# Design of Dynamic and Personalized Deception: A Research Framework and New Insights for Cyberdefense

Cleotilde Gonzalez
Carnegie Mellon University
coty@cmu.edu

Palvi Aggarwal
Carnegie Mellon University
palvia@andrew.cmu.edu

Edward A. Cranford
Carnegie Mellon University
cranford@cmu.edu

Christian Lebiere
Carnegie Mellon University
cl@cmu.edu

## Abstract

*Deceptive defense techniques (e.g., intrusion detection, firewalls, honeypots, honeynets) are commonly used to prevent cyberattacks. However, most current defense techniques are generic and static, and are often learned and exploited by attackers. It is important to advance from static to dynamic forms of defense that can actively adapt a defense strategy according to the actions taken by individual attackers during an active attack. Our novel research approach relies on cognitive models and experimental games: Cognitive models aim at replicating an attacker's behavior allowing the creation of personalized, dynamic deceptive defense strategies; experimental games help study human actions, calibrate cognitive models, and validate deceptive strategies. In this paper we offer the following contributions: (i) a general research framework for the design of dynamic, adaptive and personalized deception strategies for cyberdefense; (ii) a summary of major insights from experiments and cognitive models developed for security games of increased complexity; and (iii) a taxonomy of potential deception strategies derived from our research program so far.*

## 1. Introduction

Cyberattacks fundamentally occur by taking advantage of the power of deception—the act of intentionally inducing and suppressing signals to cause behavioral changes in a target to benefit the deceiver [1]. Through the years, attackers have perfected their use of deception by taking advantage of the ease with which it is possible to conceal their identities, their actions, and their intentions in cyberspace. Furthermore, attackers have become experts in their use of social engineering—psychological manipulations to trick people into disclosing sensitive information or unlawfully granting access to a secure system [2]. However, deception has also been used as a defense strategy benefiting defenders

and attackers alike [3]. For example, common deception strategies used in cyberdefense include: masking and decoying [4, 5]. Masking is a technique of making a real object undetectable; used to hide information behind benign programs (e.g., hiding information behind an image in an email message) while decoying presents a false object to grab attention by showcasing fake but relevant information (e.g., honeypots may grab an attacker's attention by showcasing essential data that is of value to the attacker). Mechanisms such as honeypots—a fake system that is valuable by being attacked—have been used extensively to secure information, stop and detect spam, and enhance network defense [6]. Honeypots are used for detection to catch illicit interactions; in prevention, to assist in slowing attackers down; and many other defense possibilities [7]. However, the effectiveness of honeypot techniques is questionable, as they often rely on static allocations that can often be easily discovered by attackers.

A successful approach to address the optimal allocation of defense resources relies on optimization solutions that use game theoretic models and in particular *Stackelberg Security Games* (SSGs) [8, 9]. As we will discuss below, these solutions have proven to be effective in a large number of physical security cases; and although they have not been practically applied to cybersecurity, work is underway to address this gap [8].

Our research program on cyberdeception contributes to addressing this gap by investigating the use of deceptive signals and their interaction with the allocation of defense resources through SSGs in the context of cybersecurity. Finding the right balance of deceptive signals so that the attacker continues to believe the signal is crucial to the success of deceptive strategies [10]. Our research program aims at advancing our understanding of how deceptive signals can be designed and presented to attackers in order to maximize their effectiveness.

As we discuss in the following sections, fundamental to the success of a dynamic and personalized deception strategy is the integration of computational

representations of human behavior with the optimization solutions for SSGs. Designing effective defense algorithms must make use of the knowledge of human behavior: the way humans make decisions, how they explore an environment, how they take risks, and how they use their experience. Insights on human behavior often emerge from laboratory experiments where we can study would-be attackers. This behavior is then captured by theoretically grounded cognitive models that explain the dynamics of human behavior in computational forms to predict the actions of a human attacker. Cognitive models are important given that human behavior is often far different from what is predicted under the assumptions of perfect rationality models [11, 12]. Our recent findings indicate that signaling algorithms, while optimized for perfectly rational adversaries do improve defense compared to not signaling at all, are less effective than expected for boundedly rational humans [10, 13]. Therefore, the study of human behavior in experiments and the computational representation of their decision process is essential to advance current deception strategies of defense.

## 1.1. Deception as a Defense Strategy: Use of Game Theory and Signaling in Cybersecurity

Researchers of the broad field of security have addressed the important question of how to assign limited defense resources to potential targets, by using SSGs and game theoretic optimal solutions. SSGs have two players, the defender and the attacker; a defender must defend a set of targets using a limited number of resources, whereas the attacker is able to observe and learn the defenders strategy and attack after surveillance. A defender commits to a mixed strategy (e.g., how to allocate defense resources), and then the attacker conducts surveillance of these mixed strategies and responds with an attack on a target that optimizes her reward [8]. The objective in this line of research is to find an optimal mixed strategy for the defender (i.e., an optimal allocation of defense resources), called the *Strong Stackelberg Equilibrium* (SSE) [8, 9].

This work has led to an impressive set of practical cases of the theoretical insights from the SSE [8]. Cases include: the allocation of canine patrols and vehicle checkpoints at LAX [14], as well as the allocation of rangers over National Parks to protect wildlife from poachers [15].

Although the insights of SSGs research have not been practically applied to Cybersecurity this research is underway [16, 17]. Researchers have aimed at advancing the insights from physical

security to cybersecurity [18] by creating *interactive* "security games" through which strategies of the defender can be paired against humans acting as attackers. These human-in-the-loop security games are abstract representations of the essential elements of cybersecurity, often used in laboratory experiments in order to understand human attackers' behavior against particular defense strategies [15, 19, 20].

Interactive security games are helpful to investigate cyberdeception given the limitations of the theoretical insights of SSGs to practical cybersecurity applications. In contrast to physical security, cybersecurity is a complex context that challenges current SSG research and the SSE strategy computations [8]. Cybersecurity is a more dynamic and complex context compared to any physical security problem. A computer network, and the targets that a defender must protect, can change dynamically. This is a problem with current SSE calculations, as an optimal strategy must be recomputed in real-time. Furthermore, the common assumption in SSE that adversaries know and are able to survey among the defense mixed strategies is unreasonable in cybersecurity. In the cyberworld, the adversaries know little about the defense strategies and similarly the defender has little information about the possible attackers' actions.

Given the importance of deception and the multiple challenges in applying SSGs directly to naturalistic cybersecurity problems, researchers have focused on the investigation of *signaling* strategies [11, 21] with the use of interactive simulations in human experiments. In order to test the effectiveness of these algorithms, researchers compare their results against human attackers' actions in laboratory experiments [11]. In this new line of research, signaling is investigated in conjunction with the SSGs and SSE. The strategic exploitation of information by the defender can influence and deceive the adversary, and this is formalized by incorporating a signaling game model into the SSGs, where the defender strategically reveals information about a defensive strategy to the attacker, in order to influence the attackers decision making [16]. Signaling appears to improve the defender utility against a perfectly rational attacker compared to the traditional SSE model without signaling, but more research is needed to understand the use of deceptive signaling as a proactive mechanism of defense.

Various forms of signaling have been proposed to increase a human attacker's compliance with signals [17], but much remains to be investigated regarding how humans process and act according to such signals. Our recent research shows that humans behave far differently than predicted under
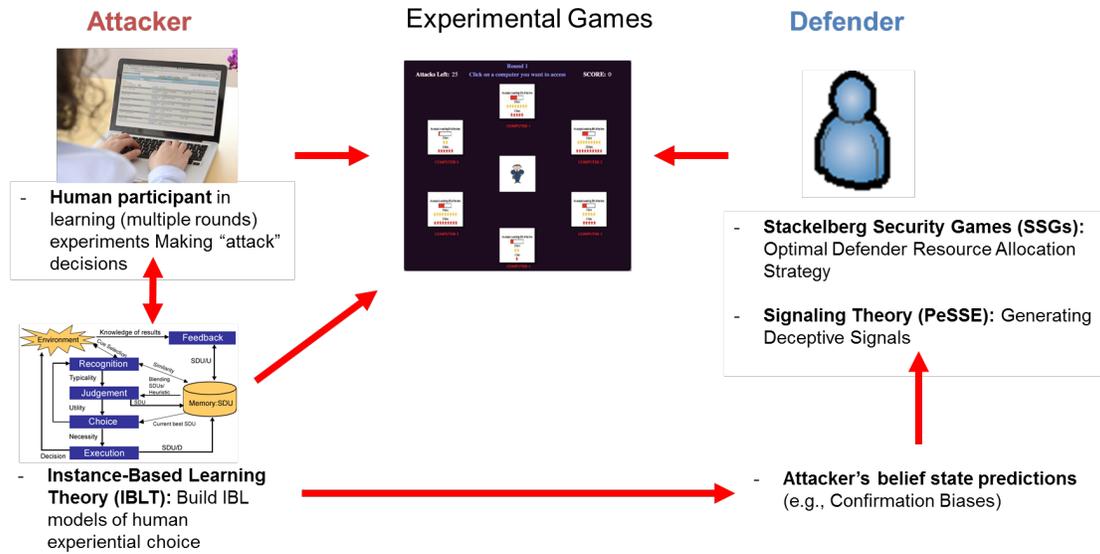
**Figure 1. A Research Framework for Dynamic, Adaptive, and Personalized Defense Strategies**

the assumption of perfect rationality [17]; humans exhibit nominally irrational behaviors that result in cognitive biases (e.g., confirmation bias) [12, 13]; and new adaptive and personalized theories that increase attacker's compliance are possible through cognitive modeling and human-in-the-loop experiments [10]. Our research framework proposes how experimental games of increased complexity and cognitive models that represent human behavior contribute to the generation of adaptive and personalized deception strategies.

## 2. A Research Framework for the Design of Dynamic, Adaptive, and Personalized Deception

The goal of our research program is to provide personalized, dynamic and adaptive deception algorithms for effective and agile defense capabilities. Our approach shown in Figure 1 uses an innovative combination of SSG algorithms for distribution of limited defense resources, optimization methods from game-theory and signaling theory (e.g., SSE), experimentation with human-in-the-loop interactive security games, and adaptive cognitive modeling using Instance-Based Learning Theory [22].

A defender allocates resources in an interactive computer game according to a policy defined by the optimization algorithms (e.g., SSE). Then, a human attacker makes decisions over multiple rounds about surveillance and attack of available targets. A signaling scheme is defined either independently or dependently over the allocation algorithm to optimize the frequency

of deceptive and truthful signals sent to the attacker. A deceptive signal provides information that modifies the ground truth. The attacker observes the signal and then decides whether to proceed or not with the attack.

We contribute to the line of SSG research, by (1) providing insights from human experiments regarding human trust to signals that are deceptive and truthful and (2) creating cognitive models that represent the decisions made by would-be human attackers that can inform the algorithms for allocation of defense resources.

### 2.1. Experimental Games: Scaling up Complexity

To demonstrate applicable dynamic and personalized deception strategies, we rely on interactive security games to investigate basic principles of deception across increasing levels of complexity and realism of the games. The interactive security games presented in Figure 2 are representative of key dimensions that guide deception against human attackers (see section 3). Through the analyses of human actions across these games, we identify cognitive biases that can be exploited in deception operations and how those biases can be modeled computationally. In what follows we summarize the major insights from current experimental work in each of these 4 exemplar games.

**2.1.1. Box Game.** The Box Game (Figure 2a) is a simple, 2-stage, 2-alternative SSG. In stage 1, a defender allocates resources to one of two boxes with 0.5 probability according to the optimal resource
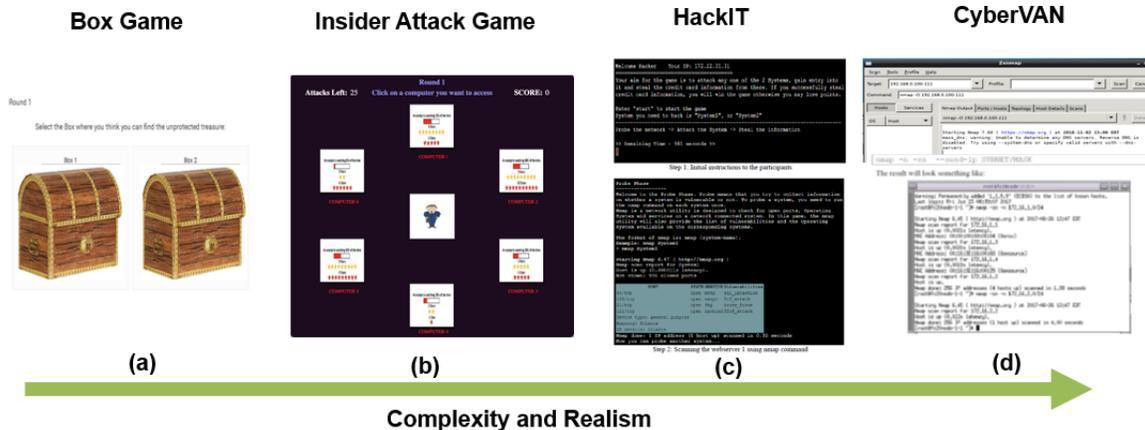
**Figure 2. Different levels of sophistication in interactive security games**

allocation SSE algorithm (i.e., pure strategy) [16]. In stage 2, the defender sends a signal to influence the attackers decision making to her benefit by exploiting the fact that the attacker is unaware of the pure strategy at any given time.

Each round of the two-stage game formalized in [10] plays out as follows: (1) The defender allocates her resources (i.e., a "treasure" of 100 points) covering a random subset of the targets based on her mixed strategy (probability of 0.5); (2) The attacker chooses a target, t, to attack accordingly; (3) The defender sends a (possibly deceptive) signal to the attacker regarding the current protection status of t; (4) Based on the information given in the signal, the attacker chooses to either continue attacking or withdraw his attack. If the attacker attacks and the node is protected, the attacker looses 50 points, if the node is not protected, the attacker gains 100 points. If the attacker withdraws the attack, it yields 0 points.

Our goal is to find out how humans acting as "attackers" (i.e., treasure hunters) behave under various frequencies of deceptive signals. In other words, the question is how often should a defender send a deceptive signal in order to gain the most benefit? A recent (unpublished) study manipulated the frequency of sending a truthful signal when the node is protected (1, 0.75, 0.5) and the type of signal sent (positive or negative frame). This probability also controls the proportion of truthful or deceptive signals for unprotected nodes according to the procedure formalized in [10]. We found that, overall, the signaling scheme has no effect on attacks to unprotected nodes; however, as the frequency of a truthful signal for protected nodes decreased, the proportion of attacks on protected nodes increased. Signaling increased the attackers losses in protected nodes, while there was no significant increase in the attackers gains from attacks to

unprotected nodes. We also found that deceptive signals generate less attacks than truthful signals. Furthermore, positive (i.e., hopeful) signals produce more attacks when they are truthful than deceptive. The negative (i.e., "suspicious") signals cause more deterrence compare to positive signals.

**2.1.2. Insider Attack Game.** The Insider Attack Game (IAG) (Figure 2b) is an escalation of the box game as it increases the number of nodes to six, and adds more contextual information on the nodes (e.g., gains and losses, and probabilities of coverage) [11]. The allocation of defense resources is more complicated as only two out of the six nodes can be protected at a time. The attacker has access to the abstract information about the node i.e., value of each node, losses if the node is protected and probability information about each node being protected. As an example, the IAG has been used to test the signaling strategy for deception in a cybersecurity scenario where the participants play a role of an attacker (a company's employee) who tries to attack the computers to gain points. The company has six computers and only two security defenders to monitor these computers. The defenders could only protect two computers at a time. To secure more computers, we use signaling to send warnings to deceive attackers [10]. In each trial, the player analyzes the information on each node and selects a computer to attack. They may receive a signal from the defender and then decide whether to proceed with the attack or withdraw it. The signaling algorithm used determines whether to send a truthful signal or a deceptive signal. [10] suggested that instead of using an exact proportion of deceptive or truthful signals, "Goldilocks" zones works better: the algorithm that sends least proportion of signals achieves the lower proportion of attacks. However, the proportion of attacks increase to almost

95% when there is no signal. This result suggests, as cognitive science would predict, that human behavior is dynamic and adaptable, and having one static proportion of signaling will ultimately be ineffective if it does not adapt to the attacker actions. Using cognitive models that emulate the attacker's actions, we have developed adaptive and personalized signaling schemes that learn about attacker's actions and adjust signals accordingly [12]. An adaptive signaling scheme starts by sending truthful signals to the attacker (i.e., to gain trust), then, according to the attacker's level of trust, the algorithm adapts to whether the next signal should be a truthful or deceptive. We demonstrate how this scheme reduces the probability of attack, although at the expense of giving up more attacks in the first few trials [13]. Current work is ongoing to advance on these insights.

**2.1.3. HackIT Game.** HackIT (Figure 2c) is a generic web-based framework for cybersecurity to study human learning and decision-making of attackers and defenders [23]. HackIT advances the IAG by including more semantic information such as network nodes, representing the characteristics of real nodes; deception tactics: masking, decoying; and commands, which are used for communication with the network. The defender protects the real nodes using deception tactics and the attacker's goal is to identify the real network nodes and exploit them. In HackIT, the attacker gathers information (pull information) such as operating systems, open and closed ports, services on the network nodes, and vulnerabilities from the network using probing action. Attackers could communicate with the network in HackIT using tools such as nmap and gain information about network nodes, topologies and configurations. However, attackers are not aware of the strategies used by defenders and they must learn those strategies overtime by playing different rounds.

HackIT has the potential to simulate many real-world dynamic situations in the laboratory: manipulating deception tactic (e.g., decoying and masking); frequency of deceptive signals (e.g., using different proportion of honeypots in the network or testing optimal placement of honeypots); and manipulating the content of the signal (e.g. use different configurations of honeypots).

A recent (unpublished) study, builds on insights from Achleitner et al. (2016). They simulated the following reconnaissance strategies in deceptive and non-deceptive networks: Uniform Scanning, Local Preference Scanning, Preference Sequential Scanning, Non-Preference Sequential Scanning and Preference Parallel Scanning. Using HackIT, we are investigating whether humans exhibit behavior described by these strategies while probing. In order to do this, we first

simulated a network of 40 computers in two network topologies: RDS (Reconnaissance deceptive server) and Non-RDS. Out of 40 computers, only 25% of the hosts were real while the remaining hosts were honeypots. In a RDS topology only a subset of nodes are visible to the attackers and other computers are only visible after they exploit any computer at the first layer. This provides a deceptive view of the network which would be different from each node they exploit. In a Non-RDS configuration all nodes are connected to a single central host. The hosts in the Non-RDS configuration can easily be exploited if the central host in the network is compromised. The goal of a simulated attacker is to maximize the number of real systems exploited during a fixed time. Our current results suggests that the real host detection rate in the RDS topology was lower than the on in the Non-RDS topology. Preliminary results of an on-going experimental study including 18 participants (6 in each condition), suggest a humans follow a uniform scanning strategy in both conditions. The number of real systems exploited were higher in the non-RDS condition than in the RDS, suggesting that hiding real computers in a layered network configuration (RDS) reduces the probability of exploitation.

**2.1.4. CyberVAN.** CyberVAN (Figure 2d) is a security testbed build on top of Virtual Ad hoc Network (VAN) for cybersecurity research [24]. CyberVAN is capable of speedy creation of high-fidelity strategic and tactical network scenarios using virtual machines, simulated networks, physical nodes and physical networks. These scenarios could be controlled by either GUI or commands on a console. CyberVAN is capable of generating realistic cyber experimentation environments which includes the simulated cyberattacks, cyberdefense, providing synthetic users for creating realistic network traffic and creating human-in-the-loop environments for validating various defense algorithms. Specifically, for cyberdeception experiments, CyberVAN can provide different deception tactics such as masking (by hiding/faking the configuration of nodes) and decoying (by using honeypots, honeynets, honeytokens etc.). The information manipulated for creating deception includes network structure, number of nodes in the network, operating system, ports, services, vulnerabilities, network round trip time, network traffic etc. The proportion of deception could be controlled using different defense algorithms which could be integrated in CyberVAN. Attackers could interact with virtual machines using various network scanning tools (e.g. nmap) to gather information during a probing phase.

As an example, we are using CyberVAN in an ongoing experiment where we test the effectiveness of

an optimal masking strategy against random masking strategy in reducing the utility of attackers. True configurations (TC) of virtual machines are masked to a fake observable configuration (OC). The configuration of virtual machines are different in each round and participants are able to attack only one machine in every round. This scenario is presented to human participants who probe the machines to gather information (e.g., operating systems, ports, services). Once attackers explore the network and learn the features of network (some of them were masked), they decide what machine to attack and what type of exploit to use.

An ongoing experiment consists of 6 rounds (1 practice round and 5 actual rounds) and each round consists of 12 virtual machines. We expect that attackers would likely choose medium utility nodes for attack. Also, given insights from the other interactive security games, we expect that attackers may be naive or more experienced. A naive attacker would attack based on fixed preference of OCs, but an experienced attacker would utilize the information about probability distribution of TCs being mapped to OCs. Overall, we expect that the optimal masking strategy would be helpful to improve defenders utility compared to the random masking strategy.

## 2.2. Cognitive Models of Attacker Behavior

In addressing unrealistic assumptions of attacker's rationality in SSE, researchers have used models from behavioral economics; including Quantal response and Subjective Utility Quantal response [25]. Although these models are common and well-known for their statistical properties that reflect human choices, these are not process models: they cannot explain the cognitive mechanisms by which humans make decisions, presenting many limitations to design dynamic, adaptive technologies that support cyberdefense processes, such as those involved in the new signaling schemes [26].

Cognitive models are dynamic and adaptable computational representations of the cognitive structures and mechanisms involved in cognitive tasks such as processing information for decision making. Cognitive modeling technologies rely on attention, memory, and decision making theories, that allow for the construction of generative models to be eventually tested against behavioral, physiological, or neural data. The advantage of cognitive models resides in their ability to dynamically learn from experience, to adjust to new inputs, environments, and tasks in similar ways as humans do, and to predict performance in situations that haven't been encountered

and for which data is not yet available [26, 22]. In this regard, cognitive models differ from purely statistical approaches, such as machine learning, that are often capable of evaluating only stable, long-term sequential dependencies from existing data but fail to account for the dynamics of human cognition and human adaptation to novel situations.

We developed cognitive models to replicate the attacker's actions in the presence or absence of signals in three SSGs: the Box Game, the Insider Attack Game, and the HackIT game. We have not yet developed cognitive models for CyberVan, but we expect that the same general approach that we have followed for the other SSGs will apply to our current work in CyberVan. Our cognitive models aim at turning static deception strategies into dynamic and personalized deception. Our cognitive models of an individual opponent rely on the theoretical principles defined in IBLT [22], which state that humans make decisions from experience according to the similarity, frequency and recency of experienced events and the value of the actions taken (i.e., their utility).

Cognitive models and the major insights in the SSGs have been reported in various recent publications [11, 12, 13]. From the results of the cognitive models compared to human behavior in the SSGs, we have learned that: (1) humans behave far differently than predicted under the assumption of perfect rationality [17, 11]; (2) humans exhibit nominally irrational behaviors (e.g., confirmation bias) that reflect capacity and information limitations and the need to resort to heuristic strategies; and that (3) while signaling algorithms optimized for perfectly rational adversaries do improve defense compared to not signaling at all [16], they are less effective than expected for boundedly rational humans [10].

In our current work [13] we improved upon traditional game-theoretic signaling schemes (the peSSE signaling scheme [10] by developing a cognitive signaling scheme that is adaptive and based on cognitive principles. The cognitive model predicts human decisions are made by aggregated retrieval across past experiences based on the similarity to the current situation as predicted by IBLT [22]. From this recent work we observe that: (1) human decisions are strongly influenced by confirmation bias, and (2) it is important to consider the dynamics of signal in the individual decisions: Continued attacks given truthful signals strengthen the expectation of a loss given a signal. This and other a potential strategies of defense need to be investigated further, to determine how to maintain the compliance to the signals to the benefit of the defender.

In general, our current work has outlined an initial

approach to deceptive signaling for cyberdefense that relies on cognitive models of attacker behavior to balance the rate of deception in an attempt to keep the attacker's belief in the signal high. Importantly, our cognitive signaling scheme is adaptive and personalized, and can therefore be used to induce biases and influence attackers to comply with the signal beyond the capabilities of any static scheme.

## 3. Key Dimensions for the Design of Dynamic and Personalized Cyberdefense

A recent review of the literature regarding deception in cybersecurity proposes a taxonomy of the types of deception that correspond to the game-theoretic notions of private information, actors, actions, and duration [27]. The authors use these game-theoretic notions to describe a set of 6 types of deception: perturbation, moving target defense, obfuscation, mixing, honey-x, and attacker engagement. Perturbation refers the application of noise in the information itself. Moving target defense refers to the idea of changing attack surfaces and creating random configurations; for example, by using mixed strategies. Obfuscation, refers to hiding valuable information using external noise; for example, deceptive routing of traffic. Mixing, relates to a technique of hiding valuable information in an attempt to make the entry and exit nodes unlinkable. Honey-x, refers to deception which uses common techniques such as honeypot, honeynet, honeybot, etc. Finally, attacker engagement refers to multi-period, dynamic games in which deception techniques must be adapted to the actions of the attacker.

In this section, we concentrate on the attacker engagement dimension proposed by [27]. Specifically, we propose a set dimensions for the design of dynamic and personalized deception derived from our own results from experimentation and modeling connected to the interactive SSGs reviewed above. We present three major dimensions of deception defenses in attacker engagement in dynamic SSGs: (1) Deception tactic; (2) Signaling strategy; (3) Interaction mechanisms.

**Deception Tactics.** Traditional security defense tactics are often static and reactive: a defender monitors network traffic and uses technology for intrusion detection that supports the detection of cyber-attacks. A recent data breach investigation report suggest this form of defense is highly unsuccessful, as only a low percentage of breaches are detected [28]. Deception-based tactics of defense can provide many advantages over traditional methods [7, 3]: they can induce attackers to take actions that benefit defenders.

In our current work with SSGs, we demonstrate the use of two deception tactics in over multiple periods in cybersecurity games: Masking and Decoying. Table 1 summarizes the way these two deception tactics are being used in our four examples of SSGs. Masking has been used to hide the facts about the reality (e.g., A defender can mask vulnerabilities to showcase that the computer is secure). For example, before even implementing the real defense, a defender could mask the *Server Message Block* service version to showcase it is patched and secure the network from WannaCry ransomware attack. We also use mimicking where software and services may imitate the ground truth [21]. The intention is to showcase a system more or less valuable to the attacker. For example, a system may respond as if it is running a version of Windows XP while actually running Windows 7 [4].

Table 1: Deception Tactics in different interactive cybersecurity games

| Security Game | Deception Tactics | |
|---|---|---|
| | Masking | Decoy |
| Box Game | Information about the protected node is hidden from the attacker. | Some nodes in the network could be made more attractive than others by using appropriate combination for rewards, penalty and probability or being protected. |
| Insider Attack Game | Masking involves hiding information such as: probability of node being protected, rewards and penalities for attacking the node. | Some nodes in the network could be made more attractive than others by using appropriate combination for rewards, penalty and probability or being protected. |
| HackIT | Features such as operating system, service version, vulnerability could be masked with fake features. | Systems could be configured as honeypots to lure attackers with limited features. |
| CyberVAN | Features such as operating system, service version, network traffic, response time, vulnerability could be masked with fake features. | Honeypots, honeytokens, honeyfiles could be used for decoying strategies. |

A decoy tactic is another popular concept used by defenders to identify attackers, gather information about their techniques and for securing the real network by luring hackers into the honeypots [29]. Honeypots, honeynets and honeytokens are examples of decoy deception tactic. A honeytoken could be a record in the database which is not relevant to any database user. Any access to such record would suggest a malicious user and defender could investigate this user. Honeytokens are a honey thing in a normal or honeypot system.

**Signaling Strategy.** Signaling theory addresses a fundamental problem in the communication between a sender (the signaler) and a receiver: whether the sender's *message* is conveying the truth or manipulating the information to her benefit [30]. Signaling has been used in SSGs in a way that it is incentive-compatible for a sender to transmit a message that partially reveals her private information, since the receiver cannot know the underlying information with certainty [27].

Table 2: Signaling Strategy in different interactive cybersecurity games

| Security Game | Signaling Strategy | | |
| --- | --- | --- | --- |
| | Frequency of Deceptive Signals | Level of Information | Type and Content of Signal |
| Box Game | Proportion of deceptive signals could be controlled by creating balance between deceptive signals and truthful signals. | Abstract messsage: Defender pushed information to the attacker by sending a signal. | Messages that provide information about defender protecting a node or not. |
| Insider Attack Game | Proportion of deceptive signals could be controlled by creating balance between deceptive signals and truthful signals. | Abstract messsage: Defender pushed information to the attacker by sending a signal. | Messages that provide information about defender protecting a node or not. |
| HackIT | Proportion of deceptive signals could be controlled by proportion of honeypots in the network. | Contextual Information: Attacker pull information by probing the network. | Contextual Signal: Information about network structure, number of nodes, operating system, ports, services, and unpatched |
| CyberVAN | Proportion of deceptive signals could be controlled by proportion of honeypots in the network, and number of computers being masked in the network. | Contextual Information: Attacker pull information by probing the network. | Contextual signal: Information about network structure, number of nodes in the network, operating system, ports, services, vulnerabilities, network round trip time, and network traffic. |

In the context of cybersecurity, attackers may gather information from scanning nodes in the network (i.e., "pull information"); but also, defenders may strategically use signals to provide deceptive information to the attacker (i.e., "push information"). The goal is to waste attackers' resources and time and lead them to "the light," revealing their intentions and identity. As discussed above, we have investigated signaling strategies in SSGs, where these strategies essentially identify a proportion of times in which a deceptive signal could be sent from the defender to the attacker (e.g., how often to say that an unprotected node is protected or say that a protected node is unprotected). Research regarding signaling is a promising current area of research in SSGs [17, 16]. However, a methodical cognitive approach is required to experimentally investigate the dimensions important for signaling in SSGs.

Deceptive warning messages or explicit information such as network structure, number of nodes in the network, operating systems, ports, services, network traffic, round trip time information, and unpatched vulnerabilities in the network could all be used by the defender to deceive the attacker. Here, we consider three relevant dimensions for the investigation of signaling in SSGs: (1) The frequency of deceptive signals, (2) the level of information revealed to the attacker, and (3) the type of signal and content of the signal. Table 2 summarizes the way these three signaling strategies are being used in our four examples of SSGs.

The frequency of deceptive signals is a common theme of current research in SSGs, however current algorithms optimizing the signal frequency are less successful than expected [17, 16]. The reason is that humans are not rational, they learn from their experience, and they adapt accordingly. For example, if the defender deceives too frequently, the attacker will get to learn this tendency making the defense strategy astray. Generally, any non-adaptive algorithm of defense will tend to be ineffective against human attackers, and SSG researchers believe that there is a *Goldilocks Zone*, an optimal level of deception that could be more effective to improve the attacker's compliance in cybersecurity games [10]. In our current work our cognitive models of attacker's behavior helps to inform the development of adaptive and personalized signaling algorithms in cyberdefense [12, 13], although the effectiveness of all these algorithms needs to be tested empirically.

Addressing different forms of uncertainty is one of the major current challenges in SSG research [18]. The amount of information provided in a signal can influence the effectiveness of the deception and ultimately the actions an attacker can take. A common assumption in SSG research is that an attacker and a defender have perfect information about the state of the world: the payoff matrix and the opponents' strategies (e.g. probabilities for choosing available options) [18]. However, in naturalistic tasks this is an unreasonable assumption. Past research of behavioral game theory has addressed the effects of information in traditional social dilemmas [31]; but research should systematically experiment with uncertainty levels in SSGs in the context of cybersecurity: revealing more or less information in the signals sent to the attacker to determine their effectiveness on deterring attacks.

In addressing the effectiveness of deceptive signals, the content of the signals must also influence how the attacker reacts to the information. For example, the polarity of a message (i.e., presented as positive

Table 3: Interaction Mechanism in different interactive cybersecurity games

| Security Game | Interaction Mechanism | | |
| --- | --- | --- | --- |
| | Probing | Adaptive | Personalized |
| Box Game | Current version assume that attacker already probed and gathered all the information. | The signaling algorithms could easily adapt to the attackers' actions using IBL models of the attacker. | The signaling algorithms used in this game could easily create personalized signals using IBL models of the attacker. |
| Insider Attack Game | Current version assume that attacker already probed and gathered all the information. | The signaling algorithms could easily adapt to the attackers' actions using IBL models of the attacker. | The signaling algorithms used in this game could easily create personalized signals using IBL models of the attacker. |
| HackIT | Attackers probe the network using nmap command and gather information about network nodes. | The signaling algorithms could learn the attackers actions and adapt the configuration. However, adapting the configuration could increase the complexity of the task. | The deception strategy personalized could create different configurations for different attackers. However, the number of configurations available would be limited. |
| CyberVAN | Attackers probe the network using nmap command and gather information about network nodes. | CyberVAN use ACyDS, an adaptive cyber deception system to learn the attackers actions and generate the adptive network configuration on fly. | Making the deception strategy personalized to different attackers can be very hard and costly in CyberVAN. |

or negative) may greatly influence how attackers react to the signal. Consider for example, the well-known human bias called the *framing effect* [32]. Regardless of choice options being of equal expected values, the frame will elicit systematically different choices when presented as gains or losses. This robust human bias can be used in the design of content of signals that could deter or encourage the attacker's actions. As alluded to earlier we have observed such asymmetric effects in a simple binary choice task (e.g., the "box game").

**Interaction Mechanisms.** An attacker interacts with a network to gather information about the network structure, number of nodes in the network, their configuration, protocols and unpatched vulnerabilities, by passively or actively probing the network [7]. Active and passive probing leaves information about attackers in the network which could be used by defenders to learn about attackers and improve their defense based on the attackers activities. Table 3 summarizes the way in which interaction mechanisms are being used in our four examples of SSGs; the interaction of attacker and defender could happen in three ways:(1) probing, where defender already used deception tactics and the attacker gather information about the network before attack; (2) adaptive, as attackers communicate with the network through his probing and attack actions, defender adapts the deception or the resource allocation strategy [33]; (3) personalized, a powerful defender may generate personalized signals or network configuration based on attackers activities in the network.

## 4. Conclusion

Our research program on cyberdeception contributes to elucidating ways to use deceptive signals in SSGs in the context of dynamic and personalized cyberdefense. We contribute to SSG research program, by providing insights from human experiments regarding human trust to truthful or deceptive signals, and creating cognitive models that represent the decisions made by would-be human attackers that can inform the algorithms for allocation of defense resources.

Across four levels of complexity in interactive security games and using the insights of cognitive models of attacker behavior, we find that: (1) signaling algorithms optimized for perfectly rational attackers improve defense compared to no signaling at all; (2) humans behave far differently than predicted under the assumption of perfect rationality [17]; (3) humans exhibit nominally irrational behaviors that result in cognitive biases (e.g., confirmation bias) [12, 13]; and (4) new adaptive and personalized theories that increase attacker's compliance are possible through cognitive modeling and human-in-the-loop experiments [10].

## 5. Acknowledgments

## References

[1] N. C. Rowe and J. Rrushi, *Introduction to Cyberdeception*. Springer, 2016.

[2] P. Rajivan and C. Gonzalez, "Creative persuasion: A study on adversarial behaviors and strategies in phishing attacks," *Frontiers in psychology*, vol. 9, p. 135, 2018.

[3] K. J. Ferguson-Walter, D. S. LaFon, and T. Shade, "Friend or faux: deception for cyber defense," *Journal of Information Warfare*, vol. 16, no. 2, pp. 28–42, 2017.

[4] M. A. McQueen and W. F. Boyer, "Deception used for cyber defense of control systems," in *Proceedings of the 2nd conference on Human System Interactions, HSI*, vol. 9, 2009.

[5] B. Whaley, "Toward a general theory of deception," *The Journal of Strategic Studies*, vol. 5, no. 1, pp. 178–192, 1982.

[6] L. Spitzner, *Honeypots: tracking hackers*, vol. 1. Addison-Wesley Reading, 2003.

[7] M. H. Almeshekah and E. H. Spafford, "Cyber security deception," in *Cyber deception*, pp. 23–50, Springer, 2016.

[8] A. Sinha, F. Fang, B. An, C. Kiekintveld, and M. Tambe, "Stackelberg security games: Looking beyond a decade of success.," in *IJCAI*, pp. 5494–5501, 2018.

[9] M. Tambe, *Security and game theory: algorithms, deployed systems, lessons learned.* Cambridge university press, 2011.

[10] S. Cooney, K. Wang, E. Bondi, T. Nguyen, P. Vayano, H. Winetrobe, E. A. Cranford, C. Gonzalez, C. Lebiere, Tambe, and Milind, "Learning to signal in the goldilocks zone: Improving adversary compliance in security games," in *Joint European Conference on Machine Learning and Knowledge Discovery in Databases (in press)*, Springer, 2019.

[11] E. A. Cranford, C. Lebiere, C. Gonzalez, S. Cooney, P. Vayanos, and M. Tambe, "Learning about cyber deception through simulations: Predictions of human decision making with deceptive signals in stackelberg security games.," in *CogSci*, 2018.

[12] E. A. Cranford, C. Gonzalez, P. Aggarwal, S. Cooney, M. Tambe, and C. Lebiere, "Towards personalized deceptive signaling for cyber defense using cognitive models.," in *International Conference on Cognitive Modeling (in press)*, 2019.

[13] E. A. Cranford, C. Lebiere, P. Aggarwal, C. Gonzalez, and M. Tambe, "Adaptive cyber deception: Cognitively-informed signaling for cyber defense," in *Proceedings of the 53rd Hawaii International Conference on System Sciences (submitted)*, IEEE, 2020.

[14] J. Pita, M. Jain, F. Ordónez, C. Portway, M. Tambe, C. Western, P. Paruchuri, and S. Kraus, "Using game theory for los angeles airport security," *AI magazine*, vol. 30, no. 1, pp. 43–43, 2009.

[15] F. Fang, P. Stone, and M. Tambe, "When security games go green: Designing defender strategies to prevent poaching and illegal fishing," in *Twenty-Fourth International Joint Conference on Artificial Intelligence*, 2015.

[16] H. Xu, Z. Rabinovich, S. Dughmi, and M. Tambe, "Exploring information asymmetry in two-stage security games," in *Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015.

[17] S. Cooney, P. Vayanos, T. H. Nguyen, C. Gonzalez, C. Lebiere, E. A. Cranford, and M. Tambe, "Warning time: Optimizing strategic signaling for security against boundedly rational adversaries," in *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems*, pp. 1892–1894, International Foundation for Autonomous Agents and Multiagent Systems, 2019.

[18] A. Sinha, T. H. Nguyen, D. Kar, M. Brown, M. Tambe, and A. X. Jiang, "From physical security to cybersecurity," *Journal of Cybersecurity*, vol. 1, no. 1, pp. 19–35, 2015.

[19] Y. Abbasi, D. Kar, N. Sintov, M. Tambe, N. Ben-Asher, D. Morrison, and C. Gonzalez, "Know your adversary: Insights for a better adversarial behavioral model.," in *CogSci*, 2016.

[20] P. Aggarwal, C. Gonzalez, and V. Dutt, "Cyber-security: role of deception in cyber-attack detection," in *Advances in Human Factors in Cybersecurity*, pp. 85–96, Springer, 2016.

[21] A. Schlenker, O. Thakoor, H. Xu, F. Fang, M. Tambe, L. Tran-Thanh, P. Vayanos, and Y. Vorobeychik, "Deceiving cyber adversaries: A game theoretic approach," in *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, pp. 892–900, International Foundation for Autonomous Agents and Multiagent Systems, 2018.

[22] C. Gonzalez, J. F. Lerch, and C. Lebiere, "Instance-based learning in dynamic decision making," *Cognitive Science*, vol. 27, no. 4, pp. 591–635, 2003.

[23] P. Aggarwal, A. Gautam, V. Agarwal, C. Gonzalez, and V. Dutt, "Hackit: A human-in-the-loop simulation tool for realistic cyber deception experiments," in *International Conference on Applied Human Factors and Ergonomics*, pp. 109–121, Springer, 2019.

[24] R. Chadha, T. Bowen, C.-Y. J. Chiang, Y. M. Gottlieb, A. Poylisher, A. Sapello, C. Serban, S. Sugrim, G. Walther, L. M. Marvel, *et al.*, "Cybervan: A cyber security virtual assured network testbed," in *MILCOM 2016-2016 IEEE Military Communications Conference*, pp. 1125–1130, IEEE, 2016.

[25] T. H. Nguyen, R. Yang, A. Azaria, S. Kraus, and M. Tambe, "Analyzing the effectiveness of adversary modeling in security games," in *Twenty-Seventh AAAI Conference on Artificial Intelligence*, 2013.

[26] C. Gonzalez, N. Ben-Asher, A. Oltramari, and C. Lebiere, "Cognition and technology," in *Cyber defense and situational awareness*, pp. 93–117, Springer, 2014.

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

[28] "2019 data breach investigations report," may 2019.

[29] E. Al-Shaer, J. Wei, K. W. Hamlen, and C. Wang, "Honeypot deception tactics," in *Autonomous Cyber Deception*, pp. 35–45, Springer, 2019.

[30] D. Gambetta, *Signaling*, p. 168194. Oxford University Press, 2011.

[31] J. M. Martin, C. Gonzalez, I. Juvina, and C. Lebiere, "A description–experience gap in social interactions: Information about interdependence and its effects on cooperation," *Journal of Behavioral Decision Making*, vol. 27, no. 4, pp. 349–362, 2014.

[32] C. Gonzalez, J. Dana, H. Koshino, and M. Just, "The framing effect and risky decisions: Examining cognitive functions with fmri," *Journal of economic psychology*, vol. 26, no. 1, pp. 1–20, 2005.

[33] M. Gutierrez, J. Černỳ, E. Aharonov, N. Ben-Asher, B. Boansk, C. Kiekintveld, and C. Gonzalez, "Evaluating models of human adversarial behavior against defense algorithms in a contextual multi-armed bandit task.," in *International Conference on Cognitive Modeling (in press)*, 2019.