

# Partially Observable Games for Secure Autonomy\*

Mohamadreza Ahmadi<sup>1</sup>, Arun A. Viswanathan<sup>2</sup>, Michel D. Ingham<sup>2</sup>, Kymie Tan<sup>2</sup>, and Aaron D. Ames<sup>1</sup>

**Abstract**—Technology development efforts in autonomy and cyber-defense have been evolving independently of each other, over the past decade. In this paper, we report our ongoing effort to integrate these two presently distinct areas into a single framework. To this end, we propose the two-player partially observable stochastic game formalism to capture both high-level autonomous mission planning under uncertainty and adversarial decision making subject to imperfect information. We show that synthesizing sub-optimal strategies for such games is possible under finite-memory assumptions for both the autonomous decision maker and the cyber-adversary. We then describe an experimental testbed to evaluate the efficacy of the proposed framework.

## I. INTRODUCTION

The growing ubiquity of autonomous systems, their use in ever more remote and unknown environments, and the increasing sophistication of cyber threats are driving a need for unprecedented system resilience, coupling robust autonomy with efficient cyber-defense strategies [10], [7]. Consider the push to develop swarms of smallsats in low Earth orbit. Cost-effective operations of such swarms require improved autonomy capabilities, both onboard and on the ground. However, complex autonomous behavior makes such systems susceptible to malicious tampering. Similarly, current unmanned air/ground/underwater systems rely on various signals for communication and localization and are already vulnerable to spoofing attacks. A GPS spoofing attack against such systems could result in malicious GPS coordinates being fed to the vehicle, causing it to be (mis)guided on an adversary's behest [8]. A resilient autonomous system should be able to detect attacks against itself, diagnose the probable causes, and automatically take corrective actions while ensuring the system's low/high-level goals and objectives are achieved.

However, a primary challenge to achieving this vision of integrated cyber and physical resilience is that technology development efforts in autonomy and cyber-defense are presently evolving independently of each other. Our work aims to reverse this trend. Our overall goal is to develop and demonstrate resilient autonomy for autonomous agents, by extending existing risk-aware planning and execution capabilities [14] with a combination of state-of-the-art model-

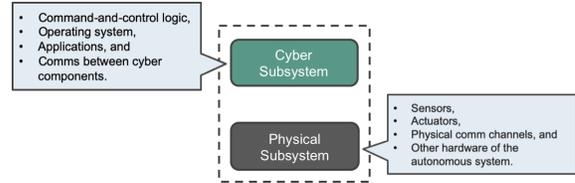


Fig. 1. A simplified model of an autonomous system.

based reasoning for situational and self-awareness and active cyber-defense mechanisms.

Current cyber adversaries can study the defender's behavior, identify security caveats, and modify their actions adaptively [15]. To tackle these security challenges, cyber-agents require adversarial decision making under uncertainty. Furthermore, agents cannot directly observe their adversary's true state and/or intention. Hence, active cyber-defense methods necessitate dealing with partial observations [2] and imperfect/incomplete information. A game-theoretic framework known as partially observable stochastic games (POSG) [11] provides a promising mathematical formalism for these capabilities.

In this paper, we report our preliminary methodology based on POSGs to integrate high-level autonomy and adversarial decision making. Our method based on POSGs is aimed at addressing cyber-physical threats caused by active cyber-adversaries, for example, as seen in the Stuxnet attack [12], wherein the attacker modifies their strategy in reaction to defensive actions. We show that the solution to the POSG can be cast as an optimization problem. Then, we propose an experimental setup to evaluate our technique. In summary, we hope to make the following contributions:

- Novel high-level resilient autonomy in the presence of active cyber-attacks leveraging the POSG framework;
- Demonstration of an integrated "defense-in-depth" capability for secure autonomy of cyber-physical systems.

The rest of the paper is organized as follows: Section II discusses the threat model for a cyber-physical system such as an UAV, an autonomous robot, or a swarms of spacecrafts; Section III discusses our proposed methodology using POSGs; Section IV discusses our experimental evaluation methodology followed by our conclusions and future work in Section V.

## II. CYBER-PHYSICAL THREAT MODEL

In this section, we first describe a model of an autonomous system, followed by a description of adversarial goals and a high-level taxonomy of threats.

\*The work described in this paper was performed at the California Institute of Technology, and at the Jet Propulsion Laboratory, California Institute of Technology, under a contract with the National Aeronautics and Space Administration (NASA).

<sup>1</sup>M. Ahmadi and A. D. Ames are with the California Institute of Technology, 1200 E. California Blvd., Pasadena, CA 91125. [mrahmadi, ames]@caltech.edu.

<sup>2</sup>A. Viswanathan, M. Ingham, and K. Tan are with the Jet Propulsion Laboratory, California Institute of Technology, 4800 Oak Grove Dr, Pasadena, CA 91109. [arun.a.viswanathan, michel.d.ingham, kymie.tan]@jpl.nasa.gov.

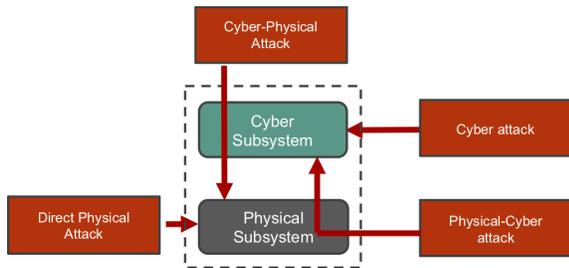


Fig. 2. Cyber-Physical Threat Model

Figure 1 shows a simplified model of an autonomous system (agent), containing two subsystems: *cyber* and *physical*. The cyber subsystem encapsulates functionality such as command and control logic, operating system, applications and any communications between the cyber components. Cyber components may be located on the agent or be external to the agent. Multi-agent systems may have a centralized cyber subsystem coordinating the agents. The physical subsystem encapsulates entities such as sensors, actuators, physical communication channels, and any other hardware comprising the autonomous system. An attacker would want to gain malicious control, cause damage, or deny service to prevent the autonomous system from achieving its goals. Referring to Figure 2, there are four different kinds of attacks an adversary could use to achieve their goals.

#### Cyber Attack

A cyber attack directly targets the components in the cyber subsystem. For example, a denial-of-service attack against the communication network of an autonomous system is an example of a cyber attack.

#### Physical Attack

A physical attack targets the components in the physical subsystem. For example, a ballistic impact is a type of physical attack which could damage physical components of an autonomous system. A physical attack often requires physical proximity to the system.

#### Cyber-Physical Attack

In a cyber-physical attack, an attacker leverages a cyber vulnerability with the intent to affect the physical subsystem. For example, malicious input injection attacks such as the malicious command or data injection seen in recent car hacks [9]. Cyber-physical attacks are often the most devastating as they can be initiated remotely, and cause serious damage to the physical subsystem.

#### Physical-Cyber Attack

In a physical-cyber attack, an attacker influences the cyber subsystem by attacking the components in the physical subsystem. For example, an attack on the physical sensors of an autonomous system (say the IMU), may cause inaccurate data to be sent upstream to the cyber components (for example,

incorrect location information), thereby causing incorrect decision-making and response by the cyber component.

In our work, we focus on the cyber-physical and physical-cyber kinds of attacks, as these attacks cross boundaries and as such, are often more subtle and difficult to diagnose, and consequently pose significant risk to missions. In addition, existing cyber or physical defenses generally do not protect against these attacks.

In the next section, we describe a mathematical formalism considering cyber-physical and physical-cyber attacks.

### III. METHODOLOGY: TWO-PLAYER POSG

A POSG is formally defined as follows.

A *probability distribution* over a finite or countably infinite set  $X$  is a function  $\mu: X \rightarrow [0, 1] \subseteq \mathbb{R}$  with  $\sum_{x \in X} \mu(x) = \mu(X) = 1$ . The set of all distributions on  $X$  is  $\text{Distr}(X)$ . The support of a distribution  $\mu$  is  $\text{supp}(\mu) = \{x \in X \mid \mu(x) > 0\}$ . A distribution is *Dirac* if  $|\text{supp}(\mu)| = 1$ .

**Definition 1:** A stochastic game (SG) is a tuple  $G = (S, s_I, \text{Act}, \mathcal{P})$  with a finite set  $S = S_\circ \cup S_\square$  of states, a set  $S_\circ$  of Player 1 states, a set  $S_\square$  of Player 2 states, the initial state  $s_I \in S$ , a finite set  $\text{Act} = \text{Act}_\circ \cup \text{Act}_\square$  of actions, and a transition function  $\mathcal{P}: S \times \text{Act} \rightarrow \text{Distr}(S)$ . We define costs using a state-action cost function  $C: S \times \text{Act} \rightarrow \mathbb{R}_{\geq 0}$ .

A Markov decision process (MDP) is an SG in which  $S_\circ = \emptyset$ , and consequently  $S = S_\square$ . A *path* of an SG  $G$  is an (in)finite sequence  $\pi = s_0 \xrightarrow{a_0} s_1 \xrightarrow{a_1} s$ , where  $s_0 = s_I$ ,  $s_i \in S$ ,  $a_i \in \text{Act}$ , and  $\mathcal{P}(s_i, a_i) \neq 0$  for all  $i \in \mathbb{N}$ . For finite  $\pi$ ,  $\text{last}(\pi)$  denotes the last state of  $\pi$ . The set of (in)finite paths of  $G$  is  $\text{Paths}_{\text{fin}}^G$  ( $\text{Paths}^G$ ).

To define a probability measure over the paths of an SG  $G$ , the non-determinism needs to be resolved by *strategies*.

**Definition 2 (SG strategy):** A strategy  $\sigma$  for  $G$  is a pair  $\sigma = (\sigma_\circ, \sigma_\square)$  of functions  $\sigma_i: \{\pi \in \text{Paths}_{\text{fin}}^G \mid \text{last}(\pi) \in S_i\} \rightarrow \text{Distr}(\text{Act})$  such that for all  $\pi \in \text{Paths}_{\text{fin}}^G$ ,  $\{a \mid \sigma_i(\pi)(a) > 0\} \subseteq \text{Act}$ ,  $i \in \{\circ, \square\}$ .

A Player- $i$  strategy  $\sigma_i$  (for  $i \in \{\circ, \square\}$ ) is *memoryless* if  $\text{last}(\pi) = \text{last}(\pi')$  implies  $\sigma_i(\pi) = \sigma_i(\pi')$  for all  $\pi, \pi' \in \text{dom}(\sigma_i)$ . It is *deterministic* if  $\sigma_i(\pi)$  is a Dirac distribution for all  $\pi \in \text{dom}(\sigma_i)$ .

A strategy  $\sigma$  for an SG resolves all non-deterministic choices, yielding an *induced MC*, for which a *probability measure* over the set of infinite paths is defined by the standard cylinder set construction [5]. These notions are analogous for MDPs.

In our framework,  $S_\circ$  consists of the physical and mission states, e.g. robot(s) location and obstacles, or the autonomous decision maker; whereas,  $S_\square$  corresponds to the internal states of the cyber-adversary. These states are not directly observable to either player; the players must infer the probability of their opponent being at different states based on the observations received at every step of the game. Thus, we have a POSG as follows (see Figure 3).

**Definition 3:** A partially observable stochastic game (POSG) is a tuple  $\mathcal{G} = (G, Z_\circ, Z_\square, O_\circ, O_\square)$ , with  $G = (S, s_I, \text{Act}, \mathcal{P})$  the underlying SG of  $\mathcal{G}$ ,  $Z_\circ$  and  $Z_\square$  are

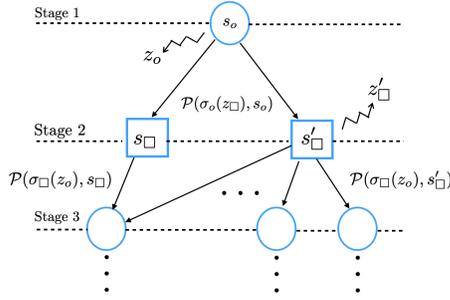


Fig. 3. Three stages of an example POSG. The states of the players need to be estimated based on the observations, and in the case of the attacker  $\square$ , counteracted. The game starts at  $s_0$  with an initial observation  $z_0$ .

finite set of observations for Player 1 and 2, respectively, and  $O_\circ: S \rightarrow Z_\circ$  ( $O_\square: S \rightarrow Z_\square$ ) the observation function for Player 1 (Player 2).

We lift the observation function to paths: For  $\pi = s_0 \xrightarrow{a_0} s_1 \xrightarrow{a_1} s_n \in \text{Paths}_{fin}^M$ , the associated *observation sequence* is  $O(\pi) = O(s_0) \xrightarrow{a_0} O(s_1) \xrightarrow{a_1} O(s_n)$ .

**Definition 4 (POSG Strategy):** An observation-based strategy  $\sigma_i$  for Player  $i$  in POSG  $\mathcal{G}$  is a strategy  $\sigma_i$  for Player  $i$  in the underlying SG  $G$  such that  $\sigma_i(\pi) = \sigma_i(\pi')$  for all  $\pi, \pi' \in \text{Paths}_{fin}^G$  with  $O_i(\pi) = O_i(\pi')$ .

Applying the strategy  $\sigma = (\sigma_\circ, \sigma_\square)$  to a POSG  $\mathcal{G}$  resolves all nondeterminism and partial observability, resulting in the induced Markov chain  $\mathcal{G}^\sigma$ .

However, since POSGs simply extend POMDPs to multiple players, computing optimal strategies requires infinite memory [6]. To circumvent this difficulty, we represent observation-based strategies with *finite* memory and we use *finite-state strategies* (FSSs) (see also FSSs in Delay Games [16]). If such an FSS has  $n$  memory states, we say the memory size for the underlying strategy  $\sigma$  is  $n$ .

**Definition 5 (FSS):** A finite-state strategy (FSS) for Player  $i$  in POSG  $\mathcal{G}$  is a tuple  $\mathcal{A}_i = (N_i, n_i^I, \gamma_i, \delta_i)$ , where  $N_i$  is a finite set of memory states,  $n_i^I \in N_i$  is the initial memory state,  $\gamma_i$  is the action mapping  $\gamma_i: N_i \times Z_i \rightarrow \text{Distr}(\text{Act})$ , and  $\delta_i$  is the memory update  $\delta_i: N_i \times Z_i \times \text{Act} \rightarrow \text{Distr}(N_i)$ . The set  $\text{FSS}_k^G$  denotes the set of FSSs with  $k$  memory states, called  $k$ -FSSs.

At each stage of the game, for each player, from a node  $n$  and the observation  $z$  in the current state of the POSG, the next action  $a$  is chosen from  $\text{Act}(z)$  randomly as given by  $\gamma(n, z)$ . Then, the successor node of the FSS is determined randomly via  $\delta(n, z, a)$ .

**A POSG for Secure Autonomy:** With the FSS assumption, the goal is then to maximize the probability of satisfying mission specifications, e.g. reach goal region while avoiding obstacles in the presence of cyber-adversarial activity. Next, we formally define the game objective.

**Game Objective:** For a POSG  $\mathcal{G}$  and a mission specification defined by a temporal logic formula  $\varphi$ , we consider the probability  $\Pr^{\mathcal{G}}(\varphi)$  to satisfy  $\varphi$ .

The specification  $\varphi$  is satisfied for a strategy  $\sigma = (\sigma_\circ, \sigma_\square)$  and the POSG  $\mathcal{G}$  with probability  $\lambda \in [0, 1]$ , if the probability  $\Pr^{\mathcal{G}^\sigma}(\varphi) = \lambda$  or simply if the induced Markov chain by applying strategy  $\sigma$  satisfies the specification with probabil-

ity  $\lambda$ . At this point, we have the following game formulation of secure autonomy problem.

**Problem 1:** Given a POSG  $\mathcal{G} = (G, Z_\circ, Z_\square, O_\circ, O_\square)$ , mission specification defined by a temporal logic formula  $\varphi$ , memory bounds  $n_\circ$  for the decision maker and  $n_\square$  for the cyber-adversary, compute a FSS  $\sigma_\circ^*$  such that

$$\sigma_\circ^* = \operatorname{argmax}_{\sigma_\circ \in \text{FSS}_{n_\circ}^G} \min_{\sigma_\square \in \text{FSS}_{n_\square}^G} \Pr^{\mathcal{G}^\sigma}(\varphi).$$

In Problem 1, we look for worst-case resilient strategies such that the probability of satisfying the specifications is maximized. Alternatively, we can search for resilient strategies that maximize the expected value of meeting the specifications in the presence of adversarial activity. Indeed, we can approximate  $\Pr^{\mathcal{G}^\sigma}(\varphi)$  with an expected total cost type constraint [3]. Then, for reachability type formulae such as  $\varphi = \diamond T$  (eventually reach a goal region represented by the states in  $T$ ), where  $T \subset S$ . The solution to Problem 1 can be found by solving an optimization problem as follows (see [1] for the derivation for one-sided POSGs).

For  $s_\circ \in S_\circ$ , and  $s_\square \in S_\square$ , we define the *cost variables*  $c_{s_\circ} \geq 0$   $c_{s_\square} \geq 0$  that represent the expected cost of reaching  $T \subseteq S$  with  $c_{s_I}$  being the expected cost of reaching to  $T$  from the initial state  $s_I$ . Let  $\gamma \in [0, 1]$  be the discount factor to ensure finite total expected cost. We then have the optimization problem:

$$\begin{aligned} & \text{minimize} && \text{maximize} && c_{s_I} \\ & c_{s_\circ}, \sigma_\circ && c_{s_\square}, \sigma_\square \end{aligned} \quad (1)$$

subject to

$$c_s = 0, \quad \forall s \in T, \quad (2)$$

$$\sum_{a \in \text{Act}_\circ} \sigma_a^z = 1, \quad \forall z \in Z_\square, \quad (3)$$

$$\sum_{a \in \text{Act}_\square} \sigma_a^z = 1, \quad \forall z \in Z_\circ, \quad (4)$$

$$c_{s_\square} = C(s_\square, a) + \gamma \sum_{a \in \text{Act}_\circ} \sigma_a^{O(s_\square)} \sum_{s'_\square \in S_\square} \mathcal{P}(s_\square, a, s'_\square) c_{s'_\square}, \quad \forall s_\square \in S_\square \setminus T, \forall \sigma_a^{O(s_\square)} \in \sigma_\square, \quad (5)$$

$$c_{s_\circ} = C(s_\circ, a) + \gamma \sum_{a \in \text{Act}_\square} \sigma_a^{O(s_\circ)} \sum_{s'_\square \in S_\square} \mathcal{P}(s_\circ, a, s'_\square) c_{s'_\square}, \quad \forall s_\circ \in S_\circ \setminus T, \forall \sigma_a^{O(s_\circ)} \in \sigma_\circ. \quad (6)$$

The objective in (1) implies the decision maker  $\circ$  is minimizing the cost of reaching  $T$  from the initial state; whereas, the cyber-adversary  $\square$  is trying to maximize the cost. We assign the expected cost of the states in the target set  $T$  to 0 by the constraints in (2). We ensure that the strategies of the decision maker and the cyber-adversary are well-defined with the constraints in (3) and (4). The constraints in (5)–(6) give the computation for the expected cost in the states of the POSG via dynamic programming.

We will develop methods based on heuristics and nonlinear programming to solve the resultant POSGs algorithmically and we will study trade-offs between resilience (cyber side)



Fig. 4. Three robots involved in experimental evaluations at CAST: (left) quadruped, (center) Segway, and (right) Flipper.

and mission goals (physical side). Preliminary work in solving POSGs was carried out in [1] for the case when only the adversary is subject to partial observation with application to network security. Instead of solving the full game, we used model checking to synthesize a set of strong (sub-optimal) strategies for the adversary and then composed robust defensive strategies.

#### IV. EXPERIMENTAL EVALUATION

The efficacy of the developed methods will be evaluated through experiments with three autonomous agents (a Segway, a quadruped, and a Flipper robot) in Caltech's Center for Autonomous Systems and Technologies (CAST) as depicted in Figure 4. The quadruped and the Flipper robot will be tasked to locate the target and the obstacles, respectively; whereas, the Segway is able to retrieve the target once the quadruped and Flipper explore the area. Flipper is equipped with a 3D LIDAR and a router. The quadruped robot is equipped with a high-resolution camera, an Inertial Measurement Unit (IMU), and a router. The Segway only has wheel odometry, an IMU, and a router. The centralized decision making is carried out through a computer connected to the robots via a wifi network. The sensor signals of each robot are also sent back to the computer via the same network.

Our previous experiments in this setting were concerned with safe autonomy enforced by discrete-time barrier functions [4], i.e., in the absence of cyber-adversaries (watch the experimental demonstration at [13]).

The goal of our next set of experiments is to find and retrieve the target in the presence of cyber adversarial activity. This experimental setup is described next.

The states of the POSG for Player  $\circ$  (the decision-maker) correspond to the locations of each agent, obstacles, and the goal. The actions for Player  $\circ$  include moving *Left*, *Right*, *Up*, *Down* for each agent. The two states of Player  $\square$  (cyber-adversary) are *Quadruped*, *Flipper* corresponding to the two surveying agents. The actions of the attacker are to *TakeDown* or *Wait*. If *TakeDown* is chosen at one stage of the game, for example, for the Flipper robot, the robot will not move in the next step and its observation cannot be used for path planning. On the other hand, *Wait* means no action is taken by the adversary.

The objective of Player  $\circ$  is then to maximize the probability of retrieving the target and avoiding obstacles; whereas,

the Player  $\square$  attempts to minimize this probability. This POSG fits in the framework of Section II and can be used to assure high-level mission autonomy as well as cyber-resilience. This initial abstract problem formulation will provide a basis for more realistic (high-fidelity) solutions to the real-world problem in future work, e.g., examining real injected cyber-attacks and practical defensive responses.

#### V. CONCLUSIONS

We described our ongoing research on the fusion of autonomous decision making and active cyber-resilience. We proposed a POSG that can capture high-level mission specifications, uncertainty, partial observation, and adversarial decision making. Although finding optimal strategies for POSGs is undecidable, we discussed finite-memory strategies as computationally tractable alternatives. Finally, we presented an experimental testbed, methodology and a case study to evaluate our secure autonomy techniques in the future.

#### ACKNOWLEDGMENT

The authors thank Prof. Richard M. Murray at Caltech and Dr. Nils Jansen at the Radboud University Nijmegen.

#### REFERENCES

- [1] M. Ahmadi, M. Cubuktepe, N. Jansen, S. Junges, J.-P. Katoen, and U. Topcu. The partially observable games we play for cyber deception. In *2019 American Control Conference*, 2019.
- [2] M. Ahmadi, M. Ono, Ingham, R. M. Murray, and Aaron D Ames. Risk-averse planning under uncertainty. In *2020 American Control Conference*, 2020.
- [3] M. Ahmadi, R. Sharan, and J. W. Burdick. Stochastic Finite State Control of POMDPs with LTL Specifications. *arXiv:2001.07679*, Jan 2020.
- [4] M. Ahmadi, A. Singletary, J. W. Burdick, and A. D. Ames. Safe Policy Synthesis in Multi-Agent POMDPs via Discrete-Time Barrier Functions. *58th Conference on Decision and Control*, Dec 2019.
- [5] Christel Baier and Joost-Pieter Katoen. *Principles of Model Checking*. MIT Press, 2008.
- [6] Krishnendu Chatterjee, Martin Chmelík, and Mathieu Tracol. What is decidable about partially observable markov decision processes with  $\omega$ -regular objectives. *Journal of Computer and System Sciences*, 82(5):878–911, 2016.
- [7] Gregory Falco. The vacuum of space cyber security. In *2018 AIAA SPACE Forum and Exposition*, page 5275, 2018.
- [8] S. M. Giray. Anatomy of unmanned aerial vehicle hijacking with signal spoofing. In *2013 6th International Conference on Recent Advances in Space Technologies (RAST)*, pages 795–800. IEEE, 2013.
- [9] Andy Greenberg. Hackers remotely kill a jeep on the highway – with me in it. <https://www.wired.com/2015/07/hackers-remotely-kill-jeep-highway/>, July 2015.
- [10] A. Kott. Intelligent autonomous agents are key to cyber defense of the future army networks. *The Cyber Defense Review*, 3(3):57–70, 2018.
- [11] Akshat Kumar and Shlomo Zilberstein. Dynamic programming approximations for partially observable stochastic games. In *Twenty-Second International FLAIRS Conference*, 2009.
- [12] Ralph Langner. Stuxnet: Dissecting a cyberwarfare weapon. *IEEE Security & Privacy*, 9(3):49–51, 2011.
- [13] J. Burdick M. Ahmadi, A. Singletary and A. Ames. Demonstration of safety-shield synthesis for multi-agent autonomy. <https://youtu.be/gmLNN8yA-oI>.
- [14] C. McGhan, R. Murray, T. Vaquero, B. Williams, M. Ingham, M. Ono, T. Estlin, R. Lanka, O. Arslan, and M. Elaasar. The resilient spacecraft executive: An architecture for risk-aware operations in uncertain environments. In *2016 AIAA SPACE Forum and Exposition*, 2016.
- [15] Colin Tankard. Advanced persistent threats and how to monitor and deter them. *Network security*, 2011(8):16–19, 2011.
- [16] Sarah Winter and Martin Zimmermann. Finite-state strategies in delay games. *Information and Computation*, page 104500, 2019.