

# A risk-security tradeoff in graphical coordination games

Keith Paarporn, Mahnoosh Alizadeh, and Jason R. Marden

**Abstract**—A system relying on the collective behavior of decision-makers can be vulnerable to a variety of adversarial attacks. How well can a system operator protect performance in the face of these risks? We frame this question in the context of graphical coordination games, where the agents in a network choose among two conventions and derive benefits from coordinating neighbors, and system performance is measured in terms of the agents’ welfare. In this paper, we assess an operator’s ability to mitigate two types of adversarial attacks - 1) broad attacks, where the adversary incentivizes all agents in the network and 2) focused attacks, where the adversary can force a selected subset of the agents to commit to a prescribed convention. As a mitigation strategy, the system operator can implement a class of distributed algorithms that govern the agents’ decision-making process. Our main contribution characterizes the operator’s fundamental trade-off between security against worst-case broad attacks and vulnerability from focused attacks. We show that this tradeoff significantly improves when the operator selects a decision-making process at random. Our work highlights the design challenges a system operator faces in maintaining resilience of networked distributed systems.

## I. INTRODUCTION

Networked distributed systems typically operate without centralized planning or control, and instead rely on local interactions and communication between the comprising agents. These systems arise in a variety of engineering applications such as teams of mobile robots and sensor networks [1]–[3]. They are also prevalent in social dynamics [4], [5] and biological populations [6].

The transition from a centralized to a distributed architecture may leave a system vulnerable to a variety of adversarial attacks. An adversary may be able to manipulate the decision-making processes of the agents. Such dynamical perturbations can potentially lead to unwanted outcomes. For example in social networks, individual opinions can be shaped from external information sources, resulting in a polarized society [7], [8]. When feasible, a system operator takes measures to mitigate adversarial influences. The literature on cyber-physical system security studies many aspects of this interplay. For instance, optimal controllers are designed to mitigate denial-of-service, estimation, and deception attacks [9]–[13].

This paper investigates security measures that a system operator can take against adversarial influences when the

underlying system is a graphical coordination game [5], [14], where agents in a network decide between two choices,  $x$  or  $y$ . One may think of these choices as two competing products, e.g. iPhone vs Android, two conflicting social norms, or two opposing political parties. Each agent derives a positive benefit from interactions with coordinating neighbors, and zero benefits from mis-coordinating ones. The system’s efficiency is defined by the ratio of total benefits of all agents to the maximal attainable benefits over all configurations of choices.

The goal of the system operator is to design a local decision-making rule for each agent in the system so that the emergent collective behavior optimizes system efficiency. One algorithm that achieves this goal is known as log-linear learning [15]–[17]. More formally, the agents follow a “perturbed” best reply dynamics where the agents’ local objectives are precisely equal to their local welfare. We seek to address the question of whether this particular algorithm is robust to adversarial influences. That is, does this algorithm preserve system efficiency when the agents’ decision-making processes are manipulated by an adversary? If not, can the operator alter the agents’ local objectives to mitigate such attacks?

We consider two adversarial attack models - *broad* and *focused* attacks. In broad attacks, the adversary incentivizes every agent in the network (hence broad) with a convention, influencing their decision-making process. This could depict distributing political ads with the intention of polarizing voters. In focused attacks, the adversary targets a specific set of agents in the network, forcing them to commit to  $x$  or  $y$ . These targeted, or fixed agents consequently do not update their choices over time but still influence the decisions of others. For instance, they could portray loyal consumers of a brand or product, or staunch supporters of a political party. Fixed agents and their effects on system performance have been extensively studied in the context of opinion dynamics and optimization algorithms [13], [18], [19].

The first contribution of this paper is a characterization of worst-case risk metrics from both adversarial attacks as a function of the operator’s algorithm design parameter (Section III). We define risk in this paper as the system’s distance to optimal efficiency. By worst-case here we mean the maximum risk among all connected network topologies subject to any admissible adversarial attack. Hence, our analysis identifies the network topologies on which worst-case risks are attained (Section V). We extend this analysis to randomized operator strategies (Sections IV, VI).

The second contribution of this paper answers the question “if the operator succeeds in protecting the system from one type of attack, how vulnerable does it leave the system to the other?” We identify a fundamental tradeoff between security against broad attacks and risks from focused attacks. We then show randomized operator strategies significantly improves

K. Paarporn (kpaarporn@ucsb.edu), M. Alizadeh (alizadeh@ucsb.edu), and J. R. Marden (jrmarden@ece.ucsb.edu) are with the Department of Electrical and Computer Engineering at the University of California, Santa Barbara, CA.

A preliminary version of this paper has been submitted for conference publication (CDC 2019) and can be found at <https://ece.ucsb.edu/~jrmarden/files/MardenC9.pdf>. The analysis of randomized operator strategies extends the previous results. Complete proofs are also provided here. This work is supported by UCOP Grant LFR-18-548175, ONR grant #N00014-17-1-2060, and NSF grant #ECCS-1638214.

the set of attainable risk levels and their associated tradeoffs (Section IV).

By characterizing this interplay, we contribute to previous work that studied the impact of adversarial influence in graphical coordination games [20]–[22]. These works analyze worst-case damages that can be inflicted by varying degrees of adversarial sophistication and intelligence in the absence of a system operator. However, these results were derived only in specific graph structures, namely ring graphs, whereas our analysis considers adversarial influence in any graph topology.

## II. PRELIMINARIES

### A. Graphical coordination games

A graphical coordination game is played between a set of agents  $\mathcal{N} = \{1, \dots, N\}$  over a connected undirected network  $G = (\mathcal{N}, \mathcal{E})$  with node set  $\mathcal{N}$  and edge set  $\mathcal{E} \subset \mathcal{N} \times \mathcal{N}$ . Agent  $i$ 's set of neighbors is written as  $\mathcal{N}_i = \{j : (i, j) \in \mathcal{E}\}$ . Each agent  $i$  selects a choice  $a_i$  from its action set  $\mathcal{A}_i = \{x, y\}$ . The choices of all the agents constitutes an action profile  $a = (a_1, \dots, a_N)$ , and we denote the set of all action profiles as  $\mathcal{A} = \prod_{i=1}^N \mathcal{A}_i$ . The local interaction between two agents  $(i, j) \in \mathcal{E}$  is based on a  $2 \times 2$  matrix game, described by the payoff matrix  $V : \{x, y\}^2 \rightarrow \mathbb{R}$ ,

$$\begin{array}{cc|cc} & & \text{Player } j & & \\ & & x & y & \\ \text{Player } i & x & 1 + \alpha_{\text{sys}}, 1 + \alpha_{\text{sys}} & 0, 0 & \\ & y & 0, 0 & 1, 1 & \end{array} \quad (1)$$

where  $\alpha_{\text{sys}} > 0$  is the system *payoff gain*. It indicates that  $x$  is an inherently superior product over  $y$  when users coordinate. Here, agents would rather coordinate than not, but prefer to coordinate on  $x$ . Agent  $i$ 's *benefit* is the sum of payoffs derived from playing the game (1) with each of its network neighbors:

$$W_i(a_i, a_{-i}) := \sum_{j \in \mathcal{N}_i} V(a_i, a_j). \quad (2)$$

A measure of system *welfare* defined over  $\mathcal{A}$  is

$$W(a) := \sum_{i=1}^N W_i(a_i, a_{-i}), \quad (3)$$

which is simply the sum of all agent benefits. The *system efficiency* for action profile  $a \in \mathcal{A}$  is defined as

$$\frac{W(a)}{\max_{a' \in \mathcal{A}} W(a')}. \quad (4)$$

For  $\mathcal{A} = \{x, y\}^N$ , the all- $x$  profile  $\vec{x}$  maximizes welfare. This does not necessarily hold for arbitrary action spaces.

### B. Log-linear learning algorithm

Log-linear learning is a distributed stochastic algorithm governing how players' decisions evolve over time [14]–[16]. It may be applied to any instance of a game with each player having a well-defined local utility function  $U_i : \mathcal{A} \rightarrow \mathbb{R}$  over a set of action profiles  $\mathcal{A}$  with an underlying interaction graph  $G$ . That is, agent  $i$ 's local utility is a function of its action  $a_i$  and actions of its neighbors in  $G$ .

Agents update their decisions  $a(t) \in \mathcal{A}$  over discrete time steps  $t = 0, 1, \dots$ . Assume  $a(0)$  is arbitrarily determined. For step  $t \geq 1$ , one agent  $i$  is selected uniformly at random from the population. It updates its action to  $a_i(t) = z \in \mathcal{A}_i$  with probability

$$\frac{\exp(\beta U_i(z, a_{-i}(t-1)))}{\sum_{z' \in \mathcal{A}_i} \exp(\beta U_i(z', a_{-i}(t-1)))}, \quad (5)$$

where  $\beta > 0$  is the rationality parameter. All other agents repeat their previous actions:  $a_{-i}(t) = a_{-i}(t-1)$ . For large values of  $\beta$ ,  $i$  selects a best-response to the previous actions of others with high probability, and for values of  $\beta$  near zero,  $i$  randomizes among its actions  $\mathcal{A}_i$  uniformly at random. This induces an irreducible Markov chain over the action space  $\mathcal{A}$ , with a unique stationary distribution  $\pi_\beta \in \Delta(\mathcal{A})$ . The *stochastically stable states* (SSS)  $a \in \mathcal{A}$  are the action profiles contained in the support of the stationary distribution in the high rationality limit: they satisfy  $\pi(a) = \lim_{\beta \rightarrow \infty} \pi_\beta(a) > 0$ . Such a limiting distribution exists and is unique [15], [23], [24]. We write the set of stochastically stable states as

$$\text{LLL}(\mathcal{A}, \{U_i\}_{i \in \mathcal{N}}; G). \quad (6)$$

For graphical coordination games, the log-linear learning algorithm specified by the action set  $\mathcal{A} = \{x, y\}^N$  and utilities  $\{W_i\}_{i \in \mathcal{N}}$  selects the welfare-maximizing profile  $\vec{x}$  as the stochastically stable state irrespective of the graph topology  $G$ . This can be shown using standard potential game arguments [16] (we provide these details in Section V). That is,  $\vec{x} = \text{LLL}(\mathcal{A}, \{W_i\}_{i \in \mathcal{N}}; G)$  for all  $G \in \mathcal{G}_N$ , where  $\mathcal{G}_N$  is the set of all connected undirected graphs on  $N$  nodes.

However, if an adversary is able to manipulate the agents' local decision-making rules, this statement may no longer hold true. A system operator may be able to alter the agents' local utility functions with the goal of mitigating the loss of system efficiency in the presence of adversarial influences. In particular, we consider the class of local utility functions  $\{U_i^\alpha\}_{i \in \mathcal{N}}$  parameterized by  $\alpha > 0$ . Specifically,  $U_i^\alpha$  takes the same form as the benefit function (2) where  $\alpha_{\text{sys}}$  is replaced with a *perceived gain*  $\alpha$  that is under the operator's control. We next introduce models of adversarial attacks in graphical coordination games. We then evaluate the performance of this class of distributed algorithms in the face of adversarial attacks.

## III. MODELS OF ADVERSARIAL INFLUENCE

In this section, we outline two models of adversarial attacks in graphical coordination games - *broad* and *focused* attacks. The system operator specifies the local utility functions  $\{U_i^\alpha\}$  that govern the log-linear learning algorithm by selecting the perceived payoff gain  $\alpha > 0$ . Our goal is to assess the performance of this range of algorithms on two corresponding worst-case risk metrics, which we define and characterize. We then identify fundamental tradeoff relations between these two risk metrics.

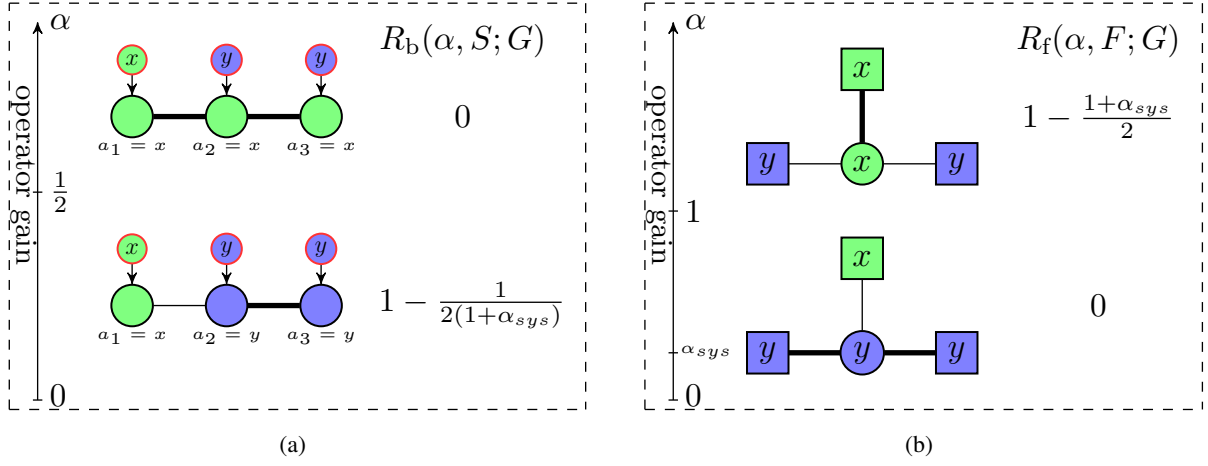


Fig. 1: (Left) An example three-node line network under a broad adversarial attack. The imposter nodes are depicted as the labelled smaller circles and agents in the network are the bigger circles. The color of each circle indicates the node's action - green for  $x$ , blue for  $y$ . In this example, maximum welfare is  $\max_{a \in \mathcal{A}} W(a) = 4(1 + \alpha_{\text{sys}})$ , achieved when all three agents play  $x$ . The adversary's target set  $S$  attaches an  $x$ -imposter to node 1 and  $y$ -impsters to nodes 2 and 3. For operator gains  $\alpha \leq \frac{1}{2}$ ,  $a = (a_1, a_2, a_3) = (x, y, y)$  is the welfare-minimizing SSS, i.e. it satisfies  $a = \arg \min_{a \in \text{LLL}(\mathcal{A}, \alpha, S; G)} W(a)$ . This gives a risk of  $R_b(\alpha, S; G) = 1 - \frac{1}{2(1 + \alpha_{\text{sys}})}$ . For  $\alpha > \frac{1}{2}$ , the welfare-minimizing SSS is  $(x, x, x)$ . This gives optimal efficiency, i.e. a risk of 0. (Right) An example of a four node star network under a focused attack where a subset  $F$  of three nodes are targeted to be fixed (squares). Only the center node is unfixed. In this example, the maximum welfare is  $\max_{a \in \mathcal{A}_F} W(a) = 4$ , achieved when the center plays  $y$ . This is because the alternative action (when center plays  $x$ ) gives the suboptimal welfare  $2(1 + \alpha_{\text{sys}}) < 4$  due to  $\alpha_{\text{sys}} < 1$ . For operator gains  $\alpha < 1$ , the center node plays  $y$  in the SSS. This yields optimal efficiency, i.e. the risk is  $R_f(\alpha, F; G) = 0$ . For  $\alpha \geq 1$ , the center node plays  $x$ , giving a risk of  $R_f(\alpha, F; G) = 1 - \frac{1 + \alpha_{\text{sys}}}{2}$ . The methods to calculate stochastically stable states under both types of attacks follow standard potential game arguments and are detailed in Section V.

#### A. Broad attacks and worst-case risk metric

We consider a scenario where the system is subject to broad adversarial attacks. For each agent in the network, the adversary attaches a single imposter node that acts as a neighbor that always plays  $x$  or  $y$ . These nodes are not members of the network but affect the decision making of agents that are. Let  $S_x \subseteq \mathcal{N}$  ( $S_y$ ) be the set of agents targeted with an imposter  $x$  ( $y$ ) node. We call the *target set*  $S = (S_x, S_y)$ . Any target set satisfies  $S_x \cap S_y = \emptyset$  and  $S_x \cup S_y = \mathcal{N}$ . We call  $\mathcal{T}(G)$  the set of all possible target sets  $S$  on the graph  $G$ . Given  $\alpha > 0$ , the agents' *perceived* utilities are

$$\tilde{U}_i^\alpha(a_i, a_{-i}) := \begin{cases} U_i^\alpha(a_i, a_{-i}) + \mathbb{1}(a_i = y) & i \in S_y \\ U_i^\alpha(a_i, a_{-i}) + (1 + \alpha)\mathbb{1}(a_i = x) & i \in S_x \end{cases} \quad (7)$$

In the notation of (6), the set of stochastically stable states is written  $\text{LLL}(\mathcal{A}, \{\tilde{U}_i^\alpha\}_{i \in \mathcal{N}}; G)$ . However for more specificity, we will refer to it in this context as  $\text{LLL}(\mathcal{A}, \alpha, S; G)$ . The induced network *efficiency* is defined as

$$J_b(\alpha, S; G) := \frac{\min_{a \in \text{LLL}(\mathcal{A}, \alpha, S; G)} W(a)}{\max_{a' \in \mathcal{A}} W(a')} = \frac{\min_{a \in \text{LLL}(\mathcal{A}, \alpha, S; G)} W(a)}{(1 + \alpha_{\text{sys}})|\mathcal{E}|}, \quad (8)$$

which is the ratio of the welfare induced by the welfare-minimizing SSS to the optimal welfare. The second equality above is due to the fact that optimal welfare is attained at  $\bar{x}$  (all play  $x$ ). We re-iterate that the imposter nodes serve only to modify the stochastically stable states, and do not contribute

to the system welfare  $W(a)$  (3). The *risk* from broad attacks faced by the system operator in choosing gain  $\alpha$  is defined as

$$R_b(\alpha, S; G) := 1 - J_b(\alpha, S; G). \quad (9)$$

Risk measures the distance from optimal efficiency under operating gain  $\alpha$ . Fig. 1a illustrates an example of a three-node network subject to a broad adversarial attack. The extent to which systems are susceptible to broad attacks is captured by the following definition of worst-case risk.

**Definition 1.** *The worst-case risk to broad attacks is given by*

$$R_b^*(\alpha) := \max_{N \geq 3} \max_{G \in \mathcal{G}_N} \max_{S \in \mathcal{T}(G)} R_b(\alpha, S; G), \quad (10)$$

The quantity  $R_b^*(\alpha)$  is the cost metric that the system operator wishes to reduce given uncertainty of the network structure and target set.

**Theorem 1.** *Let  $\alpha > 0$ . The worst-case broad risk is*

$$R_b^*(\alpha) = \begin{cases} 1 - \left(\frac{k}{k+1}\right) \left(\frac{1}{1 + \alpha_{\text{sys}}}\right) & \text{if } \alpha \in I_k, \text{ for } k = 1, 2, \dots \\ 1 - \frac{1}{1 + \alpha_{\text{sys}}} & \text{if } \alpha \in \left[1, \frac{3}{2}\right] \\ 0 & \text{if } \alpha > \frac{3}{2} \end{cases} \quad (11)$$

where

$$I_k := \left(\frac{k-1}{k}, \frac{k}{k+1}\right]. \quad (12)$$

It is a piecewise constant function on half-open intervals that is monotonically decreasing in  $\alpha$ . An illustration is given in Figure 2a, along with the graphs and target sets that achieve the worst-case risks. For sufficiently high gains  $\alpha > 3/2$ , the system is safeguarded from any broad adversarial attack,

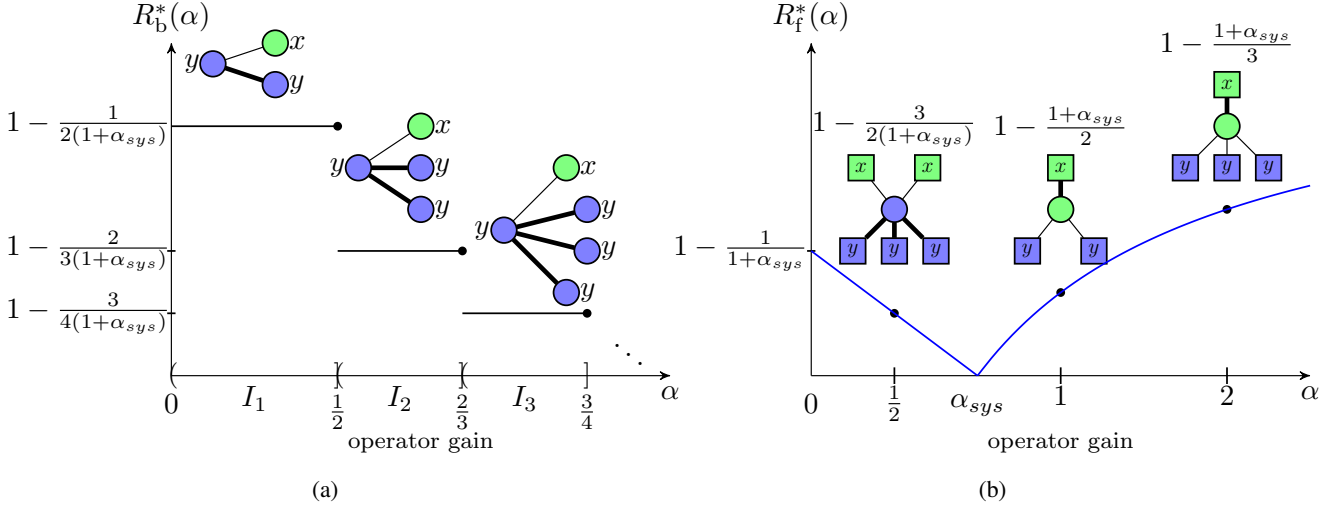


Fig. 2: (a) The worst-case risk from broad attacks  $R_b^*(\alpha)$  (11) is a piecewise constant function defined over countably infinite half-open intervals. The graphs and their corresponding target set which attain each level of worst-case broad risk are illustrated for  $\alpha < 1$ . Here, the  $x, y$  labels indicate the type of imposter influence on the agents (circles) in the network, and the color of the circles depict the action played in the welfare-minimizing SSS (green= $x$ , blue= $y$ ). If  $\alpha \in I_k$ ,  $k = 1, 2, \dots$  (recall (12)), the worst-case risk is achieved on a star graph of  $k + 2$  nodes where all nodes but one are targeted with a  $y$  imposter. The one leaf node has an  $x$  imposter attached, giving a single miscoordinating link in the network. (b) The worst-case risk from focused attacks  $R_f^*(\alpha)$  (16). The graphs and their corresponding fixed sets which attain the worst-case focused risks are illustrated for  $\alpha = \frac{1}{2}, 1$ , and  $2$ . The nodes' color represents the worst-case SSS at  $\alpha$  (blue =  $y$ , green =  $x$ ). The targeted fixed agents are represented as squares and the unfixed agents as circles. Here  $\frac{1}{2} < \alpha_{sys} < 1$ . The proofs establishing all worst-case graphs are detailed in Section V.

i.e. the worst-case risk is zero. By inflating the value of the  $x$ -convention, the adversary is unable to induce any miscoordinating links or agents to play  $y$ . The technical results needed for the proof are given in Section V.

### B. Focused attacks and worst-case risk metric

An adversary is able to choose a strict subset of agents and force them to commit to prescribed choices. This causes them to act as *fixed agents*, or agents that do not update their choices over time. One could consider this as allowing the adversary an unlimited number of imposter nodes (instead of one) at its dispatch to attach to each agent in the subset, thereby solidifying their choices. This focused influence on a single agent is stronger than the influence a broad attack has on a single agent in the sense that the latter type does not require the agent to commit to a choice - it merely incentivizes the agent towards one particular choice.

Let  $F_x \subset \mathcal{N}$  ( $F_y$ ) be the set of fixed  $x$  ( $y$ ) agents. We call the *fixed set*  $F = (F_x, F_y)$ , which satisfies  $F_x \cap F_y = \emptyset$  and  $F_x \cup F_y \subset \mathcal{N}$ . We call  $\mathcal{F}(G)$  the set of all feasible fixed sets on a graph  $G \in \mathcal{G}_N$ . A fixed set  $F \in \mathcal{F}(G)$  restricts the action space to  $\mathcal{A}(F)$ , where  $\mathcal{A}_i(F) = \{x\}$  ( $\{y\}$ )  $\forall i \in F_x$  ( $F_y$ ) and  $\mathcal{A}_i(F) = \{x, y\}$   $\forall i \notin F$ . We assume the adversary selects at least one fixed agent. The strict subset assumption avoids pathological cases (e.g. alternating  $x$  and  $y$  fixed nodes for an entire line network yields an efficiency of zero).

The set of stochastically stable states given a fixed set  $F$  is written as  $\text{LLL}(\mathcal{A}(F), \{U_i^\alpha\}_{i \in \mathcal{N}}; G)$ . However for brevity, we will refer to it as  $\text{LLL}(\mathcal{A}(F), \alpha; G)$ . The induced efficiency is

$$J_f(\alpha, F; G) := \frac{\min_{a \in \text{LLL}(\mathcal{A}(F), \alpha; G)} W(a)}{\max_{a \in \mathcal{A}(F)} W(a)}, \quad (13)$$

which is the ratio of the welfare induced by the worst-case stable state to the optimal welfare given the fixed set  $F$ . The risk faced by the system operator in choosing  $\alpha$  is defined as

$$R_f(\alpha, F; G) := 1 - J_f(\alpha, F; G). \quad (14)$$

Again, risk measures the distance from optimal efficiency when choosing  $\alpha$ . The fixed nodes here differ from the imposter nodes in that they contribute to the true measured welfare (3) in addition to modifying the SSS by restricting the action set and influencing the decisions of their non-fixed neighbors. Figure 1b provides an illustrative example of a network with three fixed agents and one unfixed agent. The extent to which the system is susceptible to focused attacks is defined by the following worst-case risk metric.

**Definition 2.** The worst-case risk from focused attacks is given by

$$R_f^*(\alpha) := \max_{N \geq 3} \max_{G \in \mathcal{G}_N} \max_{F \in \mathcal{F}(G)} R_f(\alpha, F; G). \quad (15)$$

The quantity  $R_f^*(\alpha)$  is the cost metric that a system operator wishes to reduce given uncertainty on the graph structure and composition of fixed agents in the network.

**Theorem 2.** The worst-case risk from focused attacks is

$$R_f^*(\alpha) = \begin{cases} 1 - \frac{1+\alpha}{1+\alpha_{sys}}, & \text{if } \alpha < \alpha_{sys} \\ 0, & \text{if } \alpha = \alpha_{sys} \\ 1 - \frac{1+\alpha_{sys}}{1+\alpha}, & \text{if } \alpha > \alpha_{sys} \end{cases}. \quad (16)$$

The technical results needed for the proof are given in Section V. An illustration of this quantity as well as the graphs that induce worst-case risk are portrayed in Figure 2b. We observe the choice  $\alpha = \alpha_{sys}$  recovers optimal efficiency for

any  $G \in \mathcal{G}_N$  and  $F \in \mathcal{F}(G)$ . In other words, by operating at the system gain  $\alpha_{\text{sys}}$ , the system operator safeguards efficiency from any focused attack. Furthermore,  $R_f^*(\alpha)$  monotonically increases for  $\alpha > \alpha_{\text{sys}}$ , approaching 1 in the limit  $\alpha \rightarrow \infty$ . Intuitively, the risk in this regime comes from inflating the benefit of the  $x$  convention, which can be harmful to system efficiency when there are predominantly fixed  $y$  nodes in the network. For  $\alpha < \alpha_{\text{sys}}$ ,  $R_f^*(\alpha)$  monotonically decreases. The risk here stems from de-valuing the  $x$  convention, which hurts efficiency when coordinating with fixed  $x$  nodes is more valuable than coordinating with fixed  $y$  nodes.

### C. Fundamental tradeoffs between risk and security

We describe the operator's tradeoffs between the two worst-case risk metrics. That is, given a level of security  $\gamma \in [0, 1]$  is ensured on one worst-case risk, what is the minimum achievable risk level of the other? These relations are direct consequences of Theorems 1 and 2.

**Remark 1.** *Before presenting the tradeoff relations, we first observe that since  $R_f^*(\alpha)$  is decreasing on  $\alpha < \alpha_{\text{sys}}$  and  $R_b^*(\alpha)$  is decreasing in  $\alpha$ , the operator should not select any gain  $\alpha < \alpha_{\text{sys}}$ , as it worsens both risk levels. Hence for the rest of this paper, we only consider gains greater than  $\alpha_{\text{sys}}$ .*

**Corollary 1.** *Fix  $\gamma_f \in [0, 1)$ . Suppose  $R_f^*(\alpha) \leq \gamma_f$  for some  $\alpha$ . Then*

$$R_b^*(\alpha) \geq R_b^* \left( \frac{1 + \alpha_{\text{sys}}}{1 - \gamma_f} - 1 \right). \quad (17)$$

*Proof.* From (16),  $R_f^*(\alpha) \leq \gamma_f$  implies  $\alpha \leq \frac{1 + \alpha_{\text{sys}}}{1 - \gamma_f} - 1$ . Since  $R_b^*(\alpha)$  is a decreasing function in  $\alpha$ , we obtain the result. ■

In words, as the security from worst-case focused attacks improves ( $\gamma_f$  lowered), the risk from worst-case broad attacks increases. A tradeoff relation also holds in the opposite direction.

**Corollary 2.** *Fix  $\gamma_b \in [0, 1]$ . Suppose  $R_b^*(\alpha) \leq \gamma_b$  for some  $\alpha$ . Suppose  $\alpha_{\text{sys}} \in I_{k_{\text{sys}}}$  for some  $k_{\text{sys}} \in \{1, 2, \dots\}$ . Then*

$$R_f^*(\alpha) \begin{cases} \geq 0 & \text{if } \gamma_b \in \left[ 1 - \frac{k_{\text{sys}}}{k_{\text{sys}}+1} \frac{1}{1+\alpha_{\text{sys}}}, 1 \right] \\ > R_f^* \left( \frac{k}{k+1} \right) & \text{if } \gamma_b \in \left[ 1 - \frac{k}{k+1} \frac{1}{1+\alpha_{\text{sys}}}, 1 - \frac{k-1}{k} \frac{1}{1+\alpha_{\text{sys}}} \right) \\ & \text{for } k = k_{\text{sys}}, k_{\text{sys}} + 1, \dots \\ \geq R_f^*(1) & \text{if } \gamma_b = 1 - \frac{1}{1+\alpha_{\text{sys}}} \\ > R_f^* \left( \frac{3}{2} \right) & \text{if } \gamma_b \in \left[ 0, 1 - \frac{1}{1+\alpha_{\text{sys}}} \right) \end{cases} \quad (18)$$

If  $\alpha_{\text{sys}} \in [1, 3/2]$ ,

$$R_f^*(\alpha) \begin{cases} \geq 0 & \text{if } \gamma_b \in \left[ 1 - \frac{1}{1+\alpha_{\text{sys}}}, 1 \right] \\ > R_f^* \left( \frac{3}{2} \right) & \text{if } \gamma_b \in \left[ 0, 1 - \frac{1}{1+\alpha_{\text{sys}}} \right) \end{cases} \quad (19)$$

If  $\alpha_{\text{sys}} > \frac{3}{2}$ , then  $R_f^*(\alpha) \geq 0$  for any  $\gamma_b$ .

*Proof.* All bounds are computed by finding  $\inf_{\alpha} R_f^*(\alpha)$  s.t.  $R_b^*(\alpha) \leq \gamma_b$ . The relations  $\geq$  and  $>$  follow from the fact that  $R_f^*(\alpha)$  is increasing in  $\alpha > \alpha_{\text{sys}}$ , and depending on whether  $R_f^*$  can attain the resulting value. ■

Here, as the security from worst-case broad attacks improves ( $\gamma_b$  lowered), the risk from worst-case focused attacks increases. Each of the broad risk levels can be attained for a range of focused risks. An illustration of the attainable worst-case risk levels is given in Fig. 3 (blue).

## IV. RANDOMIZED OPERATOR STRATEGIES

In this section, we consider the scenario where the operator randomizes over multiple gains. We present a definition and a characterization of worst-case expected risks. We then identify the risk-security tradeoffs available in the randomized gain setting. We observe they significantly improve upon the deterministic gain setting (Fig. 3). We then identify ways to further improve these tradeoffs through different randomizations.

### A. Worst-case expected risks

Suppose the operator selects a gain from the  $M$  distinct values  $\alpha = \{\alpha_k\}_{k=1}^M$  satisfying  $\alpha_1 < \alpha_2 < \dots < \alpha_M$  with the probability distribution  $\mathbf{p} = [p_1, \dots, p_M]^T \in \Delta_M$ . Here we denote  $\Delta_M = \{\mathbf{p} \in \mathbb{R}_+^M : \sum_{j=1}^M p_j = 1\}$  as the set of all  $M$ -dimensional probability vectors. In other words, the operator employs the payoff gain  $\alpha_j$  with probability  $p_j$ .

We consider the following natural definitions of expected risks. Given a graph  $G \in \mathcal{G}_N$  and target set  $S \in \mathcal{T}(G)$ , let  $\mathbb{E}_{\alpha, \mathbf{p}}[R_b|S, G] := \sum_{j=1}^M p_j R_b(\alpha_j, S; G)$  be the expected adversarial risk of the operator's strategy  $\alpha, \mathbf{p}$ . The worst-case expected risk from broad attacks is defined as

$$\mathbb{E}_{\alpha, \mathbf{p}}^*[R_b] := \max_{N \geq 3} \max_{G \in \mathcal{G}_N} \max_{S \in \mathcal{T}(G)} \mathbb{E}_{\alpha, \mathbf{p}}[R_b|S, G]. \quad (20)$$

Similarly, given a fixed set  $F \in \mathcal{F}(G)$ , let  $\mathbb{E}_{\alpha, \mathbf{p}}[R_f|F, G] := \sum_{j=1}^M p_j R_f(\alpha_j, F; G)$  be the expected risk from focused attacks. The worst-case expected risk from focused attacks is defined as

$$\mathbb{E}_{\alpha, \mathbf{p}}^*[R_f] := \max_{N \geq 3} \max_{G \in \mathcal{G}_N} \max_{F \in \mathcal{F}(G)} \mathbb{E}_{\alpha, \mathbf{p}}[R_f|F, G]. \quad (21)$$

**Theorem 3.** *Suppose the operator randomizes with gains  $\alpha = \{\alpha_k\}_{k=1}^M$  according to  $\mathbf{p} \in \Delta_M$ . Then the worst-case expected broad risk is*

$$\mathbb{E}_{\alpha, \mathbf{p}}^*[R_b] = \max_{k=1, \dots, M} \left\{ \left( \sum_{j=1}^k p_j \right) R_b^*(\alpha_k) \right\}. \quad (22)$$

*The worst-case expected focused risk is*

$$\mathbb{E}_{\alpha, \mathbf{p}}^*[R_f] = \max_{k=1, \dots, M} \left\{ \left( \sum_{j=k}^M p_j \right) R_f^*(\alpha_k) \right\}. \quad (23)$$

The proofs are given in Section VI. The characterization of worst-case expected risk is a discounted weighting of a deterministic worst-case risk level. This suggests that the risk levels achievable by randomization can improve upon the risks induced from a deterministic gain.

### B. Risk tradeoffs under randomized operator strategies

Given a level of security  $\gamma \in [0, 1]$  is ensured on one expected worst-case metric, what is the the minimum achievable risk level on the other? We find this can be calculated through a linear program. We formalize these tradeoffs in the following two statements, which are analogous to Corollaries 1 and 2.

**Corollary 3.** Fix  $\gamma_f \in [R_f^*(\alpha_1), 1]$  and a set of gains  $\alpha = \{\alpha_j\}_{j=1}^M$ . Suppose  $\mathbb{E}_{\alpha, \mathbf{p}}^*[R_f] \leq \gamma_f$  for some  $\mathbf{p} \in \Delta_M$ . Then

$$\mathbb{E}_{\alpha, \mathbf{p}}^*[R_b] \geq v_b(\gamma_f, \alpha), \quad (24)$$

where  $v_b(\gamma, \alpha)$  is the value of the following linear program.

$$\begin{aligned} v_b(\gamma_f, \alpha) = \min_{\mathbf{p}', v} v \\ \text{s.t. } \sum_{i=1}^M p'_i = 1, p_i \geq 0 \forall i = 1, \dots, M \\ v \in [0, 1] \\ A_{\text{LP}} \begin{bmatrix} \mathbf{p}' \\ v \end{bmatrix} \preceq \begin{bmatrix} 0_M \\ \gamma_f \mathbb{1}_M \end{bmatrix} \end{aligned} \quad (25)$$

where  $\preceq$  denotes elementwise  $\leq$ ,  $0_M$  and  $\mathbb{1}_M$  are column  $M$ -vectors of zeros and ones respectively, and  $A_{\text{LP}}$  is the  $2M \times (M+1)$  matrix

$$A_{\text{LP}} = \begin{bmatrix} R_b^*(\alpha_1) & 0 & \cdots & 0 & \cdots & -1 \\ R_b^*(\alpha_2) & R_b^*(\alpha_2) & \cdots & \vdots & \cdots & \vdots \\ \vdots & \vdots & \ddots & \vdots & \cdots & 0 \\ R_b^*(\alpha_M) & \cdots & \cdots & R_b^*(\alpha_M) & \cdots & -1 \\ \hline R_f^*(\alpha_1) & \cdots & \cdots & R_f^*(\alpha_1) & \cdots & 0 \\ 0 & R_f^*(\alpha_2) & \cdots & R_f^*(\alpha_2) & \cdots & \vdots \\ \vdots & \vdots & \ddots & \vdots & \cdots & \vdots \\ 0 & \cdots & 0 & R_f^*(\alpha_M) & \cdots & 0 \end{bmatrix}. \quad (26)$$

Moreover,  $v_b(\gamma_f, \alpha)$  is decreasing in  $\gamma_f$ .

*Proof.* We need to show equivalence between the linear program (25) and the optimization problem

$$\min_{\mathbf{p}' \in \Delta_M} \mathbb{E}_{\alpha, \mathbf{p}'}^*[R_b] \text{ subject to } \mathbb{E}_{\alpha, \mathbf{p}'}^*[R_f] \leq \gamma_f. \quad (27)$$

Let  $A_b(\alpha) \in \mathbb{R}^{M \times M}$  be the matrix defined by the upper left block of (26) and  $A_f(\alpha)$  by the bottom left block. From Theorem 3, we can express  $\mathbb{E}_{\alpha, \mathbf{p}'}^*[R_b]$  as the maximum element of the  $M$ -vector  $A_b(\alpha)\mathbf{p}'$ , and similarly  $\mathbb{E}_{\alpha, \mathbf{p}'}^*[R_f]$  as the maximum element of  $A_f(\alpha)\mathbf{p}'$ . Hence,  $\mathbb{E}_{\alpha, \mathbf{p}'}^*[R_f] \leq \gamma_f$  is the linear constraint  $[A_f(\alpha)\mathbf{p}']_i \leq \gamma_f$  for all  $i = 1, \dots, M$ . The objective  $\min_{\mathbf{p}' \in \Delta_M} \mathbb{E}_{\alpha, \mathbf{p}'}^*[R_b]$  itself can be cast as a linear objective with linear constraints, i.e.  $\min_{\mathbf{p}' \in \Delta_M, v \in [0, 1]} v$  s.t.  $[A_b(\alpha)\mathbf{p}']_i \leq v$ . Combining these two, we obtain (25). The claim  $v_b(\gamma, \alpha)$  is decreasing in  $\gamma$  follows as a consequence of the linear program (25). ■

We note that a worst-case expected focused risk  $\mathbb{E}_{\alpha, \mathbf{p}}^*[R_f] < R_f^*(\alpha_1)$  is not attainable because  $\alpha_1$  is the smallest gain it mixes with. Hence, the linear program (25) is infeasible for  $\gamma_f < R_f^*(\alpha_1)$ . The following tradeoff relation holds in the opposite direction.

### Risk-security tradeoff

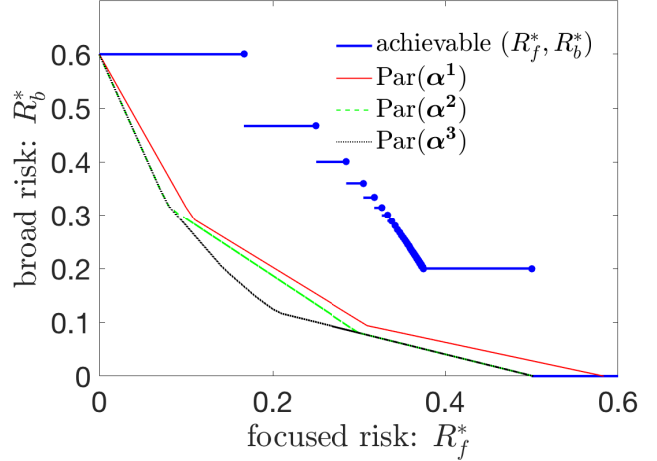


Fig. 3: Security-risk tradeoffs are depicted by the achievable worst-case risk levels from deterministic gains (blue) and randomized gains (red, green, black). The Pareto frontiers for three different randomized strategies  $\alpha^1, \alpha^2 \in \mathbb{R}_+^5$ , and  $\alpha^3 \in \mathbb{R}_+^{300}$ , are shown in increasing order of improvement. The strategies  $\alpha^1$  and  $\alpha^2$  randomize over the highest three broad risk levels in addition to the lowest two. The strategy  $\alpha^3$  randomizes over the highest 298 broad risk levels and the lowest two. We chose the values as follows. For  $k = 1, 2$ , we set  $\alpha_1^k = \alpha_{\text{sys}}$ ,  $\alpha_j^k = (1 - \epsilon_k) \frac{j-1}{j} + \epsilon_k \frac{j}{j+1} \in I_j$  for  $j = 2, 3$ ,  $\alpha_4^k = 1 + \epsilon_k$ , and  $\alpha_5^k = \frac{3}{2} + \epsilon_k$ . We have set  $\epsilon_1 = 0.5$  and  $\epsilon_2 = .01$ . Hence,  $\text{Par}(\alpha^2)$  improves upon  $\text{Par}(\alpha^1)$  via Claim 1. For  $k = 3$ , we set  $\alpha_1^3 = \alpha_{\text{sys}}$ ,  $\alpha_j^3 = (1 - \epsilon_3) \frac{j-1}{j} + \epsilon_3 \frac{j}{j+1} \in I_j$ ,  $j = 2, 3, \dots, 298$ ,  $\alpha_{299}^3 = 1 + \epsilon_3$ , and  $\alpha_5^3 = \frac{3}{2} + \epsilon_3$ . Claim 2 ensures  $\text{Par}(\alpha^3)$  improves upon  $\text{Par}(\alpha^2)$ . We chose  $\epsilon_3 = .01$  and  $\alpha_{\text{sys}} = 1/4$ .

**Corollary 4.** Fix  $\gamma_b \in [R_b^*(\alpha_M), 1]$  and a set of gains  $\alpha = \{\alpha_j\}_{j=1}^M$ . Suppose  $\mathbb{E}_{\alpha, \mathbf{p}}^*[R_b] \leq \gamma_b$  for some  $\mathbf{p} \in \Delta_M$ . Then

$$\mathbb{E}_{\alpha, \mathbf{p}}^*[R_f] \geq v_f(\gamma_b, \alpha), \quad (28)$$

where  $v_f(\gamma_b, \alpha)$  is the value of the following linear program.

$$\begin{aligned} v_f(\gamma_b, \alpha) = \min_{\mathbf{p}, v} v \\ \text{s.t. } \sum_{i=1}^M p_i = 1, p_i \geq 0 \forall i = 1, \dots, M \\ v \in [0, 1] \\ \begin{bmatrix} A_f(\alpha) & \vdots & -\mathbb{1}_M \\ \hline A_b(\alpha) & \vdots & 0_M \end{bmatrix} \begin{bmatrix} \mathbf{p} \\ v \end{bmatrix} \preceq \begin{bmatrix} 0_M \\ \gamma_b \mathbb{1}_M \end{bmatrix}, \end{aligned} \quad (29)$$

where  $A_f(\alpha)$  and  $A_b(\alpha)$  are defined as the bottom and top left blocks of (26), respectively. Furthermore,  $v_f(\gamma_b, \alpha)$  is decreasing in  $\gamma_b$ .

We omit the proof as it is similar to that of Corollary 3. Note a worst-case expected broad risk  $\mathbb{E}_{\alpha, \mathbf{p}}^*[R_b] < R_b^*(\alpha_M)$  is not attainable since  $\alpha_M$  is the highest gain it mixes with - (29) is infeasible for  $\gamma_b < R_b^*(\alpha_M)$ . Fig. 3 plots the best achievable risk levels of three randomized operator strategies (red, green, and black).

### C. Improvement of risk tradeoffs

The tradeoff relations describe the best achievable level on one risk metric given the other is subject to a security

constraint when the gains  $\alpha$  are fixed. One way to improve the achievable risks is to decrease the available gains.

**Claim 1.** Let  $\alpha, \alpha' \in \mathbb{R}^M$ . Suppose  $\alpha_j \in I_{k_j}$  (recall (12)),  $j = 1, \dots, M$  for some non-decreasing subsequence  $k_j \geq 1$ . Let  $\alpha'$  satisfy  $\alpha'_j \in I_{k_j}$  with  $\alpha'_j < \alpha_j$ . Then for all  $\gamma_b \in [R_b^*(\alpha_M), 1]$ ,  $v_f(\gamma_b, \alpha') \leq v_f(\gamma_b, \alpha)$ . Similarly, for all  $\gamma_f \in [R_f^*(\alpha_1), 1]$ ,  $v_b(\gamma_f, \alpha') \leq v_b(\gamma_f, \alpha)$ .

Randomizing over additional gains can also improve the achievable risks.

**Claim 2.** Suppose  $\alpha \in \mathbb{R}^M$  and  $\alpha' \in \mathbb{R}^{M'}$  with  $M < M'$ , and assume  $\alpha'$  contains the elements of  $\alpha$ . Then the assertion of Claim 1 holds.

The proofs of the above two Claims follow directly from the formulation of the LPs (25), (29), and hence we omit them.

Fig. 3 depicts the best achievable risk levels of three randomized operator strategies of increasing improvement due to Claims 1 and 2 (red, green, and black curves). In particular, these plots constitute the *Pareto frontier* of all attainable expected risks among distributions  $\mathbf{p}$  given a fixed set of gains.

That is, for any  $\alpha$ , we say a risk level  $\begin{bmatrix} \mathbb{E}_{\alpha, \mathbf{p}}[R_f] \\ \mathbb{E}_{\alpha, \mathbf{p}}[R_b] \end{bmatrix} \in \mathbb{R}^2$  belongs to the frontier  $\text{Par}(\alpha)$  if there does not exist a  $\mathbf{p}' \neq \mathbf{p}$  such that  $\begin{bmatrix} \mathbb{E}_{\alpha, \mathbf{p}'}[R_f] \\ \mathbb{E}_{\alpha, \mathbf{p}'}[R_b] \end{bmatrix} \prec \begin{bmatrix} \mathbb{E}_{\alpha, \mathbf{p}}[R_f] \\ \mathbb{E}_{\alpha, \mathbf{p}}[R_b] \end{bmatrix}$ . Within  $\text{Par}(\alpha)$ , the operator can only improve upon one worst-case risk metric by sacrificing performance on the other.

From Corollary 4, the frontier given gains  $\alpha$  is the set of points

$$\text{Par}(\alpha) = \left\{ \begin{bmatrix} v_f(\gamma_b, \alpha) \\ \gamma_b \end{bmatrix} \in \mathbb{R}^2 : \gamma_b \in [R_b^*(\alpha_M), R_b^*(\alpha_1)] \right\}. \quad (30)$$

The parameter  $\gamma_b$  is upper bounded here by  $R_b^*(\alpha_1)$  since any risk level with  $\mathbb{E}_{\alpha, \mathbf{p}}[R_b] > R_b^*(\alpha_1)$  is unattainable under  $\alpha$ . Hence, the values  $v_f(\gamma_b, \alpha)$  and  $v_f(R_b^*(\alpha_1), \alpha)$  are equivalent for  $\gamma_b > R_b^*(\alpha_1)$ . The frontiers in Fig. 3 are generated by numerically solving the linear program (29) for a finite grid of points  $\gamma_b \in [R_b^*(\alpha_M), R_b^*(\alpha_1)]$ .

As we have seen, the transition from deterministic to randomized gains ensures a reduction of risk levels. Randomizing over only a few different gains substantially improves upon the attainable deterministic worst-case risks. However, a detailed quantification of such improvements remains a challenge due to the high dimensionality of the model. In particular, we have yet identified a “limit” frontier that could be obtained by repeated modifications to the gain vector detailed by Claims 1 and 2.

## V. PROOF OF THEOREMS 1 AND 2: DETERMINISTIC WORST-CASE RISKS

In this section, we develop the technical results that characterize the worst-case risk metrics  $R_b^*(\alpha)$  and  $R_f^*(\alpha)$  (Theorems 1 and 2). Before presenting the proofs, we first present some preliminaries on potential games [25], which are essential to calculating stochastically stable states. We then define relevant notations for the forthcoming analysis.

### A. Potential games

Graphical coordination games fall under the class of potential games - games where individual utilities  $\{U_i\}_{i \in \mathcal{N}}$  are aligned with a global objective, or potential function. A game is a potential game if there exists a potential function  $\phi : \mathcal{A} \rightarrow \mathbb{R}$  which satisfies

$$\phi(a_i, a_{-i}) - \phi(a'_i, a_{-i}) = U_i(a_i, a_{-i}) - U_i(a'_i, a_{-i}) \quad (31)$$

for all  $i \in \mathcal{N}$ ,  $a \in \mathcal{A}$ , and  $a'_i \neq a_i$  [25]. In potential games, the set of stochastically stable states (6) are precisely the action profiles that maximize the potential function [15], [16]. Specifically,  $\text{LLL}(\mathcal{A}, \{U_i\}_{i \in \mathcal{N}}; G) = \arg \max_{a \in \mathcal{A}} \phi(a)$ . Our analysis relies on characterizing a potential function for the graphical coordination game in the presence of adversarial influences. This allows us to compute stochastically stable states in a straightforward manner.

### B. Relevant notations for analysis

Any action profile  $a$  on a graph  $G = (\mathcal{N}, \mathcal{E}) \in \mathcal{G}_N$  decomposes  $\mathcal{N}$  into  $x$  and  $y$ -partitions. A node that belongs to a  $y$ -partition ( $x$ -partition) has  $a_i = y$  ( $x$ ). The partitions are enumerated  $\{\mathcal{P}_y^1, \dots, \mathcal{P}_y^{k_y}\}$  and  $\{\mathcal{P}_x^1, \dots, \mathcal{P}_x^{k_x}\}$ , are mutually disjoint, and cover the graph. Each partition is a connected subgraph of  $G$ . It is possible that  $k_x = 0$  with  $k_y = 1$  (when  $a = \vec{y}$ ),  $k_x = 1$  with  $k_y = 0$  (when  $a = \vec{x}$ ), or  $k_y, k_x \geq 1$ .

For any subset of nodes  $A, B \subseteq \mathcal{N}$ , let us denote

$$e(A, B) := \{(i, j) \in \mathcal{E} : i \in A, j \in B\} \quad (32)$$

as the set of edges between  $A$  and  $B$ . We write  $A^c$  as the complement of  $A$ . We extensively use the notation

$$W^\alpha(E, a) := \sum_{(i, j) \in E} V^\alpha(a_i, a_j) \quad (33)$$

as the welfare due to edge set  $E \subseteq \mathcal{E}$  in action profile  $a$ , where  $V^\alpha$  is of the form (1) with  $\alpha_{\text{sys}}$  replaced by  $\alpha$ . For compactness, we will denote  $W(E, a)$  as  $W^{\alpha_{\text{sys}}}(E, a)$  for the local system welfare generated by the edges  $E$ . Our analysis will also rely on the following median inequality.

**Fact 1.** Suppose  $n_i \geq 0$  and  $d_i > 0$  for each  $i = 1, \dots, m \in \mathbb{N}$ . Then

$$\frac{\sum_{i=1}^m n_i}{\sum_{i=1}^m d_i} \geq \min_i \frac{n_i}{d_i}. \quad (34)$$

We refer to the LHS above as the median sum of the  $\frac{n_i}{d_i}$ .

### C. Characterization of $R_b^*$ : worst-case broad risk

To prove Theorem 1, we seek a pair  $(S, G)$  with  $G \in \mathcal{G}_N$  of any size  $N \geq 3$  and  $S \in \mathcal{T}(G)$ , that minimizes efficiency  $J_b(\alpha, S; G)$  (maximizes risk  $R_b(\alpha, S; G)$ ). Our method to find the minimizer is to show any  $(S, G)$  can be transformed into a star network with a particular target set that has lower efficiency, when  $\alpha < 1$ . Thus, in this regime the search for the worst-case graph reduces to the class of star networks of arbitrary size. For  $\alpha \geq 1$ , structural properties allow us to deduce the minimal efficiency.



The graphical coordination game defined by  $\mathcal{A} = \{x, y\}^N$ , perceived utilities  $\{\tilde{U}_i^\alpha\}_{i \in \mathcal{N}}$  (7), target set  $S$ , and graph  $G$  falls under the class of potential games [25]. A potential function is given by

$$\frac{1}{2}W^\alpha(a) + (1 + \alpha) \sum_{i \in S_x} \mathbb{1}(a_i = x) + \sum_{i \in S_y} \mathbb{1}(a_i = y) \quad (35)$$

where

$$W^\alpha(a) := \sum_{i \in \mathcal{N}} U_i^\alpha(a). \quad (36)$$

Hence, the stochastically stable states  $\text{LLL}(\mathcal{A}, \alpha, S; G)$  are maximizers of (35). Suppose  $\hat{a} = \arg \min_{a \in \text{LLL}(\mathcal{A}, \alpha, S; G)} W(a)$  is

the welfare-minimizing SSS inducing the partitions  $\{\mathcal{P}_z^k\}_{k=1}^{k_z}$ ,  $z = x, y$ . We can express its efficiency from (8) as

$$\frac{\sum_{k=1}^{k_y} |e(\mathcal{P}_y^k, \mathcal{P}_y^k)| + (1 + \alpha_{\text{sys}}) \sum_{k=1}^{k_x} |e(\mathcal{P}_x^k, \mathcal{P}_x^k)|}{(1 + \alpha_{\text{sys}})(\sum_{k=1}^{k_y} |e(\mathcal{P}_y^k, \mathcal{N})| + \sum_{k=1}^{k_x} |e(\mathcal{P}_x^k, \mathcal{P}_x^k)|)}. \quad (37)$$

Note the denominator is simply the number of edges in  $G$  multiplied by  $1 + \alpha_{\text{sys}}$ . From (35), each  $y$ -partition  $\mathcal{P}_y^k$  in  $\hat{a}$  satisfies<sup>1</sup>

$$|\mathcal{P}_y^k| + |e(\mathcal{P}_y^k, \mathcal{P}_y^k)| \geq \max_{a_{\mathcal{P}_y^k} \neq \vec{y}_{\mathcal{P}_y^k}} W^\alpha(e(\mathcal{P}_y^k, \mathcal{N}), (a_{\mathcal{P}_y^k}, \hat{a}_{-\mathcal{P}_y^k})) + \sum_{i \in \mathcal{P}_y^k} \mathbb{1}(a_i = y). \quad (\text{CY})$$

In words, no subset of agents in  $\mathcal{P}_y^k$  can deviate from  $y$  to improve the collective perceived welfare of  $\mathcal{P}_y^k$ . A similar stability condition holds for each  $x$ -partition  $\mathcal{P}_x^k$ .

$$(1 + \alpha)(|\mathcal{P}_x^k| + |e(\mathcal{P}_x^k, \mathcal{P}_x^k)|) \geq \max_{a_{\mathcal{P}_x^k} \neq \vec{x}_{\mathcal{P}_x^k}} W^\alpha(e(\mathcal{P}_x^k, \mathcal{N}), (a_{\mathcal{P}_x^k}, \hat{a}_{-\mathcal{P}_x^k})) + (1 + \alpha) \sum_{i \in \mathcal{P}_x^k} \mathbb{1}(a_i = x) \quad (\text{CX})$$

The following result characterizes the threshold on  $\alpha$  above which any network is safeguarded from any imposter attack.

**Lemma 1.** *Let  $N \geq 3$ . Then  $\alpha > \frac{N}{N-1}$  if and only if*

$$\min_{G \in \mathcal{G}_N} \min_{S \in \mathcal{T}(G)} J_b(\alpha, S; G) = 1. \quad (38)$$

*Proof.* ( $\Rightarrow$ ) Let  $\alpha > \frac{N}{N-1}$ . Suppose there is a pair  $(S, G)$  with  $J_b(\alpha, G, S) < 1$ . Then there must exist a  $y$ -partition  $\mathcal{P}_y \subset \mathcal{N}$ . From (CY),

$$|\mathcal{P}_y| + |e(\mathcal{P}_y, \mathcal{P}_y)| \geq (1 + \alpha)|e(\mathcal{P}_y, \mathcal{N})| > 2|e(\mathcal{P}_y, \mathcal{N})|. \quad (39)$$

Since  $G$  is connected,  $|e(\mathcal{P}_y, \mathcal{P}_y)| \geq |\mathcal{P}_y| - 1$  and there is at least one outgoing link from  $\mathcal{P}_y$ , i.e.  $|e(\mathcal{P}_y, \mathcal{P}_y^c)| \geq 1$ . Consequently,  $|e(\mathcal{P}_y, \mathcal{N})| \geq |\mathcal{P}_y|$ , from which we obtain

$$|e(\mathcal{P}_y, \mathcal{N})| + |e(\mathcal{P}_y, \mathcal{P}_y)| > 2|e(\mathcal{P}_y, \mathcal{N})|. \quad (40)$$

<sup>1</sup>Since we are seeking worst-case pairs  $(S, G)$ , we may consider any  $y$ -partition as only having  $y$  imposters placed among its nodes. This is because any  $x$  imposters that were placed in a resulting  $y$ -partition can be replaced by  $y$ -imposters and retain stability. We reflect this generalization in (CY) and (CX), where influence from only  $y$  ( $x$ ) imposters is considered.

which is impossible.

( $\Leftarrow$ ) Assume  $\min_{G \in \mathcal{G}_N} \min_{S \in \mathcal{T}(G)} J_b(\alpha, S; G) = 1$ . Then no  $y$ -partition can exist for any graph. In particular, (CY) is violated for  $\mathcal{P}_y = \mathcal{N}$ .

$$N + |\mathcal{E}| < (1 + \alpha)|\mathcal{E}| \Rightarrow \alpha > \frac{N}{|\mathcal{E}|}. \quad (41)$$

Since  $|\mathcal{E}| \geq N - 1$ , we obtain  $\alpha > \frac{N}{N-1}$ . ■

We also deduce the following minimal efficiencies for any graph when  $1 \leq \alpha \leq \frac{N}{N-1}$ .

**Lemma 2.** *Suppose  $N \geq 3$ . Then  $\alpha \in [1, \frac{N}{N-1}]$  if and only if*

$$\min_{G \in \mathcal{G}_N} \min_{S \in \mathcal{T}(G)} J_b(\alpha, S; G) = \frac{1}{1 + \alpha_{\text{sys}}}. \quad (42)$$

*Proof.* The ( $\Rightarrow$ ) direction follows the same argument as Lemma 1.

( $\Leftarrow$ ) The assumption implies the only  $y$ -partition that is stabilizable is  $\mathcal{N}$ . Then for any  $\mathcal{P}_y \subset \mathcal{N}$ , (CY) is violated, i.e.

$$|\mathcal{P}_y| + |e(\mathcal{P}_y, \mathcal{P}_y)| < (1 + \alpha)|e(\mathcal{P}_y, \mathcal{N})|. \quad (43)$$

Since  $G$  is connected and there is at least one outgoing edge from  $\mathcal{P}_y$ , we obtain

$$\frac{2|\mathcal{P}_y| - 1}{|\mathcal{E}|} < 1 + \alpha \quad (44)$$

The above holds for any graph  $G = (\mathcal{N}, \mathcal{E})$  and subset of nodes  $\mathcal{P}_y \subset \mathcal{N}$ . From the facts that  $|\mathcal{P}_y| \leq N - 1$  and  $|\mathcal{E}| \geq N - 1$ , we have  $\alpha > \frac{N-2}{N-1}$  for any  $N \geq 3$ . Consequently,  $\alpha \geq 1$  and Lemma 1 establishes that  $\alpha \leq \frac{N}{N-1}$ . ■

The class of star graphs is central to the worst-case analysis in the interval  $0 < \alpha < 1$ .

**Definition 3.** *Let  $\mathbb{S}_N$  be the set of all  $(S, G)$  where  $G$  is the star graph with  $N$  nodes,  $S_y$  contains the center node, and  $S_x = \mathcal{N} \setminus S_y$ .*

An immediate consequence of this definition is the leaf nodes satisfy (CX). The efficiency is then proportional to the fraction of leaf nodes that are stable to  $y$ , if any. Furthermore, the stability condition (CY) of  $\mathcal{P}_y = S_y$  for members of  $\mathbb{S}_N$  simplifies to

$$2|e(\mathcal{P}_y, \mathcal{P}_y)| + 1 \geq (1 + \alpha)(N - 1). \quad (45)$$

In other words, stability of the target set  $S_y$  as a  $y$ -partition hinges on (CY) being satisfied for the selection  $a_{\mathcal{P}_y} = \vec{x}$ . The following result reduces the search space for efficiency minimizers to  $\mathbb{S}_N$  when  $\alpha < 1$ .

**Lemma 3.** *Suppose  $0 < \alpha < 1$  and  $n \geq 3$ . Consider any  $(S, G)$  with  $G \in \mathcal{G}_N$ ,  $S \in \mathcal{T}(G)$ . Then there is a  $(S', G') \in \mathbb{S}_{N'}$  such that  $J_b(\alpha, S'; G') \leq J_b(\alpha, S; G)$  for some  $N' \geq N$ .*

The idea of the proof is to construct a member of  $\mathbb{S}_{N'}$  by recasting the  $y$  and  $x$ -partitions of  $(S, G)$  as star subgraphs while preserving the same number and type of edges, thus preserving efficiency. Further efficiency reduction can be achieved by converting excess  $x$  links into  $y$  links in this star configuration. We provide the proof detailing the constructive procedure in



the Appendix. We now characterize the minimal efficiency for the star graph of size  $N$ ,  $J_N^*(\alpha) := \min_{(G,S) \in \mathbb{S}_N} J_b(\alpha, S; G)$  for  $\alpha < 1$ .

**Lemma 4.** *Suppose  $\alpha < 1$  and fix  $N \geq 3$ . Then*

$$J_N^*(\alpha) = \frac{1}{(1 + \alpha_{\text{sys}})(N - 1)} \left\lceil \frac{(1 + \alpha)(N - 1) - 1}{2} \right\rceil. \quad (46)$$

*Proof.* The goal is to find the smallest  $y$ -partition of the  $n$  star that is still stabilizable under a gain  $\alpha$ . This is written

$$J_N^*(\alpha) = \min_{N_y} \frac{1}{1 + \alpha_{\text{sys}}} \frac{N_y}{N - 1} \quad \text{(size of } y\text{-partition)}$$

$$\text{s.t. } \begin{cases} N_y \leq N - 1 \\ 2N_y + 1 \geq (1 + \alpha)(N - 1) \end{cases} \quad \text{(stability)} \quad (47)$$

The smallest integer  $N_y$  that satisfies the constraints is  $\left\lceil \frac{(1 + \alpha)(N - 1) - 1}{2} \right\rceil$  for  $\alpha \in (0, 1)$ . ■

*Proof of Theorem 1.* For  $\alpha < 1$ , by Lemma 3, the worst-case efficiency is

$$\min_{N \geq 3} \min_{(G,S) \in \mathbb{S}_N} J_b(\alpha, S; G) = \min_{N \geq 3} J_N^*(\alpha). \quad (48)$$

Using the formula of Lemma 4, we obtain the first entry in (11). Lemma 2 asserts the minimal efficiency is  $\frac{1}{1 + \alpha_{\text{sys}}}$  for  $\alpha \in [1, \frac{3}{2}]$  because the upper bound  $\frac{N}{N - 1}$  is maximized at  $N = 3$  (for  $N \geq 3$ ). This gives the second entry in (11). Lastly, Lemma 1 asserts the minimal efficiency is 1 for  $\alpha > \frac{3}{2}$ . ■

#### D. Characterization of $R_T^*$ : worst-case focused risk

Our approach for the proof of Theorem 2 differs from that of  $R_b^*$ . Instead of reducing the search of worst-case graphs, we simply provide an upper bound on  $R_T^*(\alpha, F; G)$  for any  $G$  and fixed set  $F \in \mathcal{F}(G)$ , and show one can construct a graph with fixed nodes that achieves it. .

We observe  $\frac{1}{2}W^\alpha(a) : \mathcal{A}(F) \rightarrow \mathbb{R}$  serves as a potential function (recall (36)) for the game with restricted action set  $\mathcal{A}(F)$  and utilities  $\{U_i^\alpha\}_{i \in \mathcal{N}}$ . Hence, the stochastically stable states  $\text{LLL}(\mathcal{A}(F), \alpha; G)$  are maximizers of  $\frac{1}{2}W^\alpha(a)$ . Suppose  $\hat{a} = \arg \min_{a \in \text{LLL}(\mathcal{A}(F), \alpha; G)} W(a)$  decomposes the graph into the  $x$  and  $y$ -partitions  $\{\mathcal{P}_z^k\}_{k=1}^{k_z}$ ,  $z = x, y$ . We express its efficiency (13) as

$$\frac{\sum_{k=1}^{k_y} |e(\mathcal{P}_y^k, \mathcal{P}_y^k)| + (1 + \alpha_{\text{sys}}) \sum_{k=1}^{k_x} |e(\mathcal{P}_x^k, \mathcal{P}_x^k)|}{\sum_{k=1}^{k_y} W^{\alpha_{\text{sys}}}(e(\mathcal{P}_y^k, \mathcal{N}), a^*) + \sum_{k=1}^{k_x} W^{\alpha_{\text{sys}}}(e(\mathcal{P}_x^k, \mathcal{P}_x^k), a^*)} \quad (49)$$

where  $a^* = \arg \max_{a \in \mathcal{A}(F)} W(a)$  is the welfare-maximizing action profile. Similar to (CY), each  $y$ -partition  $\mathcal{P}_y^k$  formed from  $\hat{a}$  satisfies the stability condition

$$|e(\mathcal{P}_y^k, \mathcal{P}_y^k)| \geq \max_{a_{\mathcal{P}_y^k} \neq \vec{y}} W^\alpha(e(\mathcal{P}_y^k, \mathcal{N}), (a_{\mathcal{P}_y^k}, \hat{a}_{-\mathcal{P}_y^k})). \quad \text{(CYE)}$$

To reduce cumbersome notation, it is understood the max is taken over actions of unfixed nodes,  $a_{\mathcal{P}_y^k \setminus F}$ . Likewise, each  $x$ -partition  $\mathcal{P}_x^k$  satisfies

$$(1 + \alpha)|e(\mathcal{P}_x^k, \mathcal{P}_x^k)| \geq \max_{a_{\mathcal{P}_x^k} \neq \vec{x}} W^\alpha(e(\mathcal{P}_x^k, \mathcal{N}), (a_{\mathcal{P}_x^k}, \hat{a}_{-\mathcal{P}_x^k})). \quad \text{(CXE)}$$

The following lemma asserts that agents playing  $y$  in the SSS under the gain  $\alpha$  remain playing  $y$  under a lower gain  $\alpha' < \alpha$ . The result is crucial for establishing a lower bound on efficiency for any graph  $G$  with arbitrary fixed set  $F \in \mathcal{F}(G)$ .

**Lemma 5.** *Suppose  $\alpha' < \alpha$ . Denote  $\hat{a}' = \arg \min_{a \in \text{LLL}(\mathcal{A}(F), \alpha'; G)} W(a)$  as the welfare-minimizing SSS under  $\alpha'$ . Then for any  $y$ -partition  $\mathcal{P}_y$  induced from  $\alpha$ ,  $\hat{a}'_i = y$  for all  $i \in \mathcal{P}_y \setminus F$ .*

*Proof.* Condition (CYE) asserts for all  $a_{\mathcal{P}_y} \neq \vec{y}$  that

$$W^\alpha(e(\mathcal{P}_y, \mathcal{N}), (\vec{y}_{\mathcal{P}_y}, \hat{a}_{-\mathcal{P}_y})) \geq W^\alpha(e(\mathcal{P}_y, \mathcal{N}), (a_{\mathcal{P}_y}, \hat{a}_{\mathcal{P}_y})). \quad (50)$$

It also holds for all  $a_{\mathcal{P}_y} \neq \vec{y}$  and for any  $a_{-\mathcal{P}_y} \neq \hat{a}_{-\mathcal{P}_y}$  that

$$W^\alpha(e(\mathcal{P}_y, \mathcal{N}), (\vec{y}_{\mathcal{P}_y}, a_{-\mathcal{P}_y})) \geq W^\alpha(e(\mathcal{P}_y, \mathcal{N}), (a_{\mathcal{P}_y}, a_{-\mathcal{P}_y})) \quad (51)$$

because any  $y$ -links garnered in the RHS above by changing  $\hat{a}_{-\mathcal{P}_y}$  to  $a_{-\mathcal{P}_y}$  also contribute to the LHS. In particular, the above holds for  $a_{-\mathcal{P}_y} = \hat{a}'_{-\mathcal{P}_y}$ . Lowering the gain to  $\alpha'$  preserves the above inequality as well, as it de-values  $x$ -links garnered on the RHS. ■

A dual statement holds - agents playing  $x$  in the SSS under  $\alpha$  remain so under a higher gain  $\alpha' > \alpha$ .

**Lemma 6.** *Suppose  $\alpha' > \alpha$ . Then for any  $x$ -partition  $\mathcal{P}_x$  induced from  $\alpha$ ,  $\hat{a}'_i = x$  for all  $i \in \mathcal{P}_x \setminus F$ .*

We omit the proof for brevity, as it is analogous to the proof of Lemma 5. We are now ready to prove Theorem 2.

*Proof of Theorem 2.* Consider any graph  $G \in \mathcal{G}_N$  with fixed set  $F$ . Recall that efficiency is one for  $\alpha = \alpha_{\text{sys}}$ . Thus, we first consider  $\alpha < \alpha_{\text{sys}}$ . Observe that

$$|e(\mathcal{P}_y^k, \mathcal{P}_y^k)| \geq W^\alpha(e(\mathcal{P}_y, \mathcal{N}), (a_{\mathcal{P}_y}^*, \hat{a}_{-\mathcal{P}_y})) = W^\alpha(e(\mathcal{P}_y, \mathcal{N}), (a_{\mathcal{P}_y}^*, a_{-\mathcal{P}_y}^*)) \quad (52)$$

where the inequality is due to (CYE). The equality results from Lemma 6 - the agents ( $\notin \mathcal{P}_y$ ) that neighbor any member of  $\mathcal{P}_y$  remain playing  $x$  in  $a^*$ . We then obtain

$$\frac{|e(\mathcal{P}_y^k, \mathcal{P}_y^k)|}{W(e(\mathcal{P}_y^k, \mathcal{N}), a^*)} \geq \frac{1 + \alpha}{1 + \alpha_{\text{sys}}}. \quad (53)$$

The inequality results since the expressions of the numerator and denominator garner the same edges for welfare. It occurs with equality if and only if  $a_i^* = x \forall i \in \mathcal{P}_y \setminus F$ . Applying the median inequality (34) to (49),  $J_T(\alpha, F; G) \geq \frac{1 + \alpha}{1 + \alpha_{\text{sys}}}$ . The case when  $\alpha > \alpha_{\text{sys}}$  follows analogous arguments.

From Lemma 5,  $|e(\mathcal{P}_y^k, \mathcal{P}_y^k)| = W^\alpha(e(\mathcal{P}_y^k, \mathcal{P}_y^k), a^*)$ . For  $x$ -partitions,

$$\begin{aligned} \frac{(1+\alpha_{\text{sys}})|e(\mathcal{P}_x^k, \mathcal{P}_x^k)|}{W(e(\mathcal{P}_x^k, \mathcal{N}), a^*)} &\geq \frac{1+\alpha_{\text{sys}}}{1+\alpha} \frac{W^\alpha(e(\mathcal{P}_x^k, \mathcal{N}), a^*)}{W(e(\mathcal{P}_x^k, \mathcal{N}), a^*)} \\ &\geq \frac{1+\alpha_{\text{sys}}}{1+\alpha} \end{aligned} \quad (54)$$

where the first inequality is from (CXE) and the second occurs with equality if  $a_i^* = y \forall i \in \mathcal{P}_x^k \setminus F$ . From (34) and (49),  $J_f(\alpha, F; G) \geq \frac{1+\alpha_{\text{sys}}}{1+\alpha}$ . ■

We have just shown fundamental lower bounds on efficiency for any graph with fixed agents. The bounds are tight as they can be achieved for any gain  $\alpha$  by arranging  $N_x$  fixed  $x$  and  $N_y$  fixed  $y$  leaf nodes that influence a single unfixed agent in the center of a star graph. If  $\alpha < \alpha_{\text{sys}}$ , choosing  $\frac{N_x}{N_y} = \frac{1}{1+\alpha}$  gives the minimal efficiency  $\frac{1+\alpha}{1+\alpha_{\text{sys}}}$ . If  $\alpha > \alpha_{\text{sys}}$ , choosing  $\frac{N_x}{N_y} = \frac{1}{1+\alpha}$  gives the minimal efficiency  $\frac{1+\alpha_{\text{sys}}}{1+\alpha}$ . Note that if  $\alpha$  is rational, one could choose finite integers  $N_y, N_x$  that achieve such ratios. Recall Figure 2b for illustrative examples. However if it is irrational, they must be taken arbitrarily large to better approximate the ratio.

## VI. PROOF OF THEOREM 3: WORST-CASE RISKS UNDER RANDOMIZED OPERATOR DESIGNS

Recall a randomized strategy consists of gains  $\alpha = \{\alpha_i\}_{i=1}^M$  with distribution  $\mathbf{p} \in \Delta_M$ . The gains are ordered  $\alpha_{\text{sys}} \leq \alpha_1 < \dots < \alpha_M$ . To prove Theorem 3, we outline a few technical Lemmas. The key insight is the expected efficiency of any graph  $G$  can be expressed in the form  $\sum_{i=1}^M p_i s_i$ , where the coefficient  $s_i$  is a mediant sum over local efficiencies of partitions in  $G$  when gain  $\alpha_i$  is used. The following two mathematical facts are the basis of this insight.

**Fact 2.** Let  $\nu_i < \frac{n_i}{d_i} \leq 1$  with  $r_i \geq 0$  and  $n_i, d_i > 0$  for all  $i = 1, \dots, M$ . Then for all  $\mathbf{p} \in \Delta_M$ ,

$$\sum_{i=1}^M p_i s_i \geq 1 + \min_{i=1, \dots, M} \left\{ \left( \sum_{j=1}^i p_j \right) (\nu_i - 1) \right\} \quad (55)$$

where  $s_i := \frac{\sum_{j=1}^{i-1} d_j + \sum_{j=i}^M n_j}{\sum_{j=1}^M d_j}$ ,  $i = 1, \dots, M$ .

We provide a proof in the Appendix. The following dual result follows directly.

**Fact 3.** For all  $\mathbf{p} \in \Delta_M$ ,

$$\sum_{i=1}^M p_i s'_i \geq 1 + \min_{i=1, \dots, M} \left\{ \left( \sum_{j=i}^M p_j \right) (\nu_i - 1) \right\} \quad (56)$$

where  $s'_i := \frac{\sum_{j=1}^i n_j + \sum_{j=i}^M d_j}{\sum_{j=1}^M d_j}$ ,  $i = 1, \dots, M$ .

*Proof.* The proof follows similarly to Fact 2, where the indices of the  $s_i$  coefficients are reversed. ■

We will show for any  $(S, G)$  that  $\mathbb{E}_{\alpha, \mathbf{p}}[J_b|S, G] = 1 - \mathbb{E}_{\alpha, \mathbf{p}}[R_b|S, G]$  can be expressed in the form  $\sum_{i=1}^M p_i s_i$  from the LHS of (55). The lower bounds establish worst-case expected efficiencies - and hence risks. The  $\nu_i$  correspond to the

worst-case deterministic efficiencies  $J_b^*(\alpha_i) = 1 - R_b^*(\alpha_i)$  of the  $M$  gains and  $\frac{n_i}{d_i}$  to local efficiencies of selected partitions in the graph. Fact 2 will be used to establish (22), and Fact 3 for (23) (Theorem 3). We now identify a structural property required of worst-case graphs.

**Lemma 7.** A worst-case graph, i.e. a member of  $\arg \min_{G \in \mathcal{G}_N, S \in \mathcal{T}(G)} \mathbb{E}_{\alpha, \mathbf{p}}[J_b|S, G]$ , has no active  $x$ -links in  $\alpha_1$ .

*Proof.* Any active  $x$ -links in  $\alpha_1$  remain so for all  $\{\alpha_i\}_{i=2}^M$ . The efficiency corresponding to each gain can be reduced in the following manner. Delete all such  $x$ -links and associated agents. For each mis-coordinating link between an  $x$  and  $y$  agent that existed, replace with a single link to a newly created isolated agent with an  $x$ -imposter attached. This preserves the stochastically stable states of all other nodes while reducing efficiency in each gain. ■

Intuitively, a graph that has coordinating  $x$  nodes in each gain  $\alpha_1, \dots, \alpha_M$  can be modified by removing these links, resulting in a lower efficiency. We are now ready to prove (22) (Theorem 3).

*Proof of (22) (Theorem 3).* Consider any graph  $G = (N, \mathcal{E}) \in \mathcal{G}_N$  and  $S \in \mathcal{T}(G)$ . Let us denote the  $M$  (worst-case) stochastically stable states that correspond to each gain  $\alpha_i$  with  $\hat{a}^i$ . Define for each  $k = 1, \dots, M$

$$P^k = \{n \in \mathcal{N} : \hat{a}_n^i = y \forall i \leq k, \hat{a}_n^i = x \forall i > k\} \quad (57)$$

as the set of nodes that play  $y$  in the SSS in  $\alpha_1, \dots, \alpha_k$  and play  $x$  in  $\alpha_{k+1}, \dots, \alpha_M$ . Note that  $P^k$  is possibly composed of multiple  $y$ -partitions. Also note it is possible that  $P^k = \emptyset$  for all  $k > \bar{m}$  for some  $\bar{m} \in \{2, \dots, M-1\}$  - that is,  $\hat{a}^i = \bar{x}$  for all  $i > \bar{m}$ . We first consider the case when  $P^k \neq \emptyset$  for every  $k = 1, \dots, M$ .

Let  $Q^k := \{n \in \mathcal{N} : \hat{a}_n^k = y\} = \bigcup_{i=k}^M P^i$ . Denote  $P^x := \{n \in \mathcal{N} : \hat{a}_n^1 = x\} = (Q^1)^c$  as the set of nodes stable to  $x$  for all  $\alpha_i$ . Consider the gain  $\alpha_i$  with  $i \leq k$ . Then the local efficiency  $\frac{W(e(P^k, \mathcal{N}), \hat{a}^k)}{W(e(P^k, \mathcal{N}), (a_{P^k}^*, \hat{a}_{-P^k}^*))}$  of  $P^k$  is

$$\frac{|e(P^k, Q^{k+1})| + |e(P^k, P^k)|}{(1+\alpha_{\text{sys}})(|e(P^k, P^k)| + |e(P^k, P^x)| + |e(P^k, (Q^k)^c)|)} > J_A^*(\alpha_i). \quad (58)$$

The inequality is due to Proposition 1. For gains  $\alpha_i$  with  $i > k$ , the local efficiency of  $P^k$  is

$$\frac{(1+\alpha_{\text{sys}})(|e(P^k, P^k)| + |e(P^k, P^x)| + |e(P^k, (Q^k)^c)|)}{(1+\alpha_{\text{sys}})(|e(P^k, P^k)| + |e(P^k, P^x)| + |e(P^k, (Q^k)^c)|)} = 1. \quad (59)$$

Hence, the overall system efficiency under gain  $\alpha_i$  is the mediant sum of the local efficiencies of the  $P^k$ . An application of Fact 2 gives the result. The case when  $P^k = \emptyset$  for  $k > \bar{m} \in \{2, \dots, M-1\}$  also follows directly from Fact 2. From the notation of Fact 2,  $\frac{n_k}{d_k} = 1$  for  $k > \bar{m}$ . ■

The details for the proof of (22) (Theorem 3) follow analogous arguments pertaining to focused attacks. Recall for a graph  $G \in \mathcal{G}_n$  and restricted action set  $\mathcal{A}$ , we denote  $F = F_x \cup F_y \subset \mathcal{N}$  as its set of fixed nodes. Additionally, we

restrict attention to gains  $\alpha_i \geq \alpha_{\text{sys}}$ , as these are not strictly dominated in the risk curve. The following structural property holds in a worst-case graph for focused risk.

**Lemma 8.** *A worst-case graph, i.e., a member of  $\arg \min_{G \in \mathcal{G}_N, F \in \mathcal{F}(G)} \mathbb{E}_{\alpha, \mathcal{P}}[J_{\text{f}}[F, G]]$ , has no active  $y$ -links in  $\alpha_M$ . Additionally,  $a_{F^c}^* = \vec{y}$ .*

*Proof.* A graph that has active  $y$ -links in  $\alpha_M$  remain active for all  $\alpha_1, \dots, \alpha_{M-1}$ . The efficiency corresponding to each gain can be reduced by removing all such links and keeping the border nodes as fixed  $y$  agents. This preserves the stability properties of all other nodes. The claim  $a_{F^c}^* = \vec{y}$  follows from Lemma 5. ■

We are now ready to prove (23) in Theorem 3.

*Proof of (23) (Theorem 3).* Consider any graph  $G = (\mathcal{N}, \mathcal{E}) \in \mathcal{G}_N$  and fixed nodes  $F \in \mathcal{F}(G)$ . The  $M$  stochastically stable states that correspond to each gain  $\alpha_i$  are denoted  $\hat{a}^i$ . Define for each  $k = 1, \dots, M$

$$P^k = \{n \in F^c : \hat{a}_n^i = x \ \forall i \geq k, \ \hat{a}_n^i = y \ \forall i < k\} \quad (60)$$

as the set of unfixed nodes that play  $x$  in the SSS for  $\alpha_k, \dots, \alpha_M$  and play  $y$  in  $\alpha_1, \dots, \alpha_{k-1}$ . Note that it is possible  $P^k = \emptyset$  for all  $k < \bar{m}$  for some  $\bar{m} \in \{2, \dots, M-1\}$ . That is,  $a_{F^c}^k = \vec{y}$  for  $k = 1, \dots, \bar{m} - 1$ . We first consider the case when  $P^k \neq \emptyset$  for every  $k = 1, \dots, M$ .

Let  $Q^k = \{n \in F^c : \hat{a}_n^k = y\} = \bigcup_{i=k}^M P^i$ . Consider the gain  $\alpha_i$  with  $i \geq k$ . Then the local efficiency  $\frac{W(e(P^k, \mathcal{N}), \hat{a}^k)}{W(e(P^k, \mathcal{N}), (a_{P^k}^*, \hat{a}_{-P^k}^k))}$  of  $P^k$  is

$$\frac{(1 + \alpha_{\text{sys}})(|e(P^k, P^k)| + |e(P^k, (Q^{k-1})^c)| + |e(P^k, F_x)|)}{|e(P^k, P^k)| + |e(P^k, Q^k)| + |e(P^k, F_y)|} > J_{\text{f}}^*(\alpha_i). \quad (61)$$

Here, we use the convention  $|e(P^1, (Q^0)^c)| = 0$ . For gains  $\alpha_i$  with  $i < k$ , the local efficiency of  $P^k$  is

$$\frac{|e(P^k, P^k)| + |e(P^k, Q^k)| + |e(P^k, F_y)|}{|e(P^k, P^k)| + |e(P^k, Q^k)| + |e(P^k, F_y)|} = 1. \quad (62)$$

Hence the overall system efficiency under  $\alpha_i$  is the median sum of the local efficiencies of the  $P^k$ . An application of Fact 3 gives the result. The case when  $P^k = \emptyset$  for  $k > \bar{m} \in \{2, \dots, M-1\}$  also follows directly from Fact 3. ■

## VII. SUMMARY

In this paper, we framed graphical coordination games as a distributed system subject to two types of adversarial influences. The focus of our study concerned the performance of a class of distributed algorithms against the associated worst-case risks. We identified fundamental tradeoffs between ensuring security against one type of risk and vulnerability to the other, and vice versa. Furthermore, our analysis shows randomized algorithmic designs significantly improves the available tradeoffs. Our work highlights the design challenges a system operator faces in maintaining the efficiency of networked, distributed systems.

*Proof of Lemma 3.:* This proof outlines a procedure to transform any  $(S, G)$  into a star graph with lower efficiency if  $\alpha < 1$ . We split into two cases - either  $(S, G)$  induces a single  $y$ -partition or more than one. First, assume  $(S, G)$  induces a single  $y$ -partition  $\mathcal{P}_y$ . An illustration of the constructive process is shown in Figure 4.

Construct a star subgraph  $\Gamma_y$  that has  $1 + |e(\mathcal{P}_y, \mathcal{P}_y)|$  nodes, each having a  $y$  imposter attached. Call the center node  $i_y$ . Construct similar star configurations  $\Gamma_x^k$  for each  $x$ -partition  $\mathcal{P}_x^k$ . Call their center nodes  $i_x^k$ . Connect  $\Gamma_y$  to each  $\Gamma_x^k$  with a link between  $i_x^k$  and  $i_y$ . If there are multiple edges between  $\mathcal{P}_x^k$  and  $\mathcal{P}_y$  ( $|e(\mathcal{P}_x^k, (\mathcal{P}_x^k)^c)| \geq 2$ ), create  $|e(\mathcal{P}_x^k, (\mathcal{P}_x^k)^c)| - 1$  new isolated nodes with a single  $x$  imposter attached, and connect each to  $i_y$  with a single link. At this point,  $\Gamma_y$  and  $\Gamma_x^k$  are stable  $y$  and  $x$ -partitions, and the isolated nodes are stable playing  $x$ . We have obtained a graph of  $N' \geq N$  nodes with identical efficiency to  $(S, G)$  since the number and type of edges are preserved.

We can further reduce efficiency if there are active  $x$  links, i.e. if  $|e(\Gamma_x^k, \Gamma_x^k)| \geq 1$  for at least one  $\Gamma_x^k$ . If there are none, then the graph belongs to  $\mathbb{S}_{N'}$  and we are done. Otherwise for each leaf node  $j \in \mathcal{P}_x^k$ , redirect the edge  $(j, i_x^k)$  to  $(j, i_y)$ , and replace  $j$ 's  $x$  imposter with a  $y$  imposter. Call  $m_x$  the total number of such converted nodes. The resulting graph-target pair  $(S', G')$  belongs to  $\mathbb{S}_{N'}$ . We claim the resulting (larger)  $y$ -partition  $\Gamma_y'$  is stable. For this claim to hold, (45) requires that

$$2|e(\mathcal{P}_y, \mathcal{P}_y)| + 2m_x + 1 \geq (1 + \alpha)(|e(\mathcal{P}_y, \mathcal{N})| + m_x). \quad (63)$$

From the original  $\mathcal{P}_y$ , it holds that

$$\begin{aligned} |e(\mathcal{P}_y, \mathcal{P}_y)| + |\mathcal{P}_y| &\geq (1 + \alpha)|e(\mathcal{P}_y, \mathcal{N})| \\ \Rightarrow 2|e(\mathcal{P}_y, \mathcal{P}_y)| + 2m_x + 1 &> (1 + \alpha)(|e(\mathcal{P}_y, \mathcal{N})| + m_x) \end{aligned} \quad (64)$$

due to  $|\mathcal{P}_y| \leq 1 + |e(\mathcal{P}_y, \mathcal{P}_y)|$  and  $\alpha < 1$ . All  $x$ -partitions in  $(S', G')$ , now just a collection of single nodes connected to  $i_y$  with an  $x$ -imposter, are stable. The efficiency is less than the original because active  $x$ -links increase efficiency more than active  $y$ -links do. Hence,

$$J_b(\alpha, S; G) > J_b(\alpha, S'; G'). \quad (65)$$

Now, we consider the remaining case when  $(S, G)$  induces  $k_y > 1$   $y$ -partitions  $\{\mathcal{P}_y^k\}_{k=1}^{k_y}$  and  $k_x \geq 1$   $x$ -partitions  $\{\mathcal{P}_x^k\}_{k=1}^{k_x}$ . Consider  $k_y$  such star subgraphs  $\{\Gamma_y^k\}_{k=1}^{k_y}$  with center nodes  $i_y^k$ . Recast the  $x$ -partitions into similar star subgraphs  $\{\Gamma_x^k\}_{k=1}^{k_x}$  with center nodes  $i_x^k$ . We first connect each  $\Gamma_x^k$  to some  $\Gamma_y^j$  with a single link  $(i_x^k, i_y^j)$  in any manner as long as a link between the original  $\mathcal{P}_x^k$  and  $\mathcal{P}_y^j$  exists. For each excess outgoing edge, we create an isolated node with an  $x$ -imposter attached. Each isolated node is attached to a corresponding  $i_y^k$  such that the original number of outgoing edges for each  $\mathcal{P}_y^k$  is satisfied. At this point, there are  $k' \leq k_y$  connected components  $G_k$  in the construction, and the efficiency of this construction is identical to the original. Lastly, we apply the

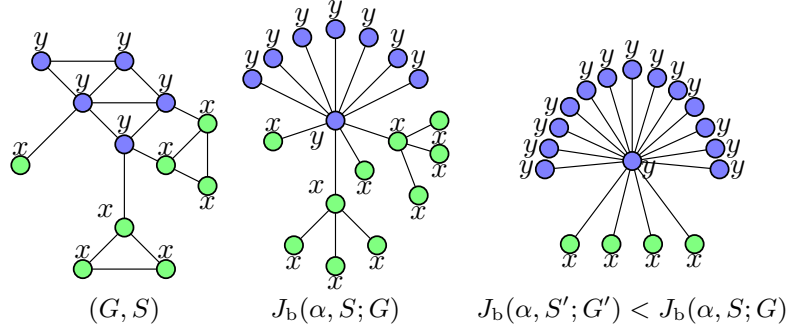


Fig. 4: An illustration of the constructive process (proof of Lemma 3) that generates a member  $(S', G') \in \mathbb{S}_m$  from any graph  $(S, G)$  with one  $y$ -partition, and  $\alpha < 1$ . Here, the labels on each node indicate the type of imposter influence. Green (blue) nodes play  $x$  ( $y$ ) in the SSS. (Left) Start with an arbitrary graph-adversary pair  $(S, G)$ . (Center) The partitions of  $(S, G)$  are re-cast as star subgraphs with the same number of edges. When there is more than one edge between the  $y$  and an  $x$ -partition, new nodes are created for the excess outgoing edges. This re-casting preserves the original efficiency  $J_b(\alpha, S; G)$ . (Right) The active  $x$ -links are converted into  $y$ -links by redirecting them to the center of the  $y$ -partition. This results in a graph  $(S', G') \in \mathbb{S}_m$ .

efficiency reduction procedure from before for each  $G_k$  to obtain  $(S'_k, G'_k) \in \mathbb{S}_{m_k}$ . From (37) and (34), we have

$$J_b(\alpha, S; G) > \frac{\sum_{k=1}^{k'} |e(\Gamma_y^k, \Gamma_x^k)|}{(1 + \alpha_{\text{sys}}) \sum_{k=1}^{k'} |e(\Gamma_y^k, \mathcal{N})|} \quad (66)$$

$$\geq \min_{k=1, \dots, k'} J_b(\alpha, S'_k; G'_k).$$

■

*Proof of Lemma 2: Technical result for expected risks.*

Let us define  $f(\mathbf{p}) := \sum_{i=1}^M p_i s_i$  and  $f_i(\mathbf{p}) := \left( \sum_{j=1}^k p_j \right) (\nu_k - 1) + 1$ . The set of probability vectors such that  $k = \arg \min_{i=1, \dots, M} f_i(\mathbf{p})$  can be written as the set

$$V_k := \left\{ \mathbf{p} \in \Delta_M : f_k(\mathbf{p}) \leq f_\ell(\mathbf{p}) \forall \ell \neq k \right\}$$

$$= \bigcap_{\ell \neq k} \left\{ \mathbf{p} \in \Delta_M : \frac{\sum_{j=1}^k p_j}{\sum_{j=1}^M p_j} \geq \frac{1 - \nu_k}{1 - \nu_\ell} \right\}. \quad (67)$$

Define  $\lambda_k = \frac{\sum_{j=1}^k d_j}{\sum_{j=1}^{k+1} d_j}$  for each  $k = 1, \dots, M-1$ . With some algebra, we can express each  $s_i$  as

$$s_i = \left[ \sum_{j=1}^{M-i+1} \frac{n_j}{d_j} (1 - \lambda_{j-1}) \left( \prod_{k=j}^{M-1} \lambda_k \right) \right] + \left( 1 - \prod_{j=M-i+1}^{M-1} \lambda_j \right) \quad (68)$$

Using the identities  $\sum_{k=1}^M (1 - \lambda_{k-1}) \left( \prod_{j=k}^{M-1} \lambda_j \right) = 1$  and  $\sum_{k=1}^\ell (1 - \lambda_{k-1}) \left( \prod_{j=k}^{M-1} \lambda_j \right) = \prod_{j=\ell}^{M-1} \lambda_j$  for  $\ell = 1, \dots, M-1$ , we obtain (omitting the algebraic steps)

$$f(\mathbf{p}) = \sum_{i=1}^M (1 - \lambda_{i-1}) \left( \prod_{j=i}^{M-1} \lambda_j \right) \left[ \left( \frac{n_i}{d_i} - 1 \right) \left( \sum_{j=1}^{M-i+1} p_j \right) + 1 \right] \quad (69)$$

Now, suppose  $\mathbf{p} \in V_k$  for  $k \in \{1, \dots, M\}$ . Using (67) and  $\nu_i \leq n_i/d_i \leq 1$ , we then have

$$f(\mathbf{p}) \geq \sum_{i=1}^M \lambda_{i-1} \left( \prod_{j=i}^{M-1} \lambda_j \right) f_k(\mathbf{p}) = f_k(\mathbf{p}). \quad (70)$$

■

## REFERENCES

- [1] J. Cortes, S. Martinez, and F. Bullo, "Robust rendezvous for mobile autonomous agents via proximity graphs in arbitrary dimensions," *IEEE Transactions on Automatic Control*, vol. 51, no. 8, pp. 1289–1298, Aug 2006.
- [2] M. Mesbahi and M. Egerstedt, *Graph theoretic methods in multiagent networks*. Princeton University Press, 2010, vol. 33.
- [3] I. Akyildiz, W. Su, Y. Sankarasubramaniam, and E. Cayirci, "Wireless sensor networks: a survey," *Computer Networks*, vol. 38, no. 4, pp. 393–422, 2002.
- [4] H. P. Young, *Individual strategy and social structure: An evolutionary theory of institutions*. Princeton University Press, 2001.
- [5] A. Montanari and A. Saberi, "The spread of innovations in social networks," *Proceedings of the National Academy of Sciences*, vol. 107, no. 47, pp. 20 196–20 201, 2010.
- [6] S. A. West, S. P. Diggle, A. Buckling, A. Gardner, and A. S. Griffin, "The social lives of microbes," *Annual Review of Ecology, Evolution, and Systematics*, vol. 38, no. 1, pp. 53–77, 2007.
- [7] M. Del Vicario, A. Bessi, F. Zollo, F. Petroni, A. Scala, G. Caldarelli, H. E. Stanley, and W. Quattrociocchi, "The spreading of misinformation online," *Proceedings of the National Academy of Sciences*, vol. 113, no. 3, pp. 554–559, 2016.
- [8] Y. Mao, S. Bouloki, and E. Akyol, "Spread of information with confirmation bias in cyber-social networks," *IEEE Transactions on Network Science and Engineering*, 2018.
- [9] S. Amin, A. A. Cárdenas, and S. S. Sastry, "Safe and secure networked control systems under denial-of-service attacks," in *Hybrid Systems: Computation and Control*, R. Majumdar and P. Tabuada, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2009, pp. 31–45.
- [10] S. Amin, X. Litrico, S. S. Sastry, and A. M. Bayen, "Stealthy deception attacks on water scada systems," in *Proceedings of the 13th ACM international conference on Hybrid systems: computation and control*. ACM, 2010, pp. 161–170.
- [11] H. Fawzi, P. Tabuada, and S. Diggavi, "Secure state-estimation for dynamical systems under active adversaries," in *2011 49th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*. IEEE, 2011, pp. 337–344.
- [12] F. Pasqualetti, F. Dorfler, and F. Bullo, "Attack detection and identification in cyber-physical systems," *IEEE Transactions on Automatic Control*, vol. 58, no. 11, pp. 2715–2729, Nov 2013.
- [13] S. Sundaram and B. Ghahserifard, "Distributed optimization under adversarial nodes," *IEEE Transactions on Automatic Control*, 2018.
- [14] V. Auletta, D. Ferraioli, F. Pasquale, and G. Persiano, "Metastability of logit dynamics for coordination games," in *Proceedings of the twenty-third annual ACM-SIAM symposium on Discrete algorithms*. Society for Industrial and Applied Mathematics, 2012, pp. 1006–1024.
- [15] L. E. Blume, "The statistical mechanics of best-response strategy revision," *Games and Economic Behavior*, vol. 11, no. 2, pp. 111–145, 1995.
- [16] J. R. Marden and J. S. Shamma, "Revisiting log-linear learning: Asynchrony, completeness and payoff-based implementation," *Games and Economic Behavior*, vol. 75, no. 2, pp. 788–808, 2012.

- [17] T. Tatarenko, "Proving convergence of log-linear learning in potential games," in *2014 American Control Conference*. IEEE, 2014, pp. 972–977.
- [18] D. Acemoglu, A. Ozdaglar, and A. ParandehGheibi, "Spread of (mis) information in social networks," *Games and Economic Behavior*, vol. 70, no. 2, pp. 194–227, 2010.
- [19] J. Ghaderi and R. Srikant, "Opinion dynamics in social networks with stubborn agents: Equilibrium and convergence rate," *Automatica*, vol. 50, no. 12, pp. 3209–3215, 2014.
- [20] H. P. Borowski and J. R. Marden, "Understanding the influence of adversaries in distributed systems," in *2015 54th IEEE Conference on Decision and Control (CDC)*, Dec 2015, pp. 2301–2306.
- [21] P. N. Brown, H. Borowski, and J. R. Marden, "Security against impersonation attacks in distributed systems," *IEEE Transactions on Control of Network Systems*, pp. 1–1, 2018.
- [22] B. Canty, P. N. Brown, M. Alizadeh, and J. R. Marden, "The impact of informed adversarial behavior in graphical coordination games," in *2018 IEEE Conference on Decision and Control (CDC)*, Dec 2018, pp. 1923–1928.
- [23] H. P. Young, "The evolution of conventions," *Econometrica: Journal of the Econometric Society*, pp. 57–84, 1993.
- [24] D. Foster and H. P. Young, "Stochastic evolutionary game dynamics," *Theoretical population biology*, vol. 38, no. 2, pp. 219–232, 1990.
- [25] D. Monderer and L. S. Shapley, "Potential games," *Games and economic behavior*, vol. 14, no. 1, pp. 124–143, 1996.