neq 0$. As the best response, solving (65) leads to $\kappa_{N-1}^* = -\frac{\tilde{\theta}_N \alpha_{N-1} \beta_{N-1} (\lambda - 1)}{(\phi_{N-1} + \tilde{\theta}_N \beta_{N-1}^2) \lambda \pi_{N-1}^*}$. For all $\pi_{N-1} \neq \pi_{N-1}^*$ and $\pi_{N-1} \neq 0$ we have

$$Q_{N-1}^N(b_{N-1}, \pi_{N-1}^*, \kappa_{N-1}^*) > Q_{N-1}^N(b_{N-1}, \pi_{N-1}, \kappa_{N-1}^*).$$

Thus, condition (64) cannot hold and hence the assumption is not true, i.e., there is no pure strategy SPE in this case.

In the case of $\varepsilon' = \varepsilon \neq 0$, Theorem 1 and Corollary 1 can be justified in the remaining stages of the backward dynamic programming by using the same analysis. ■

B. Proofs of Theorem 3 and Corollary 3

Proof: We prove the result by verifying that the given strategies form an SPE. The dominant pure strategy of the agent and the value function in the final stage are as shown in the proofs of Theorem 1 and Corollary 1. Note that any behavioral adversarial strategy satisfying (13)-(14) can be g_N^* since it has no impact on the agent reward. Therefore, Theorem 3 and Corollary 3 hold in the final stage.

For stage $N-1$, we first show that it is sufficient to consider a pure agent strategy with an affine form. A general behavioral agent strategy f_{N-1} decides an action a_{N-1} based on the belief b_{N-1} and the observation \hat{s}_{N-1} with the probability measure $f_{N-1}(a_{N-1} | b_{N-1}, \hat{s}_{N-1})$. Given a belief b_{N-1} , a behavioral adversarial strategy g_{N-1} , an observation \hat{s}_{N-1} , and an action a_{N-1} from the support set of a behavioral agent strategy f_{N-1} , the Q-function is

$$\begin{aligned} & Q_{N-1}^N(b_{N-1}, g_{N-1}, \hat{s}_{N-1}, a_{N-1}) \\ &= -\theta_{N-1} E_{b_{N-1}}(S_{N-1}^2) - \phi_{N-1} a_{N-1}^2 \\ &\quad - \tilde{\theta}_N E_{g_{N-1}}(\Lambda_\mu(b_{N-1}, \Pi_{N-1}, \Delta_{N-1}^2, \hat{s}_{N-1}, a_{N-1}))^2 \\ &\quad - \check{\theta}_N E_{g_{N-1}}(\Lambda_\nu(b_{N-1}, \Pi_{N-1}, \Delta_{N-1}^2)) \\ &= -\theta_{N-1} (\mu_{N-1}^2 + \sigma_{N-1}^2) - (\phi_{N-1} + \tilde{\theta}_N \beta_{N-1}^2) a_{N-1}^2 \\ &\quad - 2\tilde{\theta}_N E_{g_{N-1}} \left(\frac{\Pi_{N-1} \sigma_{N-1}^2 \hat{s}_{N-1} + \mu_{N-1} \Delta_{N-1}^2}{\Pi_{N-1}^2 \sigma_{N-1}^2 + \Delta_{N-1}^2} \right) \\ &\quad \alpha_{N-1} \beta_{N-1} a_{N-1} \\ &\quad - \tilde{\theta}_N E_{g_{N-1}} \left(\frac{\Pi_{N-1} \sigma_{N-1}^2 \hat{s}_{N-1} + \mu_{N-1} \Delta_{N-1}^2}{\Pi_{N-1}^2 \sigma_{N-1}^2 + \Delta_{N-1}^2} \right)^2 \alpha_{N-1}^2 \\ &\quad - \check{\theta}_N E_{g_{N-1}} \left(\frac{\alpha_{N-1}^2 \sigma_{N-1}^2 \Delta_{N-1}^2}{\Pi_{N-1}^2 \sigma_{N-1}^2 + \Delta_{N-1}^2} \right) - \check{\theta}_N \omega_{N-1}^2, \end{aligned} \quad (66)$$

which is a concave quadratic function of a_{N-1} . As the best response to maximize the agent reward, the support set of the behavioral agent strategy is a singleton, i.e., it is sufficient to use a pure agent strategy, which has an affine form as

$$\begin{aligned} a_{N-1} &= -\frac{\tilde{\theta}_N \alpha_{N-1} \beta_{N-1} E_{g_{N-1}} \left(\frac{\Pi_{N-1} \sigma_{N-1}^2}{\Pi_{N-1}^2 \sigma_{N-1}^2 + \Delta_{N-1}^2} \right)}{\phi_{N-1} + \tilde{\theta}_N \beta_{N-1}^2} \hat{s}_{N-1} \\ &\quad - \frac{\tilde{\theta}_N \alpha_{N-1} \beta_{N-1} E_{g_{N-1}} \left(\frac{\Delta_{N-1}^2}{\Pi_{N-1}^2 \sigma_{N-1}^2 + \Delta_{N-1}^2} \right)}{\phi_{N-1} + \tilde{\theta}_N \beta_{N-1}^2} \mu_{N-1}. \end{aligned} \quad (67)$$

Assume that an SPE in behavioral strategies consists of $\kappa_{N-1}^* = f_{N-1}^*(b_{N-1}) = 0$ and $\rho_{N-1}^* = f_{N-1}^*(b_{N-1}) = -\frac{\tilde{\theta}_N \alpha_{N-1} \beta_{N-1} \mu_{N-1}}{\phi_{N-1} + \tilde{\theta}_N \beta_{N-1}^2}$. Given b_{N-1} , $(\pi_{N-1}, \delta_{N-1}^2)$ in the support set of a behavioral adversarial strategy g_{N-1} , κ_{N-1}^* , and ρ_{N-1}^* ,

we have the following Q-function:

$$\begin{aligned}
& Q_{N-1}^N(b_{N-1}, \pi_{N-1}, \delta_{N-1}^2, \kappa_{N-1}^*, \rho_{N-1}^*) \\
&= - \left(\theta_{N-1} + \check{\theta}_N \alpha_{N-1}^2 - \frac{\check{\theta}_N^2 \alpha_{N-1}^2 \beta_{N-1}^2}{\phi_{N-1} + \check{\theta}_N \beta_{N-1}^2} \right) \mu_{N-1}^2 \\
&\quad - (\theta_{N-1} + \check{\theta}_N \alpha_{N-1}^2) \sigma_{N-1}^2 - \check{\theta}_N \omega_{N-1}^2 \\
&\quad + (\check{\theta}_N - \check{\theta}_N) \alpha_{N-1}^2 \frac{\pi_{N-1}^2 \sigma_{N-1}^2}{\pi_{N-1}^2 \sigma_{N-1}^2 + \delta_{N-1}^2} \sigma_{N-1}^2. \quad (68)
\end{aligned}$$

From Property 3 and the adversarial constraints (13)-(14), we have

$$\begin{aligned}
& \left(\pi_{N-1} \neq 0, \delta_{N-1}^2 = \frac{\pi_{N-1}^2 \sigma_{N-1}^2}{\lambda - 1} \right) \\
&= \arg \min_{(\pi_{N-1}, \delta_{N-1}^2)} Q_{N-1}^N(b_{N-1}, \pi_{N-1}, \delta_{N-1}^2, \kappa_{N-1}^*, \rho_{N-1}^*). \quad (69)
\end{aligned}$$

Therefore, any behavioral adversarial strategy is the best response of f_{N-1}^* if its support set consists of two or more elements of $\left(\pi_{N-1} \neq 0, \delta_{N-1}^2 = \frac{\pi_{N-1}^2 \sigma_{N-1}^2}{\lambda - 1} \right)$.

Assume that an SPE consists of a behavioral adversarial strategy $g_{N-1}^*(\cdot | b_{N-1})$, which is defined on a support set containing two or more elements of $\left(\pi_{N-1} \neq 0, \delta_{N-1}^2 = \frac{\pi_{N-1}^2 \sigma_{N-1}^2}{\lambda - 1} \right)$, and satisfies $E_{g_{N-1}^*}(\Pi_{N-1}) = 0$. Given b_{N-1} , g_{N-1}^* , κ_{N-1} , and ρ_{N-1} , we have the following Q-function:

$$\begin{aligned}
& Q_{N-1}^N(b_{N-1}, g_{N-1}^*, \kappa_{N-1}, \rho_{N-1}) \\
&= - (\theta_{N-1} + \check{\theta}_N \alpha_{N-1}^2) \mu_{N-1}^2 - \check{\theta}_N \omega_{N-1}^2 \\
&\quad - \left(\theta_{N-1} + \check{\theta}_N \alpha_{N-1}^2 - (\check{\theta}_N - \check{\theta}_N) \alpha_{N-1}^2 \frac{\lambda - 1}{\lambda} \right) \sigma_{N-1}^2 \\
&\quad - (\phi_{N-1} + \check{\theta}_N \beta_{N-1}^2) E_{g_{N-1}^*}(\Pi_{N-1}^2) \\
&\quad \left(\mu_{N-1}^2 + \frac{\lambda}{\lambda - 1} \sigma_{N-1}^2 \right) \kappa_{N-1}^2 \\
&\quad - (\phi_{N-1} + \check{\theta}_N \beta_{N-1}^2) \rho_{N-1}^2 \\
&\quad - 2 \check{\theta}_N \alpha_{N-1} \beta_{N-1} \mu_{N-1} \rho_{N-1}. \quad (70)
\end{aligned}$$

The best response of g_{N-1}^* is

$$\begin{aligned}
& \left(\kappa_{N-1} = 0, \rho_{N-1} = - \frac{\check{\theta}_N \alpha_{N-1} \beta_{N-1} \mu_{N-1}}{\phi_{N-1} + \check{\theta}_N \beta_{N-1}^2} \right) \\
&= \arg \max_{(\kappa_{N-1}, \rho_{N-1})} Q_{N-1}^N(b_{N-1}, g_{N-1}^*, \kappa_{N-1}, \rho_{N-1}). \quad (71)
\end{aligned}$$

It follows from (69) and (71) that a behavioral adversarial strategy $g_{N-1}^*(\cdot | b_{N-1})$, which is defined on a support set containing two or more elements of $\left(\pi_{N-1} \neq 0, \delta_{N-1}^2 = \frac{\pi_{N-1}^2 \sigma_{N-1}^2}{\lambda - 1} \right)$ and satisfies $E_{g_{N-1}^*}(\Pi_{N-1}) = 0$, and a pure agent strategy $(\kappa_{N-1}^*, \rho_{N-1}^*) = f_{N-1}^*(b_{N-1}) = \left(0, - \frac{\check{\theta}_N \alpha_{N-1} \beta_{N-1} \mu_{N-1}}{\phi_{N-1} + \check{\theta}_N \beta_{N-1}^2} \right)$ form an SPE in stage $N - 1$. Furthermore, the value function in this stage is

$$V_{N-1}^N(b_{N-1}) = -\check{\theta}_{N-1} \mu_{N-1}^2 - \check{\theta}_{N-1} \sigma_{N-1}^2 - \check{\theta}_N \omega_{N-1}^2.$$

Thus, Theorem 3 and Corollary 3 hold in stage $N - 1$.

In the remaining stages of the backward dynamic programming we can always justify that the solution of the SPE in behavioral strategies from Theorem 3 and the value function from Corollary 3 hold following the same analysis as used in stage $N - 1$. ■

C. Proofs of Theorems 5 and 6

Proof: To prove Theorems 5 and 6, it is sufficient to consider a two-stage problem, i.e., $N = 2$. The solution of the final stage is the same as in the proofs of Theorem 3 and Corollary 3, and therefore is omitted here. Theorems 5 and 6 hold in the final stage. Furthermore, as shown in the proofs of Theorem 3 and Corollary 3, it is sufficient to consider a pure agent strategy with an affine form in the first stage.

Given b_1 , (π_1, δ_1^2) in the support set of a behavioral adversarial strategy g_1 , κ_1 , and ρ_1 , the Q-function is

$$\begin{aligned}
& Q_1^2(b_1, \pi_1, \delta_1^2, \kappa_1, \rho_1) \\
&= - (\theta_1 + \theta_2 \alpha_1^2) \mu_1^2 - \theta_2 \omega_1^2 - (\theta_1 + \theta_2 \alpha_1^2) \sigma_1^2 \\
&\quad - (\phi_1 + \theta_2 \beta_1^2) (\pi_1 \kappa_1 \mu_1 + \rho_1)^2 - 2 \theta_2 \alpha_1 \beta_1 \mu_1 (\pi_1 \kappa_1 \mu_1 + \rho_1) \\
&\quad - (\phi_1 + \theta_2 \beta_1^2) \kappa_1^2 (\pi_1^2 \sigma_1^2 + \delta_1^2) - 2 \theta_2 \alpha_1 \beta_1 \pi_1 \kappa_1 \sigma_1^2. \quad (72)
\end{aligned}$$

This Q-function is non-increasing in δ_1^2 for any given b_1 , π_1 , κ_1 , and ρ_1 . As the best response to minimize the agent reward, it is sufficient to consider a behavioral adversarial strategy defined on a non-singleton support set of $\left(\pi_1 \neq 0, \delta_1^2 = \frac{\pi_1^2 \sigma_1^2}{\lambda - 1} \right)$.

Given b_1 , g_1 with a non-singleton support set of $\left(\pi_1 \neq 0, \delta_1^2 = \frac{\pi_1^2 \sigma_1^2}{\lambda - 1} \right)$, κ_1 , and ρ_1 , the Q-function is

$$\begin{aligned}
& Q_1^2(b_1, g_1, \kappa_1, \rho_1) \\
&= - (\theta_1 + \theta_2 \alpha_1^2) \mu_1^2 - \theta_2 \omega_1^2 - (\theta_1 + \theta_2 \alpha_1^2) \sigma_1^2 \\
&\quad - (\phi_1 + \theta_2 \beta_1^2) E_{g_1}(\Pi_1^2) \left(\mu_1^2 + \frac{\lambda}{\lambda - 1} \sigma_1^2 \right) \kappa_1^2 \\
&\quad - 2 (\phi_1 + \theta_2 \beta_1^2) E_{g_1}(\Pi_1) \mu_1 \kappa_1 \rho_1 - (\phi_1 + \theta_2 \beta_1^2) \rho_1^2 \\
&\quad - 2 \theta_2 \alpha_1 \beta_1 E_{g_1}(\Pi_1) (\mu_1^2 + \sigma_1^2) \kappa_1 - 2 \theta_2 \alpha_1 \beta_1 \mu_1 \rho_1. \quad (73)
\end{aligned}$$

This is a concave quadratic function of ρ_1 when b_1 , g_1 , and κ_1 are fixed. As the best response to maximize the agent reward, we can substitute ρ_1 with

$$\rho_1 = -E_{g_1}(\Pi_1) \mu_1 \kappa_1 - \frac{\theta_2 \alpha_1 \beta_1 \mu_1}{\phi_1 + \theta_2 \beta_1^2}. \quad (74)$$

Then the Q-function (73) reduces to

$$\begin{aligned}
& Q_1^2(b_1, g_1, \kappa_1) \\
&= - \left(\theta_1 + \theta_2 \alpha_1^2 - \frac{\theta_2^2 \alpha_1^2 \beta_1^2}{\phi_1 + \theta_2 \beta_1^2} \right) \mu_1^2 - \theta_2 \omega_1^2 - (\theta_1 + \theta_2 \alpha_1^2) \sigma_1^2 \\
&\quad - (\phi_1 + \theta_2 \beta_1^2) \\
&\quad \left(E_{g_1}(\Pi_1^2) \left(\mu_1^2 + \frac{\lambda}{\lambda - 1} \sigma_1^2 \right) - E_{g_1}^2(\Pi_1) \mu_1^2 \right) \kappa_1^2 \\
&\quad - 2 \theta_2 \alpha_1 \beta_1 E_{g_1}(\Pi_1) \sigma_1^2 \kappa_1. \quad (75)
\end{aligned}$$

This is also a concave quadratic function of κ_1 when b_1 and g_1 are fixed. As the best response to maximize the agent reward, we can substitute κ_1 with

$$\kappa_1 = \frac{-\theta_2 \alpha_1 \beta_1 E_{g_1}(\Pi_1) \sigma_1^2}{(\phi_1 + \theta_2 \beta_1^2) \left(E_{g_1}(\Pi_1^2) \left(\mu_1^2 + \frac{\lambda}{\lambda-1} \sigma_1^2 \right) - E_{g_1}^2(\Pi_1) \mu_1^2 \right)}. \quad (76)$$

Since we consider $0 \leq \varepsilon' < \varepsilon$ or $\varepsilon' < \varepsilon \leq 0$ and the behavioral adversarial strategy has a non-singleton support set, $E_{g_1}(\Pi_1) \neq 0$ and $\kappa_1 \neq 0$ in these cases.

We then study the support set of a behavioral adversarial strategy. Given b_1 , $(\pi_1 \neq 0, \delta_1^2 = \frac{\pi_1^2 \sigma_1^2}{\lambda-1})$ in the support set of a behavioral adversarial strategy g_1 , $\kappa_1 \neq 0$, and ρ_1 , the Q-function is

$$\begin{aligned} Q_1^2(b_1, \pi_1, \kappa_1, \rho_1) &= -(\theta_1 + \theta_2 \alpha_1^2) \mu_1^2 - \theta_2 \omega_1^2 - (\theta_1 + \theta_2 \alpha_1^2) \sigma_1^2 \\ &\quad - (\phi_1 + \theta_2 \beta_1^2) (\pi_1 \kappa_1 \mu_1 + \rho_1)^2 - 2\theta_2 \alpha_1 \beta_1 \mu_1 (\pi_1 \kappa_1 \mu_1 + \rho_1) \\ &\quad - (\phi_1 + \theta_2 \beta_1^2) \kappa_1^2 \frac{\lambda}{\lambda-1} \pi_1^2 \sigma_1^2 - 2\theta_2 \alpha_1 \beta_1 \pi_1 \kappa_1 \sigma_1^2. \end{aligned} \quad (77)$$

This is a concave quadratic function of π_1 given b_1 , $\kappa_1 \neq 0$, and ρ_1 . As the best response to minimize the agent reward, it is sufficient to consider a behavioral adversarial strategy g_1 with the following support set: $\left\{ \left(\pi_1 = \varepsilon', \delta_1^2 = \frac{\varepsilon'^2 \sigma_1^2}{\lambda-1} \right), \left(\pi_1 = \varepsilon, \delta_1^2 = \frac{\varepsilon^2 \sigma_1^2}{\lambda-1} \right) \right\}$.

When $0 = \varepsilon' < \varepsilon$ or $\varepsilon' < \varepsilon = 0$, an SPE in behavioral strategies does not exist since $0 = \varepsilon'$ or $\varepsilon = 0$ will lead to a singleton support set of the behavioral adversarial strategy; and meanwhile a pure strategy SPE does not exist since $\varepsilon' \neq \varepsilon$. This proves Theorem 5.

When $0 < \varepsilon' < \varepsilon$ or $\varepsilon' < \varepsilon < 0$, we assume that an SPE in behavioral strategies consists of

$$\begin{aligned} g_1^* \left(\pi_1 = \varepsilon', \delta_1^2 = \frac{\varepsilon'^2 \sigma_1^2}{\lambda-1} \middle| b_1 \right) &= p^*; \\ g_1^* \left(\pi_1 = \varepsilon, \delta_1^2 = \frac{\varepsilon^2 \sigma_1^2}{\lambda-1} \middle| b_1 \right) &= 1 - p^*; \end{aligned}$$

$$\kappa_1^* = f_1^*(b_1)$$

$$\begin{aligned} &= \frac{-\theta_2 \alpha_1 \beta_1 E_{g_1^*}(\Pi_1) \sigma_1^2}{(\phi_1 + \theta_2 \beta_1^2) \left(E_{g_1^*}(\Pi_1^2) \left(\mu_1^2 + \frac{\lambda}{\lambda-1} \sigma_1^2 \right) - E_{g_1^*}^2(\Pi_1) \mu_1^2 \right)}; \\ \rho_1^* &= f_1^*(b_1) = -E_{g_1^*}(\Pi_1) \mu_1 \kappa_1^* - \frac{\theta_2 \alpha_1 \beta_1 \mu_1}{\phi_1 + \theta_2 \beta_1^2}. \end{aligned}$$

Since the assumed f_1^* is the best response of the assumed g_1^* , we only need to testify if $0 < p^* < 1$ exists such that g_1^* is the best response of f_1^* , i.e., both $\pi_1 = \varepsilon'$ and $\pi_1 = \varepsilon$ are minimizers of the Q-function $Q_1^2(b_1, \pi_1, \kappa_1^*, \rho_1^*)$. There is a unique solution

$$p^* = \frac{\varepsilon}{\varepsilon' + \varepsilon}. \quad (78)$$

Therefore, there is a unique SPE in behavioral strategies for the two-stage ALQG game. The strategies and the value function of the SPE in Theorem 6 can then be obtained easily. ■

REFERENCES

- [1] V. Mnih et al., "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529-533, 2015.
- [2] S. Huang, N. Papernot, I. Goodfellow, Y. Duan, and P. Abbeel, "Adversarial attacks on neural network policies," arXiv:1702.02284.
- [3] Y.-C. Lin, Z.-W. Hong, Y.-H. Liao, M.-L. Shih, M.-Y. Liu, and S. Min, "Tactics of adversarial attack on deep reinforcement learning agents," in *Proc. of IJCAI*, 2017.
- [4] V. Behzadan and A. Munir, "Vulnerability of deep reinforcement learning to policy induction attacks," in *Proc. of MLDM*, 2017, pp. 262-275.
- [5] A. Russo and A. Proutiere, "Optimal attacks on reinforcement learning policies," arXiv:1907.13548.
- [6] T. Osogami, "Robust partially observable Markov decision process," in *Proc. of ICML*, 2015.
- [7] L. Pinto, J. Davidson, R. Sukthankar, and A. Gupta, "Robust adversarial reinforcement learning," in *Proc. of ICML*, 2017.
- [8] A. Gleave, M. Dennis, C. Wild, N. Kant, S. Levine, and S. Russell, "Adversarial policies: Attacking deep reinforcement learning," arXiv:1905.10615.
- [9] K. Horak, Q. Zhu, and B. Bosansky, "Manipulating adversary's belief: A dynamic game approach to deception by design for proactive network security," in *Proc. of GameSec*, 2017.
- [10] S. Saritas, S. Yuksel, and S. Gezici, "Nash and Stackelberg equilibria for dynamic cheap talk and signaling games," in *Proc. of ACC*, 2017.
- [11] S. Saritas, E. Shereen, H. Sandberg, and G. Dán, "Adversarial attacks on continuous authentication security: A dynamic game approach," in *Proc. of GameSec*, 2019.
- [12] Z. Li and G. Dán, "Dynamic cheap talk for robust adversarial learning," in *Proc. of GameSec*, 2019.
- [13] M.O. Sayin and T. Basar, "Secure sensor design for cyber-physical systems against advanced persistent threats," in *Proc. of GameSec*, 2017.
- [14] M.O. Sayin, E. Akyol, and T. Basar, "Hierarchical multistage Gaussian signaling games in noncooperative communication and control systems," *Automatica*, vol. 107, pp. 9-20, 2019.
- [15] R. Zhang and P. Venkatasubramanian, "Stealthy control signal attacks in linear quadratic Gaussian control systems: Detectability reward tradeoff," *IEEE Transactions on Information Forensics and Security*, vol. 12, no. 7, pp. 1555-1570, 2017.
- [16] Q. Zhang, K. Liu, Y. Xia, and A. Ma, "Optimal stealthy deception attack against cyber-physical systems," *IEEE Transactions on Cybernetics*, 2020.
- [17] Y. Chen, S. Kar, and J.M.F. Moura, "Cyber physical attacks constrained by control objectives," in *Proc. of ACC*, 2016, pp. 1185-1190.
- [18] V.P. Crawford and J. Sobel, "Strategic information transmission," *Econometrica*, vol. 50, no. 6, pp. 1431-1451, 1982.
- [19] L. Shapley, "Stochastic games," *Proc. of the National Academy of Sciences*, vol. 39, no. 10, pp. 1095-1100, 1953.
- [20] T. Soderstrom, *Discrete-Time Stochastic Systems*, Springer, 2002.
- [21] A. Baranga, "The contraction principle as a particular case of Kleene's fixed point theorem," *Discrete Mathematics*, vol. 98, pp. 75-79, 1991.