## Human and AI Decision-Making in Cybersecurity: A Multiagent Modeling Perspective

### Yinuo Du

CMU-S3D-25-111 August 2025

Software and Societal Systems Department School of Computer Science Carnegie Mellon University Pittsburgh, PA 15213

#### **Thesis Committee:**

Fei Fang (Co-chair)
Cleotilde Gonzalez (Co-chair)
Christian Lebiere
Prashanth Rajivan (University of Washington)
Tiffany Bao (Arizona State University)

Submitted in partial fulfillment of the requirements for the degree of Doctor of Philosophy in Societal Computing

Copyright © 2025 Yinuo Du

This research was sponsored by MACRO: Models of Enabling Continuous Reconfigurability of Secure Missions Cyber-Security Collaborative Research Alliance (ARL). The views and conclusions contained herein are those of the author and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the Army Research Laboratory, the U.S. Army, the Department of Defense, or the U.S. Government.



No matter how hopeless, no matter how far.

### **Abstract**

The dynamic nature of cyber threats presents significant challenges for modern defense, as sophisticated adversaries continuously adapt their strategies to evade detection and compromise valuable systems. Effective defense against these evolving threats requires multiagent interaction, where human defenders must coordinate with both other humans and AI systems to mount comprehensive responses. However, current approaches fail to adequately model the cognitive mechanisms underlying multiagent interactions in these complex environments. Without computational models of how humans adapt, collaborate, and make decisions in cybersecurity contexts, we cannot build multiagent defense systems that leverage the full potential of human and AI.

This thesis focuses on building computational cognitive models and cognitive agents for multiagent interaction in cyber defense, including designing adversarial cognitive agents (Chapter 3), modeling human decision-making in multi-defender interaction (Chapter 4), and designing human-like AI agents that can work with humans as a team (Chapter 5).

First, I investigate human behavior in cybersecurity at the individual level and build adversarial cognitive agents that capture human-like adaptivity in cyber attack, which pre-sent greater challenges to defenders than deterministic strategies. My findings show that cognitive attackers driven by Instance-Based Learning Theory can learn effective strategies that are more challenging for both human and autonomous defenders to counter than optimal but predictable attack patterns.

Second, I explore cognitive mechanisms that enable effective decision-making in multi-defender interactions. In cybersecurity, multiple defenders can share sensitive information and collaborate on threat response, however, their willingness to do so could impact the security posture of all connected defenders. I develop a novel computational model for interdependent human decision-making and investigate its validity in multi-defender interaction setting. The model incorporates three key cognitive mechanisms: dynamic prosociality, which adjusts how individuals value others' outcomes based on expectation-reality discrepancies; category learning, which efficiently organizes social experiences into behavioral prototypes; and contrast effects, which sharpen distinctions between these behavioral categories.

Finally, I examine the integration of human and AI decision-making in team defense scenarios where humans and AI collaboratively protect computer networks. I designed an AI agent that learns from experience to approximate human-like decision processes. Through empirical studies in semi-supervisory frameworks, I demonstrate that the human-like AI agent significantly enhances team performance and efficiency in cybersecurity operations compared to heuristic or random agents.

## Acknowledgments

This PhD journey has been a tremendous learning experience that has shaped me both as a researcher and as a person. I am deeply grateful to the many individuals who have supported, guided, and inspired me throughout this process.

First and foremost, I would like to express my profound gratitude to my advisors, Dr. Cleotilde Gonzalez and Dr. Fei Fang. Coty has shown me what it takes to build a meaningful research career through his brilliance, unwavering discipline, and exemplary mentorship. Fei has taught me to make every hour count and to strive for excellence in all endeavors. Every interaction with her has reminded me of the heights to which scholarship can aspire. I am deeply grateful to both for taking on far more responsibilities than any advisor should reasonably bear, and for their patience and dedication in guiding my development as a researcher.

I extend my sincere appreciation to my committee members for their invaluable contributions to this work. Thank you to Dr. Prashanth Rajivan for supporting my exploratory research when there was little evidence I could successfully execute it. Thank you to Dr. Christian Liebere for providing crucial guidance in refining both my proposal and dissertation, offering insights that substantially strengthened this research. Thank you to Dr. Tiffany Bao for bringing essential cybersecurity perspectives that enriched the interdisciplinary nature of this work.

My colleagues at DDMLab have been instrumental in my growth as a researcher. Thank you to Maria José Ferreira, for sitting through my dreadful practice talks and helping me launch new projects with enthusiasm and insight. Thank you to Ty Malloy, for revising my often-confusing early drafts but, more importantly, for teaching me the true spirit of academic collaboration. Thank you to Roderick Seow, for always being a ray of sunshine in the lab and the coffee chats. Thank you to Baptiste Prebot, for helping me to kick start my first project in PhD, and for working with me when I had little research experience. Thank you to Erin Bugbee for being the perfect reference point and demonstrating how a PhD student should be. Thank you to Chase McDonald for helping me with so many tasks, always with generosity and good humor, despite my inability to adequately reciprocate his kindness. I am also grateful to my colleagues at the AISOC lab: Chunkai Ling, Ryan Shi, Steven Jacmen, Stephanie Milani, Rex Chen, Zhicheng Zhang, Naveen Ramen, Jingwu Tang, Yixuan Xu. It has been a true privilege to have a front-row seat to observe how exceptional research is conducted and to learn from such talented researchers.

Thank you to my cohort at S3D and SDS. It has been an honor to be amongst such inspiring peers. Thank you to the WiCyS community for the support and camaraderie throughout this journey. I will always have fond memories of our Friday night meetings.

Finally, I owe an immense debt to Barbara and Halley, who shouldered the burden of listening to my frustrations and tolerating me at my worst. This PhD would not have been possible without your support.

## **Contents**

1	Introduction		
2	Background		
	2.1	The Human Element in Cybersecurity	4
	2.2	Human Cognition in Decision Making	6
	2.3	Human Model Validation in Cybersecurity	9
	2.4	Synthetic Cyber Environments for Human Experiments	10
3		gle-Agent Decision-Making in Cybersecurity: Building versarial Cognitive Agents	12
	3.1	Introduction	12
	3.2	Related work	13
	3.3	Design of the Simulated Cyber Attack Environment	18
	3.4	Experiment: Cognitive Attacker Against Human Defenders	20
	3.5	Discussion and Conclusion	23
4		alti-Agent Decision-Making in Cybersecurity: Cognitive echanisms for Multi-Defender Interaction	26
	4.1	Introduction	26
	4.2	Related work	27
	4.3	Emergent Cooperative Decision-making in Triadic Prisoner's  Dilemma: Effects of Incentives and Information	33
	4.4	Toward a Cognitive Theory of Interdependent Decisions in Groups: Dynamic Weighting, Categorization, and Contrast	40
	4.5	Discussion and Conclusion	44
5		man-AI Teaming in Cyber Defense: Enhancing Collabo- ive Performance Through Cognitive Integration	48

	5.1	Introduction	48
	5.2	Related work: Human-Autonomy Teaming	49
	5.3	Design of the Team Defense Game	52
	5.4	Experiment: Human-Autonomy Cyber Defense Team Performance	54
	5.5	Discussion and Conclusion	57
6	Co	nclusion & Future Directions	60
	6.1	Conclusion	60
	6.2	Future Work: Human-like Adversaries Modeling	62
	6.3	Future Work: Toward Autonomous Intelligent Cyber Defense	63
	6.4	Future Work: Complementary Human-AI Teaming in Cyber Defense	66
A	ppe	ndix	87
	A C	yber-War Between Bots: Cognitive Attackers are More Challenging for Defenders than Strategic Attackers (Accepted by <i>ACM Transactions of Social Computing</i> ) .	87
	Lear	ning about simulated adversaries from human defenders using interactive cyber- defense games (Accepted by <i>Journal of Cybersecurity</i> )	110
	Expo	erimental evaluation of cognitive agents for collaboration in human-autonomy cyber defense teams (Accepted by <i>Computers in Human Behavior: Artificial Human</i> ) .	123
	Eme	rgent Cooperative Decision-making in Triadic Prisoner's Dilemmas: Effects of Incentives and Information (Accepted by <i>Acta Psychologica</i> )	141
	Tow	ard a Cognitive Theory of Interdependent Decisions in Groups: Dynamic Prosociality, Categorization, and Contrast (Under Review at <i>Psychological Review</i> )	163

## Chapter 1

### Introduction

The cybersecurity landscape continues to face unprecedented challenges. The Identity Theft Resource Center reported a 78% increase in data breaches in the US, with 3,205 incidents affecting 353 million individuals in 2023. Enterprise networks have expanded to include thousands to tens of thousands of devices, creating vast attack surfaces with various entry points for adversaries. The consequences of this expanding attack surface were dramatically illustrated by the 2017 WannaCry ransomware attack, which exploited unpatched vulnerabilities that affected more than 230,000 computers in 150 countries, resulting in damages estimated at billions of dollars. As our technological infrastructure grows increasingly interconnected, the imperative to effectively detect and mitigate cyber-attacks becomes critical for the secure operation of society's essential systems.

The dynamic nature of cyber threats presents significant challenges for modern defense, as sophisticated adversaries continuously adapt their strategies to evade detection and compromise valuable systems. Unlike traditional security approaches that rely on static defenses and predictable patterns, modern attackers exhibit human-like adaptivity, learning from defensive responses and modifying their tactics accordingly. This evolution in the threat landscape has fundamentally changed the requirements for effective defense. Effective defense against these evolving threats requires multiagent interaction, where human defenders must coordinate with both other humans and AI systems to mount comprehensive responses. The complexity of modern networks and the speed of cyber operations make it impossible for individual defenders or isolated systems to maintain security. Instead, defense must emerge from the coordinated actions of multiple agents working together.

Cyber defense involves strategic interactions among attackers, defenders, and end-users, each with different objectives and incomplete information about the network and other agents. The security status at any moment depends on exogenous events and the strategies of all agents, which are rarely common knowledge. Adversaries exploit this complexity through actions that conceal their true intent, while defenders must anticipate and counter these moves. The adversarial landscape is further complicated by the diversity of threat actors, from sophisticated nation states to opportunistic script kiddies, who differ in motivation, resources, target selection, and persistence. This dynamic environment, characterized by continuous adaptation and counter-adaptation between attackers and defenders, makes the configuration of effective defense strategies extremely challenging.

Integrating human and artificial intelligence elements in decision-making processes introduces additional complexity. Determining the appropriate level of autonomy for AI systems involves complex trade-offs between immediate security responses and longer-term strategic considerations that often require human judgment. Effective human-AI collaboration requires clear communication and mutual understanding of capabilities and limitations. AI tools must be not only reliable but also interpretable and cooperative, while cybersecurity professionals require training to effectively leverage these technologies. The complexity increases further when collaboration extends across organizational boundaries, a necessity for robust defense against sophisticated threats, but hindered by varying objectives, priorities, cultures, and trust issues among different organizations. However, current approaches fail to adequately model the cognitive mechanisms underlying multiagent interactions in these complex environments. Without computational models of how humans adapt, collaborate, and make decisions in cybersecurity contexts, we cannot build multiagent defense systems that leverage the full potential of human and AI coordination.

First, I investigate human behavior in cybersecurity at the individual level and build adversarial cognitive agents that capture human-like adaptivity in cyber attacks, which present greater challenges to defenders than deterministic strategies. I developed cognitive agents based on Instance-Based Learning Theory (IBLT) that learn from interaction experiences to simulate human-like attackers. Through experimental comparisons with both strategic (Beeline) and stochastic (Meander) attackers, I demonstrated that these cognitive attackers present greater challenges to both human defenders and autonom-ous defensive agents. The experiments showed that defenders were able to effectively learn and adapt against deterministic attack strategies, but struggled significantly against cognitive attackers that dynamically adjusted their tactics based on the defender's behavior. My findings show that cognitive attackers driven by Instance-Based Learning Theory can learn effective strategies that are more challenging for both human and autonom-ous defenders to counter than optimal but predictable attack patterns. This finding has important implications for cybersecurity training, highlighting the need to prepare against adaptive human-like adversaries rather than only focusing on countering known attack patterns.

Second, I explore cognitive mechanisms that enable effective decision-making in multi-defender interactions. In cybersecurity, multiple defenders can share sensitive information and collaborate on threat response; however, their willingness to do so could impact the security posture of all connected defenders. Through empirical studies using triadic Prisoner's Dilemma scenarios framed as a cybersecurity information sharing task, I systematically varied both the structural incentives (K-index) and information availability to understand their effects on cooperation. The results demonstrated that higher structural incentives promote stable cooperation by reducing the temptation to defect, while experiential information (observing actions and outcomes) significantly enhances cooperation compared to both minimal information and overly detailed descriptive information (complete payoff matrices). I develop a novel computational model for interdependent human decision-making and investigate its validity in a multi-defender interaction setting. The model incorporates three key cognitive mechanisms: dynamic prosociality, which adjusts how individuals value others' outcomes based on expectation-reality discrepancies; category learning, which efficiently organizes social experiences into behavioral prototypes; and contrast effects, which sharpen distinctions between these behavioral categories.

Finally, I examine the integration of human and AI decision-making in team defense scenarios

where humans and AI collaboratively protect computer networks. I designed an AI agent that learns from experience to approximate human-like decision processes. I designed and implemented the Team Defense Game (TDG), an experimental platform where human participants collaborate with autonom-ous agents to protect a network against external threats. By systematically comparing three types of autonomous teammates—cognitive agents based on IBLT, heuristic agents using rule-based strategies, and random agents—I investigated how different AI approaches affect team performance and human workload. Through empirical studies in semi-supervisory frameworks, I demonstrate that the human-like AI agent significantly enhances team performance and efficiency in cybersecurity operations compared to heuristic or random agents. The results revealed that cognitive agents were more adaptive to the individual play styles of human teammates, although they were sometimes perceived as inconsistent or unpredictable. Competent agents (both cognitive and heuristic) required less human effort but sometimes led to over-reliance, highlighting important trust calibration challenges in human-autonomy teaming.

My research has made several significant contributions to the field. I developed a human-like adversary emulation method based on Instance-Based Learning Theory, accompanied by an interactive defense game testbed and human-subject experiments showing that cognitive attackers are more challenging for defenders than strategic attackers. These findings demonstrate that training against human-like adversaries is necessary to prepare against diverse adversary strategies. I created an empirical evaluation framework for cross-organizational cooperation that revealed how incentive structures and information availability influence information sharing behaviors. By systematically varying the K-index of interdependence and the levels of information provided to participants, I demonstrated that higher structural incentives promote stable cooperation, while experiential information significantly enhances cooperation compared to minimal or overly detailed descriptive information. I developed a cognitive model of interdependent decisions in groups that integrates dynamic weighting, category learning, and contrast effects to explain how individuals navigate multiple cooperative relationships simultaneously. This model successfully reproduced human behavior patterns in information sharing experiments without parameter fitting, providing insights into the psychological processes underlying group decision-making in security contexts. I designed a team defense game platform and an experimental protocol that enabled the systematic evaluation of human-AI collaboration in cybersecurity contexts. Through controlled experiments, I demonstrated that cognitive agents based on instance-based learning theory outperform both heuristic and random agents as teammates, improving both team performance and human efficiency in cyber defense tasks.

In the following chapters, I will detail these contributions, beginning with a review of related work on human and AI decision-making in cyber defense, followed by chapters describing my completed research projects and their findings. We will conclude with a discussion of the implications of this work for the future of cybersecurity operations and suggestions for future research directions.

## Chapter 2

## Background

The landscape of cybersecurity continues to evolve with increasingly sophisticated threats and defensive technologies. At the center of this evolution remains a critical component: the human element. Whether as attackers, defenders, or end-users, humans significantly influence cybersecurity outcomes through their decisions, behaviors, and cognitive processes. This chapter establishes the foundational concepts that inform our research, examining the human element in cybersecurity, the cognitive architectures that model human decision-making, and methodological approaches for studying these phenomena.

### 2.1 The Human Element in Cybersecurity

Cybersecurity has traditionally been approached as a technical challenge, with emphasis placed on developing robust security mechanisms, detection systems, and defense-in-depth strategies. However, research now recognizes that the effectiveness of these technical solutions depends critically on human factors [210, 263, 46]. The human dimension of cybersecurity includes both offensive and defensive aspects, with significant implications for system vulnerability and resilience.

Human attackers show considerable variation in capability, motivation, and behavior that influence their effectiveness. Unlike idealized computational attackers, human adversaries operate with cognitive limitations that affect their decision-making [217, 218]. These limitations, characterized as bounded rationality, often lead to suboptimal choices based on heuristics rather than exhaustive analysis. Research by Oh et al. [181] and Alsharnouby et al. [8] shows that these limitations create patterns in attack strategies that can be exploited by defensive systems that understand human decision-making processes.

The risk tolerance of human attackers affects their choice of targets and attack methods [257, 110]. As documented by Thomas and Sule [236], an adversary's risk appetite can be affected by their situational context and operational goals, emphasizing the need for continuous threat assessment. Human attackers learn from their experiences and dynamically adapt to encountered defenses, modifying their strategies accordingly [137, 89]. This adaptability makes them increasingly dangerous over time as they become more adept at evading detection and exploiting vulnerabilities. Studies by Shoetan et al. [215] provide insight into how adversaries'

risk tolerance levels can adjust based on their past successes or failures, indicating a cycle of risk evaluation and reassessment that influences future attack strategies.

Human attackers also leverage creativity when developing novel exploitation techniques. Studies of phishing campaigns show that individual creativity predicts an adversary's ability to evade detection [198]. This capacity for innovation allows human attackers to discover vulnerabilities and attack vectors that automated testing might miss. Research by Shashank et al. [213] highlights how attacker creativity manifests in developing polymorphic malware that continuously modifies its code structure and encryption patterns while maintaining functionality. This enables it to evade signature-based detection systems that rely on static patterns. Similarly, Huang and colleagues [111] documented cases where human attackers creatively chained together seemingly low-risk vulnerabilities across different system components. They exploited subtle trust relationships between authentication systems and application interfaces to achieve privilege escalation through paths that automated vulnerability scanners consistently failed to identify or prioritize.

On the defensive side, the effectiveness of cybersecurity defense depends equally on human factors. Security operations centers (SOCs) face challenges with alert fatigue, where the volume of security alerts overwhelms human analysts [176, 62]. This cognitive overload can lead to missed detections of critical threats and represents a limitation in scaling human defensive capabilities. Ban et al. [18] documented how alert fatigue leads to frustration and performance degradation among security analysts, highlighting the need for AI-assisted techniques to combat this problem. Beyond alert fatigue, defenders struggle with cognitive biases in risk perception that affect security decision-making. Pfleeger and Caputo [190] show how inattentional blindness prevents security professionals from noticing unexpected threats when focusing on primary tasks. Work overload also impacts cybersecurity behavior, with research showing that excessive job demands lead to burnout and compromised security practices [139]. Human vulnerabilities extend to organizational contexts where interdepartmental coordination failures create security gaps. Hadlington [95] found that even well-trained defenders struggle to maintain vigilance in environments with poor security culture.

Defensive decisions frequently occur under conditions of uncertainty, time pressure, and incomplete information. Research shows that these conditions affect risk assessment and response selection [24, 59], often leading to suboptimal defensive strategies when the cognitive demands exceed human capabilities. Under uncertainty, defenders must navigate an asymmetric information environment where attackers can observe defense mechanisms while defenders have limited visibility into attack methods and origins [263, 30]. This information disadvantage forces defenders to make decisions with partial situational awareness, increasing reliance on heuristic judgments rather than comprehensive analysis. When operating under severe time constraints, defenders shift toward rapid, intuitive decision processes that prioritize immediate action over analytical assessment [38, 51]. As shown by Hwang [112], time pressure alters decision strategies, increasing reliance on recognition rather than calculation. This challenge is heightened by the compressed timeframes of automated attacks. Defenders require extensive experience to effectively counter dynamic and distributed attacks [125, 92]. The development of this expertise follows patterns identified in studies of naturalistic decision making, where recognition-primed decisions based on prior experiences guide expert responses [126, 125]. Klein's model explains how cybersecurity experts leverage pattern recognition to rapidly identify threat situations without needing to compare multiple response options, enabling effective responses even with incomplete information [128, 127]. However, this expertise-driven approach faces challenges in cyberse-curity contexts due to the rapidly evolving threat landscape and the difficulty of accumulating relevant experiences for novel attack vectors [224].

The traditional approach to improving defensive capabilities has focused on cybersecurity training using games and simulations. Platforms like picoCTF and SecGen offer progressively challenging tasks structured within immersive narratives, making complex concepts more accessible [195, 171]. Studies by Hendrix et al. [104] and Tioh et al. [238] show increased knowledge retention and skill acquisition when participants actively engage in game formats compared to traditional learning approaches. However, these training platforms often rely on deterministic adversary models that fail to capture the dynamic and adaptive nature of real human attackers, potentially giving defenders false confidence in their abilities to counter actual threats [1, 97]. This limitation presents a research opportunity to develop more realistic adversary models that better reflect human cognitive processes and adaptability. By incorporating principles from cognitive science into adversary simulation, as proposed in this dissertation, training environments can better prepare defenders for the unpredictable and evolving tactics employed by human attackers in operational environments.

### 2.2 Human Cognition in Decision Making

Modeling how humans make decisions in cybersecurity contexts requires frameworks that account for cognitive limitations while capturing learning from experience. Two approaches provide the foundation for our research: cognitive architectures that implement bounded rationality and Instance-Based Learning Theory [81].

Cognitive architectures provide computational frameworks for modeling human cognition, capturing constraints like memory limitations, attention bottlenecks, and information processing capabilities [9, 173]. Unlike purely algorithmic approaches that optimize toward theoretical ideals, cognitive architectures prioritize psychological plausibility by adhering to known constraints of human cognition. As Salvucci and Taatgen [206] note, these architectures allow for the integration of multiple theoretical accounts of cognition into unified computational systems that can reproduce human behavior across diverse tasks.

The concept of bounded rationality, introduced by Simon [217, 218], recognizes that human decision-makers operate with limited information about alternatives and consequences, cognitive constraints that restrict computation and memory, and finite time for decision-making. These limitations lead humans to employ heuristics and satisficing strategies rather than exhaustively evaluating all possibilities [74, 239]. In cybersecurity contexts, bounded rationality manifests in both attackers and defenders, creating predictable patterns of behavior that diverge from theoretically optimal strategies [78, 45].

Studies by Moisan et al. [165] show that these cognitive limitations affect cooperation rates in strategic interactions, with implications for security information sharing and coordination. Research by Abbasi et al. [1] and Hamman et al. [97] has shown that understanding these limitations can inform more effective training approaches for cybersecurity professionals by aligning training scenarios with the cognitive constraints that shape real-world decision-making.

### 2.2.1 Instance-Based Learning Theory

**Activation** IBL models work by storing instances i in memory  $\mathcal{M}$ , composed of utility outcomes  $u_i$  and options k composed of features j in the set of features  $\mathcal{F}$  of environmental decision alternatives. These options are observed in an order represented by the time step t, and the time step that an instance occurred in is given  $\mathcal{T}(i)$ . IBL models predict the value of options in decision-making tasks by selecting the action that maximizes the value function. In calculating this activation, the similarity between instances in memory and the current instance is represented by summing over all attributes the value  $S_{ij}$ , which is the similarity of attribute j of instance i to the current state. This gives the activation equation as:

$$A_i(t) = \ln\left(\sum_{t' \in \mathcal{T}_i(t)} (t - t')^{-d}\right) + \mu \sum_{j \in \mathcal{F}} \omega_j (S_{ij} - 1) + \sigma \xi$$
(2.1)

The parameters that are set either by modelers or set to default values are the decay parameter d; the mismatch penalty  $\mu$ ; the attribute weight of each j feature  $\omega_j$ ; and the noise parameter  $\sigma$ . The default values for these parameters are  $(d, \mu, \omega_j, \sigma) = (0.5, 1, 1, 0.25)$ . The value  $\xi$  is drawn from a normal distribution  $\mathcal{N}(-1, 1)$  and multiplied by the noise parameter  $\sigma$  to add random noise to the activation.

**Probability of Retrieval** The probability of retrieval represents the probability that a single instance in memory will be retrieved when estimating the value associated with an option. To calculate this probability of retrieval, IBL models apply a weighted soft-max function onto the memory instance activation values  $A_i(t)$  giving the equation:

$$P_i(t) = \frac{\exp A_i(t)/\tau}{\sum_{i' \in \mathcal{M}_i} \exp A_{i'}(t)/\tau}$$
(2.2)

The parameter that is either set by modelers or set to its default value is the temperature parameter  $\tau$ , which controls the uniformity of the probability distribution defined by this soft-max equation. The default value for this parameter is  $\tau = \sigma \sqrt{2}$ .

**Blended Value** The blended value of an option k is calculated at time step t according to the utility outcomes  $u_i$  weighted by the probability of retrieval of that instance  $P_i$  and summing over all instances in memory  $\mathcal{M}_k$  to give the equation:

$$V_k(t) = \sum_{i \in \mathcal{M}_k} P_i(t) u_i \tag{2.3}$$

### 2.2.2 IBLT in Cybersecurity Applications

IBLT has been successfully applied to model various aspects of cybersecurity decision-making across both offensive and defensive contexts. The theory provides a psychologically plausible account of how cybersecurity professionals and attackers learn from experience and make decisions in dynamic environments characterized by uncertainty and incomplete information.

**Input:** default utility  $u_0$ , a memory dictionary  $\mathcal{M} = \{\}$ , global counter t = 1, step limit L, a flag delayed to indicate whether feedback is delayed.

```
repeat
   Initialize a counter (i.e., step) l=0 and observe state s_l
   while s_l is not terminal and l < L do
       Execution Loop
           Exploration Loop k \in K do
               Compute activation values A_i(t) of instances (k_i, T(i)) by Eq. (2.1)
               Compute retrieval probabilities P_i(t) by Eq. (2.2)
               Compute blended values V_k(t) corresponding to k by Eq. (2.3)
           end
           Choose an action a corresponding to option k_l \in \arg \max_{k \in K} V_k(t)
       end
       Take action a, move to state s_{l+1}, observe s_{l+1}, and receive outcome u_{l+1}
       Store t into instance corresponding to selecting k_l and achieving outcome u_{l+1} in
       If delayed is true, update outcomes using a credit assignment mechanism
       l \leftarrow l + 1 and t \leftarrow t + 1
   end
until task stopping condition
```

Algorithm 1: Pseudo Code of Instance-Based Learning Process

Dutt et al. [60] showed how IBLT can model the situational awareness of cybersecurity analysts, making concrete predictions about recognition and comprehension processes during attack scenarios. Their model captured how analysts interpret security events based on similar past experiences, explaining why analysts with different experiential backgrounds might reach different conclusions when presented with identical security data. This work was extended by Veksler et al. [244], who applied IBLT to predict attacker behavior and enhance analyst capability to anticipate future threats.

In defensive contexts, Du et al. [57] showed that IBLT-based models can effectively capture defender behaviors against various attack strategies. Their cognitive model learned to identify patterns in attack sequences and adapt defensive responses accordingly, showing performance comparable to that of human defenders with similar experience levels. Cranford et al. [45] developed these models to incorporate social factors that influence defensive decision-making, such as trust and deception, providing insights into how these psychological factors affect security outcomes.

IBLT has proven useful for modeling phishing detection and response. Cranford et al. [44] used IBLT to model how end-users identify phishing attempts, showing how variations in experience and attention to specific features lead to different vulnerability patterns across users. Similarly, research by Aggarwal et al. [4] explored how IBLT models can capture decision-making about cyber attacks, revealing how attackers gradually learn effective strategies through trial and error rather than through formal planning processes.

The application of IBLT to adaptive cyber defense has been explored by Lebiere et al.

[138], who showed how cognitive models can support the development of autonomous defense systems that reflect human-like adaptivity while overcoming human limitations in processing speed and attention. Building on this work, Gonzalez et al. [84] proposed that cognitive models can serve as the foundation for predicting both attacker and defender behavior in cybersecurity contexts, potentially enabling more effective defensive strategies through improved anticipation of adversary actions.

Researchers have also applied IBLT to understand information sharing decisions in cyberse-curity contexts. Monleon et al. [169] developed models that capture how trust and experience influence decisions to share or withhold threat intelligence, showing patterns that closely match those observed in human security professionals. Similarly, Nguyen and Gonzalez [174] extended IBLT to incorporate theory of mind capabilities, allowing models to predict not only direct adversary actions but also how adversaries might reason about defender knowledge and strategies.

Computational models based on IBLT differ from traditional machine learning or game-theoretic approaches in cybersecurity by incorporating realistic cognitive constraints rather than assuming unlimited computational resources, learning through experiences rather than requiring explicit rule formalization, and generalizing to novel situations based on similarity to previous experiences. As shown by Kheiri et al. [120] and Sanchez et al. [207], these properties make IBLT-based models particularly suitable for predicting how humans will behave in complex and dynamic cybersecurity scenarios, where conditions change rapidly and decision-makers must adapt to evolving threats and defenses.

### 2.3 Human Model Validation in Cybersecurity

Several frameworks have been proposed for validating human models and studying human behavior in cybersecurity contexts. Early frameworks focused primarily on technical aspects of security [212], but researchers increasingly recognized the need for structured approaches to studying human factors. Kraemer et al. [132] proposed a human factors taxonomy that categorized cybersecurity failures based on cognitive and organizational dimensions. Pfleeger and Caputo [190] developed a behavioral framework emphasizing psychological factors in security decision-making. More recently, Gonzalez et al. [79] introduced a comprehensive multi-dimensional framework that integrates levels of analysis, contextual representation, and cognitive complexity. This thesis primarily builds upon the Gonzalez et al. framework due to its systematic approach to experimental design and its explicit consideration of dynamic, interactive decision-making that characterizes modern cyber threats.

Building on the framework proposed by Gonzalez et al. [79], our research employs a structured approach that considers multiple dimensions and experimental paradigms. Cybersecurity research spans multiple dimensions that must be considered when designing experiments and analyzing results. These dimensions include the level of analysis (ranging from individual decision-makers to teams and organizations), the contextual representation (from abstract tasks to naturalistic environments), and the cognitive complexity (from static to dynamic and interactive decision-making).

At the individual level, research focuses on the cognitive processes that influence a single

decision-maker's actions and vulnerabilities. This includes studies of how security professionals detect threats [24, 199] and how end-users respond to phishing attempts [246, 245]. Moving to the dyadic level, research examines direct interactions between attackers and defenders, exploring how each adapts to the other's actions over time [1, 155]. This dyadic perspective provides insights into the strategic dynamics that characterize cybersecurity encounters.

Group-level analysis extends beyond dyads to examine how multiple agents interact in networked environments. Research at this level has explored phenomena such as information sharing [161, 253] and coordinated defense strategies [154, 201]. The team level focuses on coordinated action among interdependent agents with complementary roles, examining how team composition and communication patterns influence defensive effectiveness [226, 28].

The representation of cybersecurity contexts ranges from highly abstract to richly naturalistic. Abstract representations, such as game-theoretic formulations [7, 204], isolate key decision mechanisms at the cost of ecological validity. Contextual representations incorporate domain-specific elements while maintaining experimental control, as seen in simulated network defense scenarios [199, 155]. Naturalistic representations provide high-fidelity simulations that closely mirror real-world environments [34, 85], maximizing ecological validity but often reducing experimental control.

The cognitive complexity of cybersecurity tasks varies widely. Static decision-making involves one-time choices with fixed parameters, while sequential decision-making introduces path dependencies as choices unfold over time [119, 248]. Dynamic decision-making adds complexity by allowing parameters to change based on actions [120, 84], and interactive decision-making introduces multiple agents whose actions influence outcomes [45, 138].

### 2.4 Synthetic Cyber Environments for Human Experiments

Various synthetic environments have been developed to study human behavior in cybersecurity contexts, each offering different trade-offs between experimental control and ecological validity. Early platforms like CyberCIEGE [113] and CyberProtect [59] focused on resource allocation and security policy decisions. More sophisticated environments such as DETER [25] and SAIC's cyber range [68] provide high-fidelity network simulations but require significant technical expertise. Game-based platforms like HackIT [5] and abstract security games [166] offer accessible interfaces for studying fundamental decision processes. Among these approaches, interactive defense games strike a balance between experimental control and contextual relevance, making them particularly suitable for studying cognitive mechanisms in cyber defense while maintaining accessibility to diverse participant populations.

Our research employs interactive defense games as a primary experimental paradigm for studying cybersecurity decision-making. These games provide controlled environments for examining attacker-defender interactions while abstracting away technical details that might limit participation to specialized populations. By systematically varying game parameters, we can isolate the effects of specific factors on decision-making.

Interactive defense games typically involve participants making sequential decisions about system protection, intrusion detection, or information sharing in the context of simulated cyber threats. For example, the CyberCIEGE platform [113] places participants in the role of security

decision-makers who must allocate resources to protect networked systems while balancing security against usability concerns. Similarly, the TRACER platform [155] simulates network defense scenarios where participants must detect and respond to evolving attack patterns.

These games offer several methodological advantages. First, they allow for precise control over experimental variables such as threat characteristics, resource constraints, and information availability. Second, they facilitate detailed logging of decision processes, including choices, reaction times, and information access patterns. Third, they enable systematic manipulation of cognitive factors such as time pressure, uncertainty, and feedback timing.

Our approach to interactive defense games incorporates elements from both traditional experimental paradigms and more naturalistic simulations. By calibrating game parameters based on real-world cybersecurity challenges while maintaining experimental control, we aim to balance internal and external validity. This approach aligns with recommendations from Gonzalez et al. [79], who emphasize the importance of systematically varying contextual elements to identify generalizable principles of cybersecurity decision-making.

This chapter covered the foundational concepts for studying human and AI decision-making in cybersecurity from a multi-agent modeling perspective. I reviewed the critical role of human factors in both attack and defense, introduced cognitive frameworks for modeling human-like decision processes, and outlined methodological approaches for systematically investigating these phenomena. In the following chapters, I build upon these foundations to examine specific aspects of human and AI decision-making in cybersecurity, focusing on cognitive modeling of adversarial behavior, multi-agent decision-making, and human-AI teaming in defense scenarios.

## Chapter 3

# Single-Agent Decision-Making in Cybersecurity: Building Adversarial Cognitive Agents<sup>1</sup>

### 3.1 Introduction

In the modern cybersecurity landscape, understanding attacker behavior is fundamental to effective defense. As established in Chapter 1, cyber defense inherently involves strategic interactions between multiple agents with different objectives and incomplete information. Organizations frequently use cyber wargaming and adversary emulation (i.e., Red Teams) to train defenders (i.e., Blue Teams) and develop appropriate defense algorithms [41, 69]. However, traditional adversary emulation methods often rely on automated planning and underplay the role of human cognition, consequently leaving defenders underprepared for human attackers who can think creatively and adapt their strategies.

Existing automated adversary simulation methods primarily rely on deterministic patterns or static behavioral models [96, 121, 2]. While these approaches offer technical fidelity, they typically lack the dynamic adaptivity characteristic of human attackers [117]. Real adversaries vary in risk tolerance, learn from their experiences, and modify their strategies in response to defensive measures [258, 136]. As they interact with defenders, they become increasingly adept at evading detection and exploiting vulnerabilities. This adaptivity presents a significant challenge for defender training and highlights the limitations of conventional adversary emulation techniques.

To address these limitations, cognitive models offer a promising approach. Unlike traditional computational methods focused on optimal performance, cognitive models incorporate human constraints such as forgetting, limited attention, and bounded rationality [80]. These models can simulate the learning and adaptive processes characteristic of human attackers, potentially providing more realistic and challenging training scenarios for defenders. Instance-Based Learning Theory (IBLT) [81] provides a particularly suitable framework for modeling human-like

<sup>&</sup>lt;sup>1</sup>See Appendix .1 for published version of this chapter and Appendix .2 for more about the Interactive Defense Game

decision-making in dynamic cybersecurity contexts. Prior research has applied IBLT to model cyber situation awareness of human analysts [61] and to develop defensive agents capable of learning to counter deterministic attacks [58]. However, research has rarely examined the real-time social interactions between cognitive attackers and defenders in cybersecurity scenarios.

The strategic interplay between attackers and defenders creates unique dynamics that can influence defensive effectiveness [250]. Notably, human defenders have been found to handle random attacks more effectively than adaptive ones [167], suggesting that commonly used random attack algorithms may be less effective for training than approaches that capture human-like adaptivity. This finding aligns with studies of human attackers in phishing contexts, which have demonstrated that individual creativity significantly predicts an adversary's ability to evade detection [200]. Cognitive biases and emotional factors further influence attacker behavior and decision-making processes [115, 66], adding another layer of complexity that static models cannot capture.

Building on these insights, this chapter demonstrates that cognitive agents built based on the theoretical principles of Instance-Based Learning Theory make more challenging adversaries for defenders than strategically optimal attackers. Our research offers three main contributions. First, I develop a cognitive attacker model ( $IBL_{Red}$ ) and demonstrate how it can learn from experience to become as efficient as optimal strategic algorithms against a strategic defender. Second, I show that when pitted against cognitive defenders, the IBL attacker proves to be a more challenging adversary while the IBL defender can learn to counter carefully crafted optimal attack strategies. Third, I validate these findings through human experiments where participants play the defender role, confirming that cognitive attackers are more challenging for human defenders than strategic attackers.

These experiments contribute to both cybersecurity practice and cognitive science by showing how cognitive agents that capture human-like adaptivity create more effective adversaries for training purposes and by advancing our understanding of strategic interactions in cybersecurity contexts. The findings inform future adversary emulation efforts and the training of cyber defenders, demonstrating that preparation against adaptable, human-like adversaries better equips defenders for real-world threats than traditional approaches focused on optimal but predictable attack patterns.

### 3.2 Related work

### 3.2.1 Threat Modeling

Attackers in cyberspace range from novice script kiddies to highly organized state-sponsored actors, each motivated by varying goals and equipped with different levels of skill, knowledge, resources, access, and motives. This diversity in threat actors has made threat modeling essential in the cybersecurity landscape, particularly as the sophistication and frequency of cyber threats continue to increase. As cyber threats evolve in complexity and impact, traditional security approaches have proven insufficient to address their dynamic nature. Threat modeling has emerged as a critical methodology that enables organizations to systematically identify, evaluate, and prioritize potential security threats before they can be exploited. By analyzing system

architectures, data flows, and potential vulnerabilities from an attacker's perspective, threat modeling creates a structured framework for security analysis that supports proactive defense strategies. This systematic approach allows security teams to allocate resources more effectively, focusing on high-risk areas and implementing appropriate countermeasures based on a thorough understanding of the threat landscape.

The evolution of threat modeling has seen a transition from isolated methodologies to more structured and systematic frameworks. Xiong and Lagerström [222] conducted a comprehensive systematic review of the literature that underscored the importance of theoretical and practical advances in the field, analyzing 176 articles and identifying three key research groups: new methodological contributions, application of existing frameworks, and the foundational literature on threat modeling processes. Their findings revealed that a significant majority of the studies employed manual modeling techniques, highlighting an urgent need for automation tools to keep pace with evolving cybersecurity threats. Graphical representations such as attack trees and fault trees have been widely adopted in the literature as valuable tools for visualizing potential threats and vulnerabilities [197]. The hierarchical structure of attack trees allows analysts to break down complex attacks into manageable components systematically, facilitating a better understanding of attack vectors and potential defenses.

The approaches to threat modeling have diversified significantly in recent years. Tatam et al. [142] propose a classification of threat modeling methodologies into four primary categories: asset-centric, system-centric, threat-centric, and data-centric approaches. Each approach offers unique perspectives on how to identify and assess threats, highlighting the importance of context in tailoring threat analysis to specific organizational needs. The asset-centric approach focuses on the identification of critical organizational assets and their associated risks, thereby improving protection strategies tailored to the value of assets [129]. In contrast, system-centric practices analyze the architecture of software systems, providing a detailed examination of potential vulnerabilities within system components. The threat-centric or attacker-centric framework shifts the focus towards understanding adversaries' motivations and tactics, advocating for an approach that anticipates potential assault mechanisms [50].

As Advanced Persistent Threats (APTs) evolve, threat modeling techniques have been refined to effectively address these new challenges. APTs are characterized by their structured, multiphase approach, requiring cybersecurity professionals to adopt a more nuanced understanding of threat dynamics [130]. The development and integration of Cyber Threat Intelligence (CTI) into threat modeling frameworks have dramatically transformed organizations' ability to foresee and respond to evolving threats. Sun et al. [73] elaborate on a structured six-step methodology to extract CTI data, highlighting its vital role in informing threat modeling and risk assessment practices. This methodology not only improves the identification of vulnerabilities but also assists in the prioritization of threat responses based on real-time threat landscapes.

Several comprehensive frameworks have emerged that provide foundational structures for effective threat modeling. Lockheed Martin's Cyber Kill Chain and MITRE's ATT&CK framework are two of the most influential tools that standardize and systematize threat modeling processes, guiding organizations to understand and counter evolving threats [221]. The Cyber Kill Chain outlines the stages of an attack, from reconnaissance to exploitation, allowing defenders to develop focused countermeasures at each stage [47]. In parallel, MITRE's ATT&CK framework offers an extensive matrix of adversarial tactics and techniques, derived from documented real-

world incidents, enabling organizations to enhance their defensive capabilities through a better understanding of adversary behaviors and methods. Integrating these frameworks into a cohesive strategy helps organizations identify critical areas for improvement and aligns their security efforts with emerging threats [53].

### 3.2.2 Computational Adversary Simulation

Despite technical fidelity, most automated adversary simulation methods ignored the social context and lacked a dynamic behavior component [117]. Human attackers have varying levels of risk tolerance, which might affect their choice of target and attack methods [258]. Human attackers can also learn from their experiences [136], dynamically adapt to defenses they encounter, and modify their strategies accordingly, making them more dangerous over time as they become more adept at evading detection and exploiting vulnerabilities. Thus, to improve the training of defenders, the emulated adversaries need to exhibit behavior similar to that of the human attackers and have the capability to learn and adapt to the defender's actions.

Early models of adversary simulation contained static patterns prescribed for the attacker agents to follow [22]. These models eventually gave way to graph-based [29] and state-based [1] attack simulation methods, which provide a useful characterization of the attacker's profile, such as goals, starting points, and available time. This group of simulation methods models and stores generic attack patterns with preconditions and postconditions in a knowledge base. Additional attack pattern attributes include the cost of attempts, execution time, base success probability, and maximum attempts. Despite their utility, these approaches often result in deterministic and predictable attack behaviors that fail to capture the complexity of human decision-making in adversarial contexts.

An important aspect of adversary behavior is bounded rationality, which refers to the limits of human decision-making capabilities. Adversaries often operate under the constraints of information and cognitive biases, leading them to make suboptimal decisions that can affect their operational success. For example, Oh et al. discuss how reinforcement learning can help optimize responses to adversarial behavior by simulating limited rationality within attack strategies [182, 183]. Similarly, understanding the limitations of adversary decision-making processes can aid in the development of better security measures tailored to exploit these bounded rationalities [6].

Moreover, it is posited that adversaries may prioritize strategies based on perceived immediate gains rather than long-term outcomes, indicative of a bounded rationality mindset [3]. This characteristic aligns with the findings of Kure et al., which highlight the need for a comprehensive understanding of various threat elements, including the cognitive limits that drive adversarial actions [134]. Thus, the modeling of bounded rationality can inform risk management frameworks by introducing adaptive strategies that focus on the psychological and decision-making processes of attackers [229].

Risk tolerance also plays a key role in adversary behavior, representing the degree to which an adversary is willing to engage with uncertainty and potential loss in the pursuit of their objectives. It varies across individuals and contexts, influencing the types and frequencies of cyberattacks launched [237, 216]. Contextual factors, such as organizational pressures and perceived rewards, shape this tolerance. For example, Thomas and Sule argue that an adversary's risk appetite can be affected by their situational context and operational goals, emphasizing the need for continuous

threat assessment [237].

Research by Shoetan et al. provides insight into how adversaries' risk tolerance levels can adjust based on their past successes or failures, indicating a cycle of risk evaluation and reassessment that influences future attack strategies [216]. Thus, evaluating risk tolerance not only assists in threat modeling, but can also be a critical factor in designing proactive defense mechanisms that anticipate potential risks adversaries are willing to undertake [262].

The learning and adaptability of adversaries are crucial in understanding their evolving tactics and techniques in cyber warfare. Adversaries continually learn from previous engagements and adapt their strategies to improve outcomes [90]. Research suggests that adversaries use machine learning to analyze the effectiveness of their attacks or adjust to countermeasures implemented by organizations [182, 183]. Oh et al. underline the importance of dynamic models, showcasing how the integration of reinforcement learning frameworks into cybersecurity can predict adversarial adaptations over time [183].

In addition, organizational culture significantly influences this learning process. Gundu outlines the necessity for a continuous learning environment, where both adversaries and defenders must engage in unlearning outdated strategies and acquiring new knowledge to remain effective [90]. This adaptability can lead to a form of arms race, in which defenders must constantly update their defenses in response to increasingly sophisticated adversarial techniques, highlighting the fluid nature of cyber-security dynamics [184].

Adversaries often rely on cognitive biases and heuristic reasoning when making decisions about attacks, which affects their effectiveness and strategy choices. These cognitive shortcuts can lead to significant miscalculations in threat assessments and decisions about attack vectors [229]. For example, the work of Loi and Christen provides evidence that biases, such as overconfidence in attack outcomes, can skew adversary planning and execution [148].

Understanding these biases enables defenders to implement countermeasures that exploit adversaries' faulty reasoning processes. Almansoori et al. emphasize that recognizing cognitive flaws can inform the design of security systems that anticipate and mitigate likely adversarial errors [6]. In this context, employing training programs for cybersecurity personnel to recognize these biases could improve preparedness against unexpected adversarial behavior [52]. This nuanced understanding of cognitive biases enriches the overall modeling of adversary behavior, offering additional information on potential decision-making errors that can create opportunities for effective countermeasures.

### 3.2.3 Cybersecurity Training with Games

While understanding adversary behavior is crucial, effectively training cybersecurity professionals to counter these threats requires innovative training approaches. Cybersecurity training utilizing games has emerged as an innovative approach to enhance awareness and skills in an age where the frequency of cyber threats is steadily increasing. The categorization of these training methods into a taxonomy begins with understanding the various types of games deployed. Serious games, often including elements designed for educational purposes, can be further divided into simulation games, digital games, and role-playing games [40]. This categorization mirrors previous findings that focus on the depth of the game and user engagement, with an emphasis on how game mechanics can support educational goals in cybersecurity [40, 225].

The consensus among researchers is that gamified training strategies positively influence learning outcomes, particularly by promoting interaction and engagement among participants [94, 191]. For example, studies recommend incorporating gamified elements such as storytelling, team challenges, and rewards to enhance the learning experience [225, 94]. Moreover, distinguishing between formal and informal learning avenues - where games serve as informal pathways - provides information on how players identify themselves with cybersecurity topics and issues [40, 94]. Therefore, training strategies must align with the cognitive and social learning theories that underpin successful gamification [14].

Noteworthy examples of game-based cybersecurity training include platforms such as picoCTF and SecGen. These platforms have successfully used game design elements to create immersive learning experiences. For example, picoCTF is a cybersecurity challenge platform that offers users a series of progressively challenging tasks structured in a narrative, making complex concepts more accessible [160, 23]. SecGen, alternatively, provides simulations that mimic the tactics, techniques, and procedures (TTPs) used by real-world cyber adversaries, allowing students to conceptualize and operationalize their knowledge in a controlled environment [191, 160].

The efficacy of these platforms has been documented in various studies, with findings indicating that participants exhibit increased knowledge retention and skill acquisition when actively engaged in game formats compared to traditional learning environments [191, 160, 114]. Furthermore, studies reveal that combining theoretical knowledge with practical challenges in a gamified context allows learners to develop critical analytical skills essential for responding to cybersecurity incidents [94, 14, 133].

Building upon the human behavior components of adversary modeling discussed previously, incorporating adversary simulation is essential to deepen the understanding of participants of cybersecurity threats and defenses. Many platforms effectively simulate adversarial behavior through game mechanics, allowing users to experience the intricate dynamics of a cybersecurity breach from both the attacker and the defender perspectives [133, 232]. The importance of adversary simulation becomes evident as simulations mirror real-world scenarios, providing players with the necessary information to preemptively devise protective measures against potential attacks [144, 114].

A prime example of adversary simulation in cyber training is the Network Defense Training Game (NDTG), which immerses players in a series of network defense scenarios against simulated attacks. This setup fosters strategic thinking and mimics a real attack environment, enabling participants to make timely decisions under pressure [14, 234, 17]. Research has highlighted that adversary simulation games improve not only the technical skills of players but also their decision-making capabilities when faced with rapid cyber threats [23, 144].

In addition, participating in adversary simulation prepares people for future scenarios in which they may need to utilize cyber threat intelligence (CTI) effectively. Studies suggest a profound impact on situational awareness when participants practice responding to simulated cyber attacks, highlighting the need to combine academic knowledge with hands-on simulations [133, 114]. Through adversary simulations, educators aim to build a resilient cybersecurity workforce capable of anticipating, detecting, and responding to a myriad of cyber threats. This dynamic is crucial in a landscape characterized by the ever-evolving nature of cybercrime.

### 3.2.4 Research Gap

Despite advances in threat modeling, computational adversary simulation, and gamified cybersecurity training, there remains a significant gap in the development of realistic adversary models that effectively capture human cognitive processes. Most existing approaches rely on deterministic and static patterns that fail to emulate the dynamic and adaptive nature of real human attackers. This limitation results in suboptimal training for cyber defenders, as they learn to counter predictable attack strategies rather than the complex, evolving tactics employed by actual adversaries.

The absence of cognitively realistic attacker models creates several challenges for effective cyber defense training. Defenders trained against deterministic attack patterns may develop false confidence in their abilities to counter real-world attacks, as static adversary emulations fail to prepare them for the adaptive and learning behaviors exhibited by human attackers. Traditional models do not account for the bounded rationality, risk tolerance, and cognitive biases that influence decision-making by humans. Furthermore, current training platforms lack adversaries that can learn from experience and dynamically adjust their strategies based on defender responses, severely limiting their effectiveness in preparing defenders for real-world scenarios.

Cognitive architectures and theories of human decision-making have made significant progress in emulating human-like behavior in dynamic environments. Unlike typical computational algorithms that aim to make optimal decisions, cognitive architectures adhere to human constraints such as forgetting, limited attention, and bounded rationality. However, past work on cognitive modeling in cyber security systems has rarely considered the real-time social interactions of attackers and defenders together. The characteristics of human attackers have been studied in specific contexts, such as phishing experiments, where individual creativity was found to be a predictor of an adversary's ability to evade detection, but these insights have not been systematically incorporated into comprehensive adversary models for training.

Therefore, there is a compelling need for cognitive models of cyber attackers that can more realistically emulate human adversary behavior, provide more challenging training scenarios for defenders, and ultimately improve cyber defense capabilities in real-world contexts. Our research addresses this gap by developing and evaluating a cognitive attacker model based on Instance-Based Learning Theory (IBLT) that captures human learning and decision-making processes, demonstrating its effectiveness as a more challenging and realistic adversary for training cyber defenders. This approach represents a significant advancement over deterministic strategies by creating adversaries that are dynamic, adaptive, and able to learn from experience, more accurately reflecting the behavior of human attackers in real-world scenarios.

### 3.3 Design of the Simulated Cyber Attack Environment

Testing attacker and defender agents requires a simulation or training platform that encapsulates cyber elements in an integrated environment. We use the interactive defense game based on CybORG AI gym with adversarial cyber-operation scenarios. Each combat between an attacker and a defender is an episode of 25 steps, ensuring sufficient time to observe attack strategies.

Figure 3.1 illustrates the topology of the network. The network is divided into three subnets:

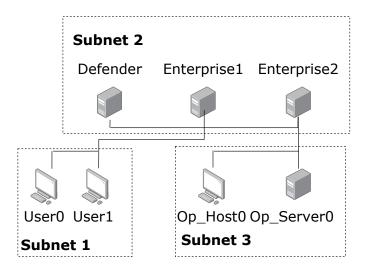


Figure 3.1: Adaptation of the Cage Challenge Network

subnet 1 consists of user hosts that are not critical, subnet 2 consists of enterprise servers supporting user activities, and subnet 3 contains the critical operational server and operational hosts. Attackers typically establish their entry point through social engineering or spear phishing, with host *User0* serving as the entry point in our scenario.

Figure 3.2 summarizes attack phases (red arrows) and defensive countermeasures (blue arrows). The Red agent starts by searching for hosts with *DiscoverRemoteSystems*, identifies vulnerabilities using *DiscoverNetworkServices*, obtains User-level access through *ExploitRemoteService*, and escalates to Root-level with *PrivilegeEscalate*. The Blue agent can *Remove* adversaries at User level, use *Restore* if the Red agent has escalated, *Analyse* activities, or passively *Monitor* the network.

### 3.3.1 Agent Types

We developed three types of red agents to test against human defenders:

 $Beeline_{Red}$  represents advanced, well-organized attackers with perfect knowledge. This agent assumes prior knowledge of network topology and moves directly to the operational server following a predetermined path  $(User0 \rightarrow User1 \rightarrow Enterprise1 \rightarrow Enterprise2 \rightarrow Op\_Server0)$  in a predictive and deterministic way.

 $Meander_{Red}$  simulates novice "script kiddie" attackers who rely on pre-made exploit programs without careful planning. This agent assumes no prior knowledge about network structure and behaves stochastically, choosing random targets to advance.

 $IBL_{Red}$  is our novel cognitive agent based on Instance-Based Learning Theory (IBLT). This dynamic agent learns from experience and adapts its actions according to the environment

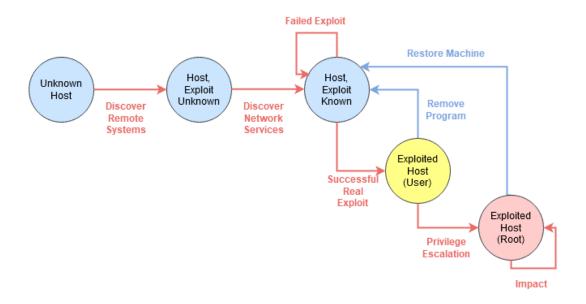


Figure 3.2: Effect of actions on the host state

conditions and defender responses. The instances that drive its decisions represent:

- State ( $s_a$ ): Features representing the attacker's knowledge of network resources in various states (Detected, Scanned, Exploited-User, Exploited-Root, Impacted)
- Action Space  $(a_a)$ : Dynamically constructed at each step based on the current network state
- Utility  $(z_a)$ : Rewards calculated based on attack success, with higher rewards for accessing significant systems

For initial training and validation, we used two simulated blue agents: (1) a passive  $Sleepy_{Blue}$  agent that only monitors the network, and (2) a dynamic cognitive  $IBL_{Blue}$  agent also based on IBLT. Our preliminary simulation experiments demonstrated that  $IBL_{Red}$  could learn effective attack strategies over time, with 55% of  $IBL_{Red}$  agents eventually outperforming the deterministic  $Beeline_{Red}$  strategy. More importantly, when facing cognitive defense agents,  $IBL_{Red}$  maintained consistently higher performance while  $Beeline_{Red}$ 's effectiveness rapidly declined as defenders learned to counter its predictable pattern.

### 3.4 Experiment: Cognitive Attacker Against Human Defenders

While our simulation results provided strong evidence that cognitive attackers present greater challenges than deterministic ones, we needed to validate these findings with human defenders. The goal of this experiment was to compare the performance of human defenders when faced with the three types of attackers ( $Beeline_{Red}$ ,  $Meander_{Red}$ , and  $IBL_{Red}^{Trained}$ ) and to assess whether humans display similar vulnerability patterns to those observed in our cognitive defender simulations.

Experimental Design Human participants completed the same cyber defense task and scenario using the Interactive Defense Game (IDG), which provides an interactive decision interface in the cyber environment. The task interface displayed the network status, observed activities on each host, and allowed participants to select hosts and defense actions (Monitor, Analyse, Remove, Restore).

Participants We recruited 186 participants (124 men, 61 women, 1 N/A) aged 21 to 65 years (M = 37.12  $\pm$  10.15) through Amazon Mechanical Turk. Approximately 9% of participants (17 individuals) reported having more than 5 years of experience in network operation and security along with at least a Master's degree in a related field. Each participant was randomly assigned to face one of the three red agents:  $Beeline_{Red}$ ,  $Meander_{Red}$ , or  $IBL_{Red}^{Trained}$ .

Procedure After providing informed consent and completing a demographic questionnaire, participants received task instructions followed by a quiz to verify their understanding. They then watched a video introduction explaining the interface, game controls, and episode dynamics.

The experiment consisted of two phases: (1) a practice session with two short episodes of 10 steps each, and (2) the main task with 7 episodes of 25 steps against the same type of adversary. The practice episodes familiarized participants with the interface and game controls using simplified scenarios. In the main task, no time restrictions were imposed, and the initial network state was identical for all participants across episodes.

After completing the main task, participants filled out a post-experiment survey about their performance, perceived strategy, and their experience in computer science and cyber defense.

#### 3.4.1 Results

Human Defenders Perform Similarly Against Both Static and Dynamic Attackers. The performance of the three attacker types against human defenders is shown in Figure 3.3. As human participants learned from experience,  $Beeline_{Red}$  agents initially performed better (M=66.25, SD=5.39) than  $IBL_{Red}^{Trained}$  agents (M=55.51, SD=4.70) in the first episode. However,  $IBL_{Red}^{Trained}$  agents posed a more persistent threat as the experiment progressed. The performance of  $Beeline_{Red}$  agents deteriorated rapidly, with  $Beeline_{Red}$  (M=47.33, SD=6.42) performing worse than  $IBL_{Red}^{Trained}$  (M=54.141, SD=6.191) in the last episode.  $Meander_{Red}$  agents (M=5.26, SD=4.42) performed significantly worse than  $IBL_{Red}^{Trained}$  across all episodes. Consistent with our simulation predictions,  $IBL_{Red}^{Trained}$  demonstrated persistent impact du-

Consistent with our simulation predictions,  $IBL_{Red}^{Trained}$  demonstrated persistent impact duration and superior performance compared to  $Meander_{Red}$  (M=0.817, SD=1.488). While  $IBL_{Red}^{Trained}$  achieved a shorter impact duration (M=1.29, SD=0.18) than  $Beeline_{Red}$  (M=3.15, SD=0.39) in the first episodes, this gap decreased significantly by the final episode.

Analysis of attack command distribution (Figure 3.4) revealed that  $IBL_{Red}^{Trained}$  used the Impact command most frequently, while  $Beeline_{Red}$  and  $Meander_{Red}$  were unable to consistently impact the operational server and resorted primarily to ExploitRemoteService and PrivilegeEscalate.

Human Defenders Employ Consistent Defense Strategies Regardless of Attacker Type. Analysis of human defensive actions (Figure 3.5) showed that participants tended to be more

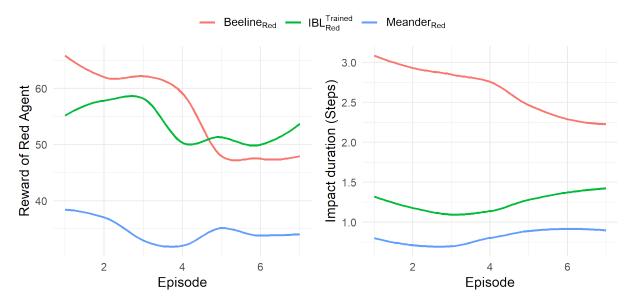


Figure 3.3: Red Agents Performance when confronted by a Human defender. The average reward per episode (left) and average impact duration (right).

passive, taking more *Analyse* and *Monitor* actions than active *Remove* and *Restore* actions. Notably, humans maintained relatively consistent action preferences throughout the experiment, although subtle adaptation occurred depending on the attacker type.

Cognitive Attackers Maintain Cognitive Load on Defenders. Figure 3.6 presents the number of options available to humans during episodes. When facing  $Beeline_{Red}$ , participants could reduce their cognitive load by narrowing the option space between the first and last episodes. In contrast, the option space remained approximately the same when facing  $Meander_{Red}$  and  $IBL_{Red}^{Trained}$ , indicating that stochastic and adaptive attackers maintain higher cognitive demands on defenders.

Dynamic Cognitive Attackers Perform Best Against The Most Efficient Human Defenders. To further investigate individual differences, we categorized participants as "Efficient Defenders" (attacker reward below mean) and "Inefficient Defenders" (attacker reward equal to or above mean). Figure 3.7 reveals that  $IBL_{Red}^{Trained}$  had significantly higher rewards against Efficient Defenders in later trials (mean: 33.19 ± 33.31; Tukey's HSD p=0.040), while  $Beeline_{Red}$  performed better against Inefficient Defenders (mean: 92.18 ± 53.36; Tukey's HSD p=0.041).

This finding is particularly striking: even skilled human defenders struggled against cognitive attackers, while less skilled defenders had more difficulty with deterministic attackers. Further analysis of defender actions showed that efficient defenders paired with  $IBL_{Red}^{Trained}$  employed significantly more active defense strategies than those paired with other attackers, demonstrating that countering cognitive attackers requires more sophisticated defensive approaches.

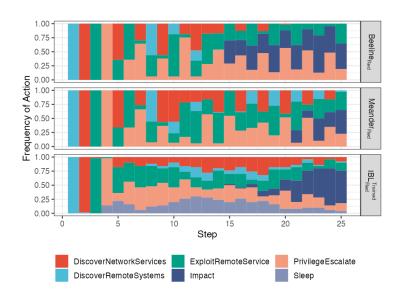


Figure 3.4: Distribution of attack commands against human defenders

### 3.5 Discussion and Conclusion

Our research demonstrates that cognitive agents emulating human-like adversaries present significantly greater challenges to cyber defenders than deterministic strategies. While human defenders effectively learned to counter predictable attack patterns, they struggled against cognitive attackers that dynamically adjusted their tactics based on defender responses. These findings have important implications for cybersecurity training and defense system evaluation.

Three key insights emerge from our research: First, training with deterministic attack patterns may inadequately prepare defenders for real-world threats from adaptive human attackers. The rapid decline in  $Beeline_{Red}$ 's performance against experienced defenders contrasts sharply with  $IBL_{Red}$ 's persistent effectiveness, suggesting that current training approaches using static patterns may create a false sense of security. Second, cognitive models based on IBLT provide a cost-effective method for producing realistic adversaries that adapt to defender actions. These models can be more effective in training cyber defense strategies than static and deterministic adversaries. The cognitive attacker agent can serve as a training partner in interactive gaming platforms, potentially addressing the scarcity of human experts for training exercises. Third, most strikingly, our finding that cognitive attackers performed best against the most skilled defenders highlights a critical vulnerability in current approaches. Even highly efficient defenders struggled against adaptive attackers, suggesting that traditional expertise may not transfer effectively to countering adaptive threats.

These results underscore the importance of training against human-like adversaries for improved cyber defense. By incorporating cognitive agents into training regimes, defenders can develop more robust strategies against the adaptive tactics employed by sophisticated human attackers.

Limitations of this work include the fact that our cognitive agent does not fully capture all aspects of human attacker behavior, such as heuristic reasoning and various cognitive biases.

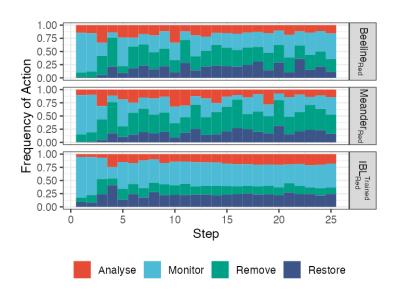


Figure 3.5: Average action frequency of Human defender

Future work will enhance the cognitive attacker agent with additional psychological mechanisms, test with expert defenders from security operation centers, and extend the approach to more complex network environments that better reflect real-world conditions.

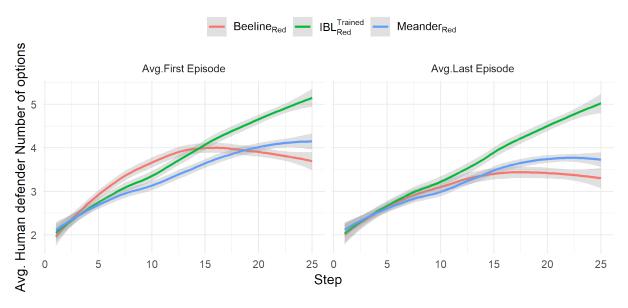


Figure 3.6: Average size of the Human defender's option space in the first (left) and the last episode (right)

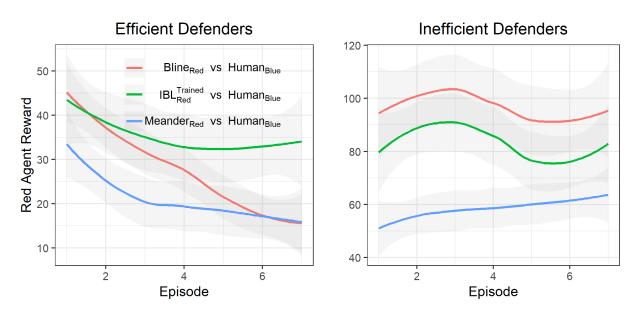


Figure 3.7: Average red agent reward by episode, split between efficient and inefficient human defenders.

## Chapter 4

# Multi-Agent Decision-Making in Cybersecurity: Cognitive Mechanisms for Multi-Defender Interaction<sup>1</sup>

### 4.1 Introduction

Cybersecurity challenges increasingly extend beyond individual defenders to complex networks of interdependent actors. Organizations must collaborate to share threat intelligence while simultaneously managing competitive risks [162]. As cyber threats grow in sophistication, no single organization can maintain comprehensive threat awareness independently, making information sharing crucial for collective defense [241]. However, this creates a strategic dilemma—sharing valuable threat intelligence benefits collective security but may expose the sharing organization to reputational damage, legal liabilities, or competitive disadvantage [71].

While research has extensively studied bilateral information sharing relationships, real-world cybersecurity contexts involve multiple simultaneous relationships under cognitive constraints. Security Operation Centers (SOCs) maintain partnerships with multiple organizations, government agencies, and Information Sharing and Analysis Centers (ISACs), each with different trust levels, priorities, and information needs [98]. Analysts must decide which information to share with each partner, considering both immediate strategic implications and long-term relationship development.

Current approaches to modeling behavior in strategic group interactions typically fall into two broad categories. Evolutionary approaches [208, 193] effectively capture population-level outcomes, but often abstract away individual cognitive processes. Game-theoretic frameworks [70, 228] provide precise mathematical formulations but rely on strong rationality assumptions that rarely match human behavior. Between these approaches, cognitive modeling has emerged as a promising third path that explicitly addresses how humans navigate social learning under inherent limitations.

This chapter investigates how individuals make strategic decisions in multi-defender cybersecurity environments through two complementary studies. First, I examine how incentive

<sup>&</sup>lt;sup>1</sup>See Appendix .4 for the Multi-Defender Game and Appendix .5 for the Computational Model

structures and information availability influence cooperation in triadic cybersecurity information sharing contexts. By systematically varying both structural incentives and information participants receive about their interdependencies, I reveal important patterns in how cooperation emerges and stabilizes in small groups.

Building on these empirical findings, I then develop a broader cognitive theory that explains how individuals navigate multiple relationships simultaneously despite cognitive limitations. This theory integrates three key psychological mechanisms: dynamic weighting of others' outcomes, categorical learning for efficient relationship management, and contrast effects for comparative evaluation of relationship partners. Together, these mechanisms enable individuals to process multiple relationships efficiently while maintaining strategic flexibility.

This progression from specific cybersecurity experiments to general cognitive theory provides crucial insights into both practical cybersecurity cooperation challenges and the fundamental cognitive processes that underlie strategic social interaction. By understanding these mechanisms, we can design more effective information sharing frameworks that align with natural cognitive processes, ultimately improving collective cybersecurity outcomes.

### 4.2 Related work

### 4.2.1 Emergence of cooperation in groups

Cooperation within groups larger than two players has been widely studied in multiple disciplines, including economics, social psychology, and evolutionary biology. This broader exploration of group cooperation extends beyond the classic Prisoner's Dilemma (PD) and includes diverse game-theoretic and real-world scenarios. For instance, information sharing in cybersecurity resembles a Prisoner's Dilemma, as entities must decide whether to share valuable threat information (cooperate) or withhold it for their advantage (defect), facing risks if others do not cooperate. In game theory, group cooperation is often studied through public goods games and collective action problems, which introduce complexities not present in dyadic interactions. Public goods games, for example, examine how individuals contribute to a shared resource pool while dealing with incentives to ride free [64, 208]. Recent research highlights factors such as group size, communication, and mechanisms for punishing free riders as critical to maintaining cooperation [188, 233]. These studies show that group cooperation is more complex than dyadic interactions due to greater mutual dependencies and coordination challenges.

The Prisoner's Dilemma has also been extended to multiplayer versions, sometimes referred to as the N-person Prisoner's Dilemma. In these settings, individuals must decide whether to cooperate with the entire group rather than just one other player. Research suggests that cooperation in such multiplayer PD scenarios is influenced by factors such as reciprocity, reputation, and social norms [179, 31]. However, multi-player dilemmas introduce unique challenges, including coordination issues and an increased impact of individual decisions on group outcomes [86, 26].

In contrast to the N-person Prisoner's Dilemma, our experimental design retains the fundamental 2-by-2 PD interactions within a triad context. This allows us to focus on pairwise decision-making while examining how these interactions aggregate to influence the group. The advantage of this design is that it enables us to capture both dyadic and triadic dynamics, providing

insights into how individual relationships affect broader group behavior. Unlike typical N-person PD scenarios, where cooperation is assessed at the collective level, our approach provides a detailed examination of the interplay between dyads within the group, revealing the conditions under which cooperation is stabilized or disrupted from individual behavior.

To understand how cooperation emerges and is sustained in groups, it is crucial to consider both structural and informational factors that influence decision-making. Here, we focus on two key elements: incentive structures and information levels on mutual interdependence.

### 4.2.2 Incentive Structure and Social Preferences

The interplay between incentive structures and social preferences significantly shapes cooperative behavior in group settings. Incentive structures, defined by the potential rewards or costs players face based on their decisions to cooperate or defect, directly influence individual motivations. Moisan et al. [168] demonstrated that as players' cooperativeness increases, there is a sharp transition from defection to cooperation, with the transition point depending on the game's payoff matrix. Their work showed that inequality aversion among players promotes cooperation by transforming perceived incentives.

A well-established measure of expected cooperation in PD games with the payoff matrix is shown in Table. 4.1 is Rapoport's K-index [203], defined as K = (R-P)/(T-S), where R represents the reward for mutual cooperation, P the punishment for mutual defection, T the temptation payoff for unilateral defection, and S the sucker payoff for unilateral cooperation. The K-index captures the expected cooperation by considering how much players benefit from defecting versus the cost of mutual defection. When K is high (i.e. when T is not much larger than S or P is numerically large), defection is less rewarding, and mutual defection is more costly, making cooperation more likely.

	P2: Cooperate	P2: Defect
P1: Cooperate	(R,R)	(S,T)
P1: Defect	(T,S)	(P,P)

Table 4.1: Payoff matrix for the Prisoner's Dilemma.

Prosociality [163] (e.g., Social Value Orientation (SVO)) adds another layer to this dynamic by reflecting how individuals weigh their results against others. SVO can be represented through a utility function  $u(\pi_{\text{self}}, \pi_{\text{opponent}}) = u_{\text{self}} + \alpha \cdot u_{\text{opponent}}$ , where  $\alpha$  represents the weight given to the opponent's payoff. For any PD game, there exists a threshold  $\bar{\alpha}$  such that players with  $\alpha > \bar{\alpha}$  will prefer cooperation regardless of their beliefs about the behavior of others, while those with  $\alpha < \bar{\alpha}$  will consistently choose defection.

## 4.2.3 Information Levels and Decision Making

The effectiveness of incentive structures in promoting cooperation is highly dependent on the information available to players about their mutual interdependence [247]. The Hierarchy of

Social Information (HSI) framework proposed by Gonzalez and Martin [82] conceptualizes three main levels of interpersonal information. At the Minimal Information level, players know that they are interdependent, but lack details about how their actions affect others. The Experiential Information level allows players to observe others' actions and outcomes over time, enabling learning through experience about their interdependencies. The Descriptive Information level provides complete information about the payoff structure upfront, in addition to experiential feedback.

This framework suggests that providing more detailed information about interaction structures can foster cooperation more effectively than limited or no social information. Gonzalez et al. [83] found that continued visibility of the payoff matrix helps clarify the trade-off between short-term and long-term rewards, while experiential feedback strengthens the understanding of reciprocal relationships.

The combination of incentive structures and information levels creates a complex decision environment where players must balance individual and collective interests. These factors have been particularly relevant in cybersecurity contexts, where organizations must decide whether to share threat information. Research has shown that rewarding and punishing certain actions can significantly affect information-sharing behavior [241]. Similarly, studies on cybersecurity information exchange have shown that clarity of feedback on interdependencies influences cooperation rates [71].

The study of repeated strategic interactions between interdependent agents has a rich research history. Early work focused on simple strategies with minimal partner modeling. Axelrod's seminal computer tournaments of the Iterated Prisoner's Dilemma [15] demonstrated the success of Tit-for-Tat (TFT), which only considers the partner's last action. Although these simple strategies proved to be remarkably effective in structured environments, subsequent research revealed their limitations in noisy or complex settings [178]. This led to increasingly sophisticated approaches incorporating richer agent modeling and learning mechanisms, culminating in modern machine learning methods and cognitively inspired strategies. Two fundamental challenges have emerged in this progression: the computational demands associated with memory and learning and the complexity of modeling diverse agent strategies. Our work addresses these challenges by incorporating cognitive mechanisms for efficient memory use and agent categorization.

# 4.2.4 Learning and Agent Modeling in Interdependent Interactions

A significant theoretical advancement in agent modeling came with [193]'s discovery of zero-determinant strategies [193], which established mathematical boundaries on strategy effectiveness and enabled unilateral control over payoff relationships. This discovery fundamentally changed our understanding of what was possible in repeated interactions, showing that agents could enforce linear payoff relationships regardless of their partners' actions. [228] extended this work by demonstrating that "generous" variants often outperform purely extortionate strategies in evolutionary settings, highlighting how successful strategies must balance exploitation with mutual benefit.

The development of agent modeling approaches has followed several trajectories. Early work focused on explicit prediction of others' actions through pattern recognition [32], while later approaches incorporated uncertainty and partial observability [77]. Modern machine learning

methods, particularly deep reinforcement learning, have demonstrated impressive success in learning implicit representations of agent behavior [101, 149]. These approaches can uncover sophisticated counterstrategies through extensive self-play and experience accumulation, often exceeding hand-crafted strategies in complex environments.

However, the increasing sophistication of learning algorithms has led to an "arms race" in strategy complexity. Neural network-based approaches can learn highly non-linear decision boundaries [140], allowing more nuanced responses, but also making strategies harder to interpret and analyze. This complexity creates challenges for theoretical analysis and raises questions about the robustness of the learned strategies. Some studies suggest that simpler strategies with clear theoretical foundations may be more robust across diverse interaction partners [249].

The tension between strategy complexity and robustness has motivated research into hybrid approaches that combine machine learning with domain knowledge. For example, [43] demonstrated how the incorporation of simple mechanisms that promote mutual benefit into learning algorithms can improve generalization between different interaction partners. Similarly, [123] showed that learning algorithms constrained by the principles of game theory often develop more stable and interpretable strategies.

Recent work has increasingly focused on multi-agent scenarios in which agents must simultaneously model and adapt to multiple partners [135]. This setting introduces additional complexities, as agents must balance their responses between different partners while maintaining coherent strategies. The challenge is compounded in settings with incomplete information or when partners may change their strategies over time [105].

### 4.2.5 Memory Constraints and Cognitive Plausibility

While machine learning approaches have demonstrated impressive performance in agent modeling, they typically assume unlimited memory capacity and computational resources. These approaches often maintain complete interaction histories or complex state representations, enabling sophisticated pattern recognition but diverging significantly from human cognitive constraints. This disconnect raises important questions about the psychological plausibility and practical applicability of such models.

Empirical studies reveal clear limitations in human memory use during strategic interactions. Research consistently shows that humans typically access only 5-10 previous interactions when making decisions [170], indicating a clear cognitive bottleneck. This limitation reflects broader constraints on working memory capacity, which affects how individuals process and utilize information in dynamic social situations. Memory traces follow systematic decay patterns [10], with recent interactions more heavily weighted while maintaining the diminishing influence of established patterns, a phenomenon known as the power law of forgetting.

The relationship between memory complexity and strategy performance follows an inverted U-shaped pattern [106], suggesting optimal performance at intermediate levels of memory complexity. This finding has profound implications for the design of strategies. Although too little memory prevents recognition of important behavioral patterns, excessive memory complexity can lead to overfitting and reduced adaptability. This balance reflects the fundamental principles of bounded rationality [219], where cognitive constraints paradoxically contribute to more robust and adaptable decision-making.

Recent work has attempted to bridge this gap between machine learning approaches and cognitive constraints. For example, [227] demonstrated how memory-constrained models can achieve performance comparable to that of more complex approaches by focusing on relevant features and efficient information encoding. Similarly, [29] showed that the incorporation of human-like memory decay mechanisms can improve the model predictions of actual behavior in repeated games.

These findings suggest that effective strategies should not simply operate within memory constraints, but actively leverage them as design principles. Memory limitations can serve as natural regularizers, promoting generalization by preventing overfitting to specific interaction patterns. This perspective aligns with the ecological rationality frameworks [240], which emphasize how cognitive constraints can improve decision-making in natural environments.

## 4.2.6 Categorical Learning and Contrast Effects

A fundamental challenge in social dilemmas is the wide space of possible peer strategies. As the diversity of peers increases, the complexity of the modeling increases exponentially [143], making the modeling of direct strategies computationally intractable. This challenge becomes particularly acute in multi-agent settings where traditional modeling approaches often fail to scale effectively or require unrealistic computational resources.

Humans address this complexity through sophisticated categorical learning mechanisms that enable efficient but flexible social learning. Research shows that people actively form and update categories based on patterns of interdependence in social interactions [156]. These categories serve not just as simplifying heuristics, but as predictive models that guide future cooperation decisions. For example, when individuals identify patterns of reciprocity or exploitation, they develop categorical representations that help them anticipate and respond to similar behaviors in new interactions [124].

What makes categorical learning particularly powerful is its ability to balance efficiency with effectiveness. Although categorization reduces the granularity of social information, it paradoxically allows more sophisticated responses by capturing essential behavioral patterns [36]. People continually refine these categories based on new experiences, maintaining a dynamic equilibrium between stable categorical knowledge and adaptability to novel patterns. This process of category refinement is strongly influenced by the social context: Individuals' classifications of "cooperative" versus "non-cooperative" behavior emerge relative to their broader social experience [252].

These categorical learning mechanisms have been demonstrated in various social dilemmas. In cybersecurity information-sharing networks, Mermoud et al. [162] found that defenders naturally categorize their peers into "regular sharers" and "free-riders" based on sharing patterns, using these categories to guide their own sharing decisions even with new peers. Similarly, in organizational contexts, studies of team-based resource allocation show that managers develop categorical representations of "reciprocators" versus "opportunists" that influence future resource-sharing decisions [99].

The power of categorical learning is particularly evident in repeated interaction settings. For example, in public goods games, participants rapidly develop categories for "consistent contributors" and "strategic free-riders," with these categories shaping not only direct responses but also reputation sharing within groups [65]. These categories prove to be remarkably stable -

once an individual is categorized as a reliable cooperator, isolated defections are often discounted as anomalies rather than prompting immediate category reassignment [12].

Experimental studies of group cooperation reveal how categorical learning enables efficient decision-making under time pressure. When faced with multiple potential cooperation peers, participants don't track detailed histories but instead maintain broader categorical assessments like "trustworthy," "unpredictable," or "exploitative" [118]. These categorical judgments are particularly influential in early interactions with new peers, where they serve as default expectations until individual-specific evidence accumulates [256].

The categorical perception of peers introduces systematic contrast effects in behavior evaluation. Rather than evaluating each pee's actions in isolation, individuals evaluate behaviors relative to their experiences with other peers [255]. These contrast effects are particularly pronounced between categorically distinct peers. For instance, [122] demonstrated that players' responses to moderately cooperative behavior become more positive when they simultaneously interact with clearly non-cooperative peers, suggesting that categorical boundaries enhance behavioral discrimination.

The sophistication of categorical human learning extends beyond simple classification. Successful players develop hierarchical category structures, with broad behavioral types that contain subtypes that capture more nuanced patterns [202]. This hierarchical organization allows players to balance computational efficiency with strategic sophistication. Moreover, these learned categories are effectively transferred between different economic games [189], suggesting that categorical learning captures fundamental aspects of strategic behavior.

### Cognitive Approaches to Social Learning

Category learning represents a fundamental cognitive mechanism that influences how individuals perceive, process, and retain information about social interactions. [109] describe how attentional mechanisms significantly impact social perception through category accentuation, where individuals exaggerate differences between groups while minimizing within-group variations. [214] further demonstrate that such cognitive biases can enhance memory for features associated with majority groups while diminishing recall for minority group characteristics, highlighting how categorization processes can systematically shape learning outcomes.

The contrast effect, a key phenomenon in category-based perception, is significantly influenced by the psychological distance between learners and their interaction partners. Research indicates that as psychological distance increases, learners tend to focus more on abstract goals rather than concrete behaviors [116, 100]. [72] found that psychological proximity promotes more faithful emulation of specific actions, whereas distance encourages goal-oriented imitation. This relationship between distance and learning style suggests that the positioning of individuals relative to their interaction partners fundamentally shapes how they process and internalize social information.

Social dynamics further influence cognitive development through the cultural framework of interactions. [107] argues that social learning facilitates cultural inheritance, highlighting the intersection between associative learning mechanisms and cognitive evolution within social contexts. [194] emphasize how internalization processes are shaped by cultural frameworks that inform individual cognitive practices in social interactions. Furthermore, [243] underscore

the importance of recognizing varying contexts in shaping cognitive outcomes, including how individuals perceive frequency and value in social behaviors.

Multiple theoretical frameworks have been proposed to explain these phenomena of social learning. Reinforcement learning models [63] focus on outcome-based behavioral adjustments but often struggle with the dynamic nature of social environments. Bayesian approaches [16] represent uncertainty through probabilistic beliefs about others' intentions, but frequently assume unrealistic inferential capabilities. Theoretical frameworks of the mind [254] emphasize meta-representational abilities but may overestimate typical cognitive capacities in complex scenarios. Heuristic approaches [75] propose that simple decision rules can achieve effective social coordination despite limited information processing.

Among these various frameworks, the Instance-Based Learning Theory (IBLT) [81] offers a particularly compelling account of social learning under cognitive constraints. Unlike approaches that either oversimplify cognitive processes or assume unrealistic computational capabilities, IBLT provides a psychologically grounded explanation for how individuals learn from specific experiences while respecting memory limitations. Through mechanisms like activation decay and similarity-based retrieval, IBLT naturally explains how categorical thinking emerges from interactive experiences. Memory constraints guide attention toward meaningful patterns rather than exhaustive details, leading to more robust and generalizable learning [106].

We built on the [83] model of dyadic interdependence because it offers a formal mechanism (through the  $\alpha$  parameter) for representing how individuals incorporate others' outcomes into their decision-making. This instance-based architecture provides a suitable foundation for implementing categorical learning while maintaining cognitive plausibility. By extending this model to incorporate categorical learning and contrast effects, we address the fundamental challenge of managing multiple relationships within realistic cognitive constraints, enabling more efficient processing of multiple social interactions.

# 4.3 Emergent Cooperative Decision-making in Triadic Prisoner's Dilemma: Effects of Incentives and Information

### 4.3.1 Experimental Design: The Multi-Defender Game

The Multi-Defender Game (MDG) was developed to study cooperation in three-person groups through the lens of cybersecurity information sharing. This experimental paradigm simulates a scenario where three defenders must decide whether to share threat intelligence with each peer separately, creating a network of pairwise interactions that form a triad. The game incorporates both immediate and long-term incentives: sharing information costs the sender 15 points but provides the receiver with 35 points, while also reducing the receiver's probability of being breached in subsequent rounds according to the formula  $\Pr_{t+1} = \Pr_t - (0.95 \cdot Z_i^t/2000)$ , where  $Z_i^t$  represents the sharing points of defender i in trial t. This structure creates a Prisoner's Dilemma within each relationship, where mutual sharing provides net positive outcomes (+20 points each), but unilateral defection offers higher immediate rewards (+35 points) at the expense of the cooperating partner (-15 points). The task ran for 50 rounds, with participants making

independent decisions in each round about whether to share with each of their two peers, allowing for the emergence of complex relationship patterns and strategic adaptations over time.

The experiment systematically manipulated two key factors hypothesized to influence cooperation. First, structural incentives were varied through the K-index, a theoretical predictor of cooperation defined as K = (R - P)/(T - S), where R represents mutual cooperation reward, P represents mutual defection punishment, T represents temptation payoff, and S represents sucker payoff. Groups were assigned to either a low incentive condition (K = 0.4) or high incentive condition (K = 0.8), with higher K-index values theoretically promoting cooperation by reducing the temptation to defect relative to the cost of mutual defection. Second, information availability was manipulated across three levels based on the Hierarchy of Social Information framework: Minimal Information (participants received only basic feedback on whether peers shared information), Experiential Information (participants observed detailed outcomes of interactions with each defender through a feedback table), and Descriptive Information (participants received the complete payoff matrix in addition to experiential feedback). This created a 2×3 experimental design with four key conditions analyzed: Game I (K = 0.4, Minimal Information), Game II (K = 0.8, Minimal Information), Game III (K = 0.4, Experiential Information), and Game IV (K = 0.4, Descriptive Information). A total of 519 participants (173 groups of 3 individuals) were recruited from Amazon Mechanical Turk and randomly assigned to these conditions.

### 4.3.2 The Surprising Effect of Information on Cooperation

Contrary to theoretical predictions, the results revealed a surprising pattern in how information availability affected cooperation rates. While the K-index influenced cooperation in the expected direction (higher cooperation under K=0.8 than K=0.4), the relationship between information and cooperation did not follow the anticipated progression. Instead, an inverse U-shaped relationship emerged where Experiential Information (Game III) produced the highest cooperation rates (71.7%), followed by Minimal Information with high K-index (Game II, 68.1%), Minimal Information with low K-index (Game I, 57.2%), and finally Descriptive Information (Game IV, 49.6%). This pattern directly contradicts the hierarchical prediction that more complete information should lead to greater cooperation. Two-way ANOVA confirmed significant main effects for both K-index [F(1, 19) = 67.656, p; 0.001] and information level [F(2, 30) = 87.179, p; 0.001]. The finding that descriptive information actually reduced cooperation compared to minimal information represents a particularly counterintuitive result that challenges fundamental assumptions about how information influences strategic decision-making in group contexts.

Analysis of decision patterns revealed distinct behavioral tendencies that help explain these surprising results. When provided with descriptive information (Game IV), participants focused more on short-term strategic calculations, becoming highly attuned to the immediate temptation payoff visible in the matrix. This heightened awareness of potential exploitation prompted more defensive behaviors and pre-emptive defection. In contrast, participants with experiential information (Game III) developed stronger reciprocity norms through direct observation of outcomes, learning the value of sustained cooperation through experience rather than abstract payoff descriptions. This learning process is evident in the conditional probabilities of cooperation: participants in Game III were more likely to reciprocate cooperation (72.1%) than to retaliate against defection (68.1%), while the opposite pattern emerged in all other conditions. The re-



Figure 4.1: Initial game interface showing player status, including available points (starting endowment: 1000 points), probability of attack, and attack status for the current round. Players choose whether to share information using Yes/No buttons.

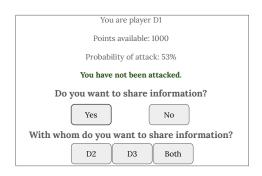


Figure 4.2: Information sharing selection interface, allowing players to choose specific recipients (D2, D3) or share with all defenders.



Figure 4.3: End-of-round feedback screen showing attack status, information sharing outcomes, and updated statistics including new points available and adjusted probability of attack for the next round.

Figure 4.4: Interfaces used during the game: (a) Initial game interface, (b) Information sharing interface, and (c) End-of-round feedback screen.

#### A. Minimal Information Level

Game Status Update:

- Defender 1 shared information with me
- Defender 2 didn't share any information

# **B.** Experiential Information Level

My Actions	Defender 1	Defender 2
<ul> <li>Did not share with Defender 1</li> <li>Shared with Defender 2</li> </ul>	<ul> <li>Shared information</li> <li>Was attacked</li> <li>My gain: +35</li> <li>Their cost: -15</li> </ul>	<ul> <li>Did not share</li> <li>No points exchanged</li> <li>My gain: 0</li> <li>Their gain: 0</li> </ul>

### C. Descriptive Information Level

Additional to experiential information, players see the payoff matrix: (left: K=0.4; right: K=0.8)

	Share	Don't Share		Share	Don't Share
Share	(+20, +20)	(-15, +35)	Share	(+30, +30)	(-15, +35)
Don't Share	(+35, -15)	(0, 0)	Don't Share	(+35, -15)	(-10, -10)

Figure 4.5: Three levels of information provided to players. (A) Minimal Information provides only basic sharing status. (B) Experiential Information shows detailed outcomes of interactions with each defender. (C) Descriptive Information adds the complete payoff matrix to help players understand potential outcomes.

gression analysis further confirmed that receiving information from peers in the preceding round increased the likelihood of cooperation by 80-87%, highlighting the importance of direct reciprocity in driving cooperative behavior. These findings suggest that experiential learning creates more robust cooperative tendencies than explicit strategic information, particularly in complex multi-agent environments where participants must manage multiple relationships simultaneously.

### 4.3.3 Emergence of Selective Cooperation

The evolution of cooperation patterns over time revealed additional insights into the learning dynamics across conditions. In all games, cooperation initially dropped sharply during the first ten rounds, a pattern commonly observed in repeated Prisoner's Dilemma studies. However,

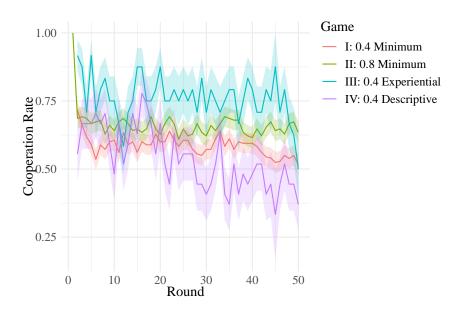


Figure 4.6: Average Individual Cooperation Rate Over Time across experimental conditions. The cooperation rate dropped sharply during the first ten rounds in all conditions, but continued to decline notably in Game IV while remaining relatively stable in other conditions.

the subsequent trajectories diverged significantly: cooperation continued to decline in Game IV (Descriptive Information) while stabilizing in the other conditions. This pattern suggests that descriptive information not only initially suppressed cooperation but also inhibited the learning processes that typically allow cooperation to re-emerge through experience. Analysis of initial strategies revealed that most participants (66%) began with a prosocial approach, sharing with both groupmates, while 33

The most distinctive feature of triadic interactions emerged in the analysis of how participants managed multiple relationships simultaneously. Over time, a clear shift occurred from universal cooperation (sharing with both peers) toward selective cooperation (sharing with only one peer). Repeated measures ANOVA confirmed significant effects of both condition  $[F(3, 515) = 2.148, p \mid 0.001]$  and round  $[F(49, 25235) = 5.148, p \mid 0.001]$  on sharing patterns. While participants initially tended to treat both peers similarly, they increasingly differentiated between them as experience accumulated, developing stronger cooperation with one peer at the expense of the other. This selective cooperation strategy was particularly pronounced in the Experiential Information condition, where participants received clear feedback about each peer's behavior. Analysis of sequential dependencies revealed sophisticated conditional strategies: after observing divergent behaviors from their peers (one cooperating, one defecting), participants in Game III often attempted to restore mutual cooperation, while those in Game IV more frequently matched their peers' previous actions, reciprocating cooperation with cooperators and defection with defectors. These patterns highlight how information conditions shape not only overall cooperation rates but also the specific strategies participants employ to navigate multiple relationships.

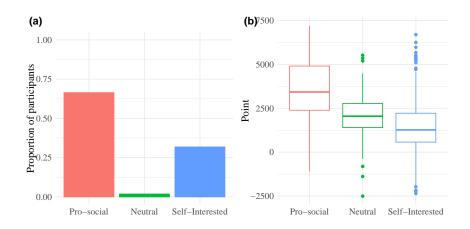


Figure 4.7: Participant Strategies and Performance. (a) The proportion of participants categorized by their initial sharing strategy: Pro-social (shared with both groupmates), Neutral (shared with one groupmate), or Self-Interested (shared with no one). (b) Cumulative points earned by participants at the end of the game, categorized by their initial strategy.

### 4.3.4 The Third-Player Effect: How Triads Differ from Dyads

Perhaps the most fascinating finding emerged in the analysis of third-player effects on dyadic relationships. When a third player adopted a selective cooperation strategy (cooperating with one peer but not the other), this had remarkably different effects depending on the existing state of the relationship being influenced. For dyads already engaged in mutual cooperation, selective cooperation by the third player strengthened and stabilized the relationship, increasing cooperation levels beyond what was observed with universal cooperation or defection strategies. For dyads engaged in mutual defection, selective cooperation by the third player similarly helped to break the defection cycle and promote cooperation, though less dramatically. However, for asymmetric relationships (one cooperating, one defecting), selective cooperation by the third player frequently destabilized the relationship, often pushing it toward mutual defection rather than mutual cooperation. This complex mediation effect reveals how third-party behaviors create feedback loops that can either reinforce or undermine cooperation between pairs, demonstrating that triadic interactions cannot be understood as simple aggregations of independent dyadic relationships. These third-player effects were strongest when structural incentives were high (K = 0.8) and participants received experiential information, highlighting the interactive nature of the experimental factors in shaping group dynamics.

Further analysis of relationship dynamics within triads revealed how initially balanced or imbalanced relationships evolved over time. When all members began with cooperation, relationships tended to stabilize with minimal differences in strength between pairs. However, mixed initial strategies often led to growing disparities between pairs, creating persistent imbalances in the triad. Even when one member cooperated while two defected, overall cooperation tended to increase over time, but with noticeable gaps between pair relationships—the cooperative member typically formed stronger bonds with whichever peer reciprocated first. These findings align with social balance theory but extend it by demonstrating how asymmetric strategies create lasting imbalances that resist equilibration, particularly in the context of information sharing decisions

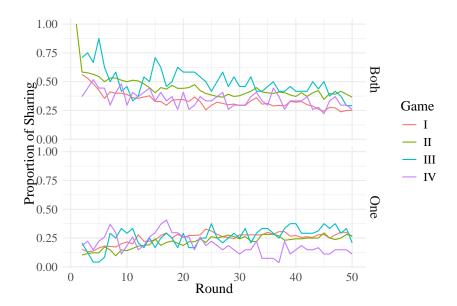


Figure 4.8: Evolution of sharing preferences over time. Participants gradually shifted from sharing with both peers to selective sharing with one peer, particularly during the first 25 trials.

with both short-term and long-term consequences.

### 4.3.5 Implications and Conclusions

The findings from this study have important implications for understanding cooperation in multiagent systems and specifically for improving cybersecurity information sharing frameworks. First, the surprising effect of information availability suggests that descriptive approaches emphasizing theoretical benefits and payoff structures may be less effective than experiential learning mechanisms that allow participants to observe the concrete outcomes of cooperation. Information sharing platforms should therefore prioritize clear feedback about successful threat mitigations resulting from shared intelligence over abstract descriptions of potential benefits. Second, the evolution toward selective cooperation indicates that tiered sharing frameworks may be more sustainable than all-or-nothing approaches, as they allow organizations to maintain different levels of information exchange with different partners based on reciprocity and trust. Finally, the third-player effects highlight how group composition influences information sharing dynamics: selectively cooperative organizations can either stabilize or destabilize existing sharing relationships depending on initial conditions, suggesting that careful attention should be paid to group formation in information sharing communities. Collectively, these insights provide a more nuanced understanding of how cooperative behavior emerges from the complex interplay of incentive structures, information availability, and social dynamics in multi-agent contexts.

The triadic framework developed in this research represents a significant advancement in the study of cooperation, moving beyond traditional dyadic analyses while maintaining analytical tractability. By systematically varying both structural incentives and information availability within the same experimental paradigm, this study provides a comprehensive examination of how these factors interact to shape cooperative behavior in small groups. The findings challenge

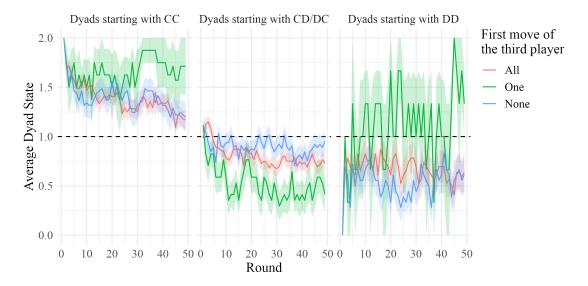


Figure 4.9: Influence of the Third Player on Dyadic Cooperation. The y-axis shows the average mutual sharing behavior from mutual defection (0) to mutual cooperation (2). Selective cooperation by the third player ("One") had the strongest positive effect on initially cooperative or defective pairs, but could destabilize mixed-strategy pairs.

prevailing assumptions about the relationship between information and cooperation, reveal sophisticated strategic adaptations to multi-agent environments, and highlight the emergent properties that arise when individuals must simultaneously manage multiple interdependent relationships. Future research should examine how these dynamics scale to larger networks, investigate the cognitive mechanisms that enable the concurrent management of multiple relationships, and test the practical applications of these findings in real-world information sharing contexts. By advancing our understanding of cooperation in small groups, this research contributes to both the theoretical foundations of strategic interaction and the practical design of systems that facilitate collective action in contexts where individual and group interests may conflict.

# 4.4 Toward a Cognitive Theory of Interdependent Decisions in Groups: Dynamic Weighting, Categorization, and Contrast

The experimental findings on triadic cooperation raised a fundamental question: How do individuals cognitively process and manage multiple strategic relationships simultaneously given limited cognitive resources? This question extends beyond cybersecurity to any domain characterized by strategic interdependence among multiple agents.

Human social systems are fundamentally characterized by strategic interdependence, where individual decisions interact with collective outcomes. In real-world scenarios, individuals must track, evaluate, and respond to multiple partners while operating with limited cognitive resources [164, 227]. The cognitive mechanisms that enable effective management of multiple relationships despite these constraints remain poorly understood.

Empirical evidence suggests that people do not maintain complete models of each interaction partner. Studies consistently show that humans typically access only 5-10 previous interactions when making decisions [170], operating under clear cognitive constraints. Yet somehow, individuals manage to navigate complex social environments effectively, suggesting the existence of efficient cognitive mechanisms for processing multiple relationships.

### 4.4.1 Theoretical Framework

Building on the experimental results, I developed a broader theoretical framework that integrates three key cognitive mechanisms to explain how people navigate complex social environments:

### **Dynamic Weighting**

The first mechanism involves dynamically adjusting how much individuals value others' outcomes based on expectation-reality discrepancies. This builds on previous work by Gonzalez et al. [83], who incorporated Social Value Orientation into Instance-Based Learning Theory through a weighted additive rule:

$$V_k = \sum_{i=1}^n p_{ik}(x_{self} + \alpha \cdot x_{other})$$
(4.1)

Where  $x_{self}$  and  $x_{other}$  are the values of the player's outcome and the peer's outcome, respectively, in instance i associated with option k;  $\alpha$  represents the extent to which a player considers others' outcomes when making choices;  $p_{ik}$  is the probability of retrieving instance i associated with alternative k from memory.

The original formulation updated  $\alpha$  based on the absolute gap between expected and actual outcomes:

$$Gap(t) = Abs(V_k(t-1) - (x_{self} + \alpha(t)x_{other}))$$
(4.2)

$$\alpha(t+1) \leftarrow (1-\eta)\alpha(t) + \eta(1 - \hat{Gap}(t)) \tag{4.3}$$

where  $\hat{Gap}(t) \in [0, 1]$  is the normalized Gap(t), and  $\eta$  is the learning rate.

However, this formulation doesn't distinguish between positive surprises (actual outcome exceeds expectations) and negative surprises (actual outcome falls short). My refined formulation addresses this limitation:

$$Gap(t) = V_k(t-1) - (x_{self} + \alpha(t)x_{other})$$
(4.4)

$$\alpha(t+1) \leftarrow \begin{cases} (1-\eta)\alpha(t) + \eta \max(\alpha(t), \hat{Gap}(t)), & \text{if } Gap(t) \ge 0\\ (1-\eta)\alpha(t) + \eta \hat{Gap}(t), & \text{if } Gap(t) < 0 \end{cases}$$
(4.5)

This asymmetric update rule increases  $\alpha$  when peers exceed expectations and decreases it when they disappoint, allowing selective cooperation based on relationship history.

Figure 4.11 shows how this refined formulation enables differentiation between cooperative and defective peers, unlike the original formulation which results in similar  $\alpha$  values for both types.

### **Category Learning**

The second mechanism addresses how people efficiently organize social experiences into behavioral prototypes. Rather than tracking each relationship individually, people categorize peers based on behavioral patterns, allowing efficient generalization across relationships.

The model incorporates five behavioral dimensions to characterize interaction partners: (1) Action tendency: Proportion of cooperative actions, (2) Entropy: Unpredictability in action sequences, (3) Responsiveness: Correlation between current action and partner's previous action, (4) Recovery propensity: Rate of return to cooperation after defection, (5) Volatility: Frequency of strategy changes.

These dimensions allow sophisticated but cognitively manageable social categorization. The model constructs and maintains a hierarchical category structure through an iterative clustering process described in Algorithm 2.

```
Function ComputeFingerprint(sequence):
   return [action_tendency, entropy, responsiveness, recovery, volatility];
end
Function MatchPrototype(fingerprint, prototypes):
   foreach prototype p in prototypes do
       similarity \leftarrow CosineSimilarity(fingerprint, p);
       if similarity > threshold then
           return p;
       end
   end
   return null;
end
Function UpdateCategories(unclassified_agents):
   if |unclassified\_agents| > min\_cluster\_size then
       clusters \leftarrow HierarchicalClustering(unclassified\_agents);
       foreach cluster c in clusters do
           if IsStable(c) then
              prototype \leftarrow ComputeCentroid(c);
              prototypes \leftarrow prototypes \cup \{prototype\};
       end
   end
end
```

**Algorithm 2:** Hierarchical Categorical Learning

When making decisions, agents retrieve instances not only from direct interactions with the

target peer but also from all peers within the same behavioral category. This categorical-based retrieval enables efficient generalization while maintaining strategic flexibility.

#### **Contrast Effects**

The third mechanism involves amplifying perceived differences between behavioral categories through contrast effects. When evaluating partners, people don't assess each peer's actions in isolation, but rather relative to their experiences with other peers [255].

The model implements this through a spreading activation mechanism:

$$A_i = B_i + s_i \sum_{j \neq A} \frac{C_{Aj}}{|O|} \tag{4.6}$$

where  $A_i$  is the total activation for instance i,  $B_i$  is the base activation from recency and frequency,  $s_i$  is a stereotype score representing how well the instance exemplifies its category,  $C_{Aj}$  is the contrast strength between categories A and j, and O is the set of other categories.

The stereotype score is calculated as:

$$s_i = \sum_{m \in M} w_m (1 - |v_{i,m} - p_{c,m}|) \tag{4.7}$$

where M is the set of behavioral metrics,  $w_m$  is the discriminative power weight for metric m,  $v_{i,m}$  is instance i's value for metric m, and  $p_{c,m}$  is category c's prototype value for metric m.

This mechanism sharpens categorical boundaries and facilitates differential responses to various relationship types by amplifying the activation of instances that strongly exemplify category-distinctive behaviors.

## 4.4.2 Model Validation and Insights

I validated this cognitive model against data from information-sharing experiments involving 150 participants (50 triads). The model successfully reproduced human behavior patterns without parameter fitting, capturing distinctive patterns observed in triadic interactions.

Figure 4.12 shows the close correspondence between model predictions and human behavior for both overall cooperation rates and specific relationship patterns. The model accurately predicted how cooperation evolved over time, including initial declines and subsequent stabilization at different levels depending on information conditions.

Analysis of the dynamic weighting parameter ( $\alpha$ ) revealed important social learning patterns:

- In mutual cooperation pairs, both partners developed high and stable  $\alpha$  values (mean = 0.78), indicating increased concern for each other's outcomes
- In mutual defection pairs,  $\alpha$  values remained consistently low (mean = 0.18)
- In asymmetric relationships,  $\alpha$  values diverged significantly, with the cooperative partner maintaining higher values than the defective partner

The most striking insights emerged from the analysis of triadic dynamics. The model captured how third-player behavior influenced relationship development between the other two members.

It correctly predicted that selectively cooperative third players would stabilize already cooperative dyads while having minimal impact on defective dyads.

The model also predicted the emergence of "social balance" effects, where triads tended toward either all-positive or mixed-sign configurations, avoiding the unstable configuration where two players have positive relationships with each other but negative relationships with the third. This phenomenon, long documented in social balance theory [103], emerged naturally from the interplay of dynamic weighting, category learning, and contrast effects without being explicitly programmed.

This theoretical framework extends beyond cybersecurity contexts to explain how individuals navigate multiple cooperative relationships in any domain characterized by strategic interdependence. By integrating dynamic weighting, categorical learning, and contrast effects, the model provides a psychologically plausible account of how humans efficiently process multiple relationships despite cognitive limitations.

The model offers several key advantages over existing approaches: (1) Cognitive plausibility: Unlike models that assume unlimited memory or processing capacity, this approach incorporates established constraints on human cognition. (2) Scalability: The categorical processing mechanism allows efficient handling of multiple relationships without computational explosion. (3) Emergent complexity: Complex social phenomena such as coalition formation and in-group/out-group dynamics emerge naturally from the interplay of basic cognitive mechanisms.

These insights suggest ways to design interventions and systems that better align with natural cognitive processes, potentially improving cooperative outcomes in domains ranging from cyber-security information sharing to organizational collaboration. For example, information sharing platforms could be designed to facilitate category-based processing by grouping organizations with similar sharing behaviors, thereby reducing cognitive load while enhancing cooperative responses.

### 4.5 Discussion and Conclusion

Our research examined cognitive mechanisms underlying multi-defender decision-making in cybersecurity through two complementary approaches: controlled experimentation with triadic information sharing and computational modeling of how individuals manage multiple simultaneous relationships. Together, these studies reveal fundamental insights about cooperation in multi-agent cybersecurity contexts.

Three key findings emerge from this work. First, the experimental results revealed a counter-intuitive inverse U-shaped relationship between information availability and cooperation. Experiential information produced the highest cooperation rates (71.7%), while descriptive information including complete payoff matrices actually suppressed cooperation (49.6%). This challenges assumptions that transparency universally promotes cooperation and suggests that abstract strategic information may trigger defensive behaviors while concrete outcome feedback fosters reciprocity norms.

Second, our analysis uncovered distinctive triadic dynamics that cannot be reduced to dyadic interactions. Participants evolved from universal cooperation toward selective strategies, with third-party behaviors creating complex mediation effects. Selectively cooperative third players

stabilized existing cooperative dyads while potentially destabilizing mixed-strategy pairs, demonstrating that *n*-person interactions exhibit emergent properties beyond paired relationships.

Third, our cognitive model demonstrates that these complex patterns can emerge from three basic mechanisms: dynamic weighting based on relationship value ( $w_i = \frac{V_i}{\sum_j V_j}$ ), categorical processing to manage cognitive load, and contrast effects in partner evaluation. The model successfully predicted empirical phenomena including social balance configurations and coalition formation without explicit programming, validating the cognitive approach to multi-agent decision-making.

These findings have important implications for cybersecurity information sharing. Rather than emphasizing theoretical benefits, platforms should provide clear feedback about concrete outcomes from shared intelligence. System designers must consider network-level effects of selective cooperation, as bilateral sharing decisions influence broader group dynamics. The cognitive mechanisms identified suggest that interfaces facilitating category-based processing of partners could reduce cognitive load while maintaining cooperation.

Several limitations should be acknowledged. Participants were not cybersecurity professionals, potentially limiting generalizability. The three-person groups, while analytically tractable, simplify real-world sharing networks. The specific payoff parameters may not reflect asymmetries in actual cybersecurity contexts where sharing risks often exceed benefits. Future work should validate findings with security practitioners, explore scaling to larger networks, and investigate how time pressure affects categorical versus individuated partner evaluation. Integration of cognitive models with autonomous agents could yield systems that leverage both human adaptability and computational power for enhanced cyber defense.

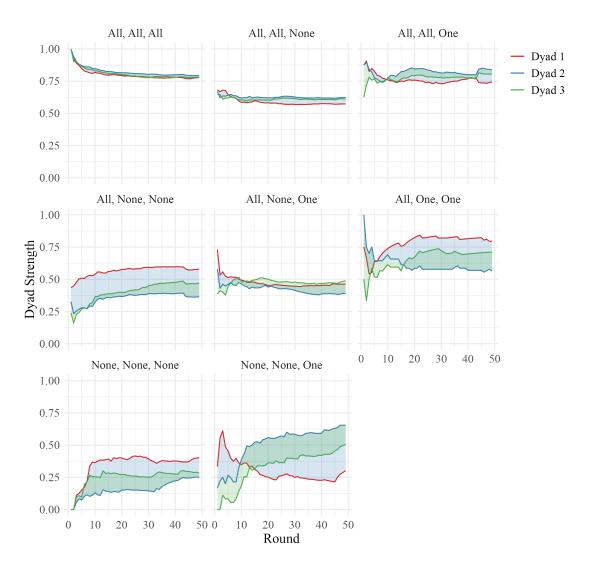


Figure 4.10: Evolution of Dyadic Relationship Strengths in Triads under different initial conditions. The shaded ribbon illustrates the growing disparity between strongest and weakest relationships over time, demonstrating how initial asymmetries tend to amplify.

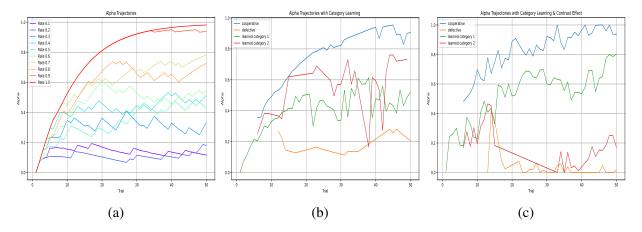


Figure 4.11: Evolution of alpha values under different cognitive mechanisms across 50 trials, all tracking interactions with peers of varying cooperation rates (0.1-1.0). (a) Baseline condition showing individual alpha trajectories for each peer. (b) With category, learning enabled alpha trajectories by learning behavioral categories rather than individual peers. (c) The combined effect of category learning and contrast mechanisms demonstrates enhanced separation between learned behavioral categories. Higher alpha values indicate greater weight given to an peer's outcomes in decision-making.

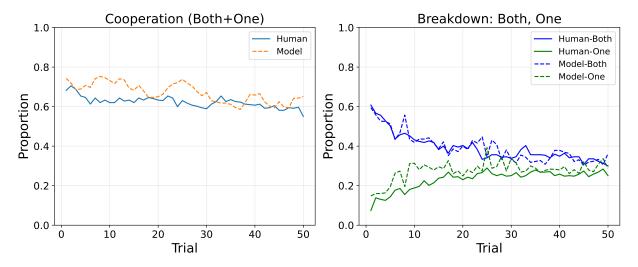


Figure 4.12: Comparison of model predictions (solid lines) and human behavior (dashed lines) for cooperation rates and relationship patterns

# Chapter 5

# Human-AI Teaming in Cyber Defense: Enhancing Collaborative Performance Through Cognitive Integration<sup>1</sup>

### 5.1 Introduction

Modern cyber defense requires capabilities that neither humans nor autonomous systems can provide in isolation. While humans excel at contextual understanding, critical thinking, and adapting to novel threats, they struggle with the speed, scale, and continuous vigilance demanded by today's security operations. Conversely, automated systems offer tireless monitoring and rapid response but lack the flexibility and judgment needed for complex, evolving threats. This fundamental limitation of single-agent approaches has driven the emergence of human-autonomy teaming (HAT) as a critical paradigm for effective cybersecurity.

HAT represents a fundamental shift from traditional human-machine interaction paradigms. As defined by [wynne2018integrative], an ideal autonomous teammate is "a highly altruistic, benevolent, interdependent, emotive, communicative and synchronized agent teammate, rather than simply an instrumental tool." This definition emphasizes that HAT involves true interdependence and coordination rather than mere parallel operation or tool use.

Despite significant advances in autonomous systems for cyber defense, including intrusion detection systems, security orchestration platforms, and automated incident response tools, these systems often function as decision support tools rather than true teammates. They may provide recommendations but typically lack agency and the ability to work interdependently with humans toward shared objectives. Meanwhile, cybersecurity operations centers face growing challenges of alert fatigue, information overload, and analyst burnout [177], creating ideal conditions for the misallocation of attention [185].

To address these challenges, autonomous systems must evolve beyond mere recommender systems and operate with higher levels of agency while maintaining appropriate human oversight [145]. As noted by [131], the cybersecurity community is increasingly recognizing the necessity of building autonomous agents that can act independently while collaborating effectively with

<sup>&</sup>lt;sup>1</sup>See Appendix .3 for published version of this chapter

human operators.

However, effective human-AI collaboration in cybersecurity contexts requires overcoming significant challenges. Autonomous agents must adapt to human working styles, properly calibrate trust, maintain appropriate transparency, and balance autonomy with human oversight. These challenges are magnified in cybersecurity operations, which are characterized by high uncertainty, time pressure, and the need for rapid adaptation to novel threats.

Although existing research has explored cognitive models of attackers and defenders in isolation, there has been limited investigation into how cognitive mechanisms influence collaboration effectiveness in human-AI defense teams. Understanding these mechanisms is crucial for designing autonomous teammates that effectively complement human strengths while compensating for human limitations.

This chapter examines the integration of human and AI decision-making in team defense scenarios where humans and AI collaboratively protect networks. Through the design and implementation of the Team Defense Game (TDG), a semi-supervisory teamwork paradigm, I investigate how different types of autonomous agents affect team performance in cybersecurity tasks. By comparing cognitive agents that learn from experience to approximate human-like decision processes with heuristic and random agents, I demonstrate that cognitive mechanisms significantly enhance team performance and efficiency in cybersecurity operations. The experiments reveal both the potential of cognitively inspired agents to improve collaborative defense and the challenges in building appropriate trust and reliance between humans and their AI teammates.

The findings advance our understanding of how to effectively integrate human and AI capabilities in cybersecurity operations, with implications for the design of future autonomous defense systems and team structures. By establishing how cognitive mechanisms influence collaboration at the team level, this research complements the individual and group-level analyses presented in previous chapters, providing a comprehensive framework for understanding human and AI decision-making in cybersecurity.

# 5.2 Related work: Human-Autonomy Teaming

Human-Autonomy Teaming (HAT) has emerged as a critical research area that examines collaborative systems where humans and autonomous agents work together to achieve shared goals. O'Neill et al. [180] conducted a comprehensive review of HAT literature, finding that effective human-autonomy teams require careful consideration of task allocation, communication protocols, and trust development. The integration of human factors into autonomous system design is crucial for successful teaming, as highlighted by Bjurling et al. [27], who proposed a framework specifically addressing design requirements for digital assistants in aviation contexts.

The theoretical foundations of HAT continue to evolve, with Lyons et al. [151] offering conceptual clarifications that differentiate HAT from traditional human-machine interaction paradigms. Their work emphasizes that effective HAT involves interdependence and coordination rather than mere parallel operation. This perspective is further supported by McNeese et al. [159], who identified trust as a fundamental component in facilitating effective collaboration between humans and autonomous systems.

Experimental paradigms for studying HAT have also advanced significantly. Schelble et

al. [211] demonstrated how reinforcement learning can be leveraged to design experimental environments that simulate realistic team interactions. Demir et al. [48] examined team coordination dynamics in HAT contexts, finding that effective interaction styles significantly influence team performance outcomes.

The growing literature on HATs in domains such as urban search and rescue [251] and hospital management [37] has identified some critical factors for successful teamwork, but little is known about HATs in cybersecurity. Teams in cybersecurity operations, especially those in 24/7 security operations centers, have specific dynamics [186]. Due to these varied and unique applications, a synthetic cyber task environment is needed to empirically evaluate HATs with different team compositions in various cyber scenarios.

### 5.2.1 Challenges and Opportunities in Cyber HAT

Today, cyber analysts are a scarce resource and are often overloaded [177]. Security Operations Centers (SOCs) combat the growing problem of alert fatigue, where the sheer volume of alerts overwhelms SOC analysts and raises the risk of overlooking critical threats [21], creating ideal conditions for misallocation of attention [185]. To address these challenges and meet the demands posed by sophisticated adversaries, autonomous systems must evolve beyond mere recommender systems and operate with higher levels of agency [145]. The cybersecurity technology community is increasingly recognizing the necessity of building autonomous agents that can act independently [131].

Recent advancements in HAT for cyber defense include the development of adaptive approaches that can respond to dynamic threat environments. Lohn et al. [147] explored how autonomous cyber defense systems can be designed to enhance human capabilities while maintaining appropriate levels of human oversight. Théron and Kott [235] examined potential future scenarios where autonomous intelligent defense systems might engage with autonomous malware, highlighting the need for robust HAT approaches in cyber warfare contexts.

It is essential, however, to explore autonomous agents that can account for the decision-maker's values or specific mission needs. For example, following a cyber attack, an AI-generated decision engine may recommend disabling an application on the compromised computer system. This action may neutralize the threat but could simultaneously endanger a mission, negatively impact a user's ability to perform critical tasks, or allow the adversary to extend the duration or scope of the attack [145]. Human experts should remain in the loop to provide intuition, critical thinking, and contextual information by approving or denying recommendations from AI decision engines that may have negative impacts [209].

Different cybersecurity scenarios also pose unique challenges for HAT. For instance, in incident response and recovery, autonomous agents might focus on information triage while leaving further analysis and strategic decision-making to humans. In adaptive defense, autonomous agents can more efficiently adjust security mechanisms based on real-time threat intelligence, with humans supervising and fine-tuning agent decisions only when necessary [145].

### 5.2.2 Research Gaps in HAT for Cyber Defense

Despite the increasing body of research on autonomous systems for cyber defense and cognitive modeling of adversarial behavior, significant gaps remain in our understanding of human-autonomy teaming for cybersecurity applications. While cognitive agents have been developed to model both attackers and defenders in isolation [57, 55], there has been limited investigation into how cognitive agents might function as teammates alongside human operators in collaborative cyber defense scenarios.

The existing research has demonstrated the effectiveness of cognitive agents in simulating adversarial behavior [54, 55] and in modeling individual cyber defense decisions [57], but has not explored how cognitive agents might be integrated into HAT frameworks for operational cybersecurity. This represents a critical gap, as the effectiveness of HAT in cybersecurity contexts depends not only on the individual capabilities of autonomous agents but also on their ability to collaborate effectively with human teammates.

Furthermore, while trust has been identified as a fundamental component of effective HAT [159], there is limited empirical evidence regarding how trust develops between human operators and cognitive agents in cybersecurity contexts. Understanding the factors that influence trust formation and maintenance in cyber HAT could provide valuable insights for designing more effective collaborative systems.

Additionally, although prior research has explored the role of cognitive models in predicting human defensive behaviors [57], there has been little investigation into how cognitive agents might adapt to the unique working styles and preferences of individual human teammates. This adaptability is likely to be crucial for effective collaboration in complex and dynamic cybersecurity environments.

Existing studies have also not sufficiently addressed the challenge of balancing agent autonomy with human oversight in cyber defense contexts. While Lohn et al. [147] and Théron and Kott [235] have explored the potential for autonomous cyber defense systems, the optimal division of responsibilities between humans and autonomous agents in different cybersecurity scenarios remains unclear.

Finally, there is a lack of empirical research comparing different types of autonomous agents (e.g., heuristic-based, machine learning-based, and cognitive agents) in HAT for cyber defense. Understanding the relative strengths and limitations of different agent architectures in collaborative cybersecurity tasks could inform the development of more effective HAT systems.

### 5.2.3 Methodological Approaches to HAT Research

The development of robust methodologies for studying HAT continues to be an active area of research. Neubauer et al. [172] introduced a Human-Autonomy Team Cohesion Scale, providing a validated instrument for measuring team dynamics in HAT contexts. This scale offers researchers a tool for assessing how different factors influence team cohesion when humans collaborate with autonomous systems.

Guidetti et al. [87, 88] explored the use of neuroergonomic approaches in cyber vigilance tasks, developing frameworks for measuring cognitive and physiological responses during network defense activities. Their work demonstrates how multidisciplinary methods can provide insights

into the cognitive demands of cybersecurity tasks and inform the design of effective HAT systems.

The integration of human factors considerations into HAT research has been emphasized by Ulusoy and Reisman [242], who argue for the importance of respecting human needs and capabilities in the design of autonomous systems. Their work suggests that successful HAT requires attention not only to technical system performance but also to human experience and well-being.

For the partnership between humans and autonomous agents to be successful, the potential benefits of HAT must be weighed against foreseeable negative human-autonomy interactions. Unintended consequences that must be addressed include creating more (not less) work for humans, failing to decrease required manpower, deskilling operators, reducing awareness, and contributing to accidents [231, 93, 150]. These concerns highlight the importance of carefully designed synthetic task environments for empirically evaluating HAT before implementation in operational settings.

With advances in computational power, network robustness, and machine learning capabilities, a new form of team is emerging in cybersecurity operations: the human-autonomy team (HAT). These teams integrate human analysts with autonomous agents, requiring both members to depend on each other to achieve collective security goals [158]. Although research has examined the effects of agent performance [19] and perceived warmth [102] on HAT effectiveness, there is increasing recognition that human-like qualities offer unique advantages for facilitating human-agent cooperation [76, 187].

Despite this growing interest, there remains limited empirical investigation regarding how humans collaborate with autonomous agents that emulate human cognitive processes [80]. This gap is particularly pronounced in cybersecurity contexts, where teams face unique challenges including high-stakes decisions, time pressure, and the need to adapt to novel threats.

Today's cyber analysts are scarce resources and frequently overloaded [177]. Security Operations Centers (SOCs) struggle with alert fatigue, where the sheer volume of security alerts overwhelms analysts and increases the risk of overlooking critical threats [22]. This creates ideal conditions for misallocation of attention and defensive failures. To address these challenges, autonomous systems must evolve beyond mere recommender systems to operate with higher levels of agency while maintaining appropriate human oversight [145].

The key questions this research addresses are: Do human-like cognitive agents have advantages over optimally performing non-cognitive agents in HAT collaborations? How do humans perceive the cooperativeness and trustworthiness of cognitive versus non-cognitive agents in cybersecurity contexts?

# 5.3 Design of the Team Defense Game

To investigate these questions, I designed the Team Defense Game (TDG), an experimental platform for studying how humans make decisions in collaboration with autonomous agents to defend a network from cyber attacks. The TDG extends previous work on interactive defense games [192] to incorporate teamwork dynamics through a semi-supervisory framework.

In TDG, human participants play the role of cyber analysts tasked with protecting a computer network against external malicious activity. Each participant is paired with an autonomous cyber

defense agent who can make decisions and partially act independently to collaborate in network defense. The human and autonomous agent must work together to monitor the network, detect suspicious activity, and take appropriate defensive actions. Figure 5.1 shows the interface through which participants interacted with the system.

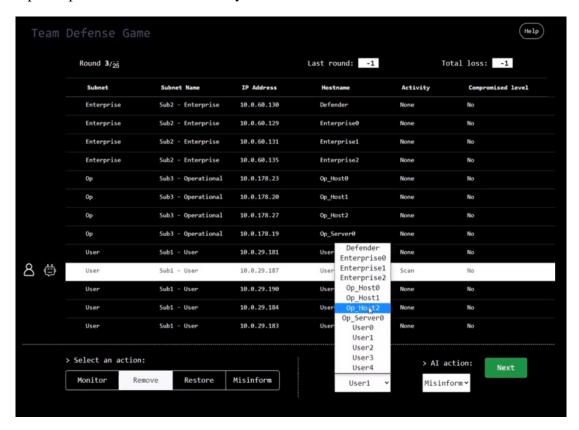


Figure 5.1: Team Defense Game interface showing network status and interaction options

The framework implements a semi-supervisory structure that mirrors real-world cyber defense teams. As shown in Figure 5.2, the autonomous agent has a set of pre-approved actions it can execute independently (Monitor, Remove) and other actions (Restore, Misinform) that require human approval before execution. This creates a hierarchical relationship where the human maintains oversight while allowing the agent to operate with partial autonomy.

In each step of the game, both the human and the autonomous agent independently decide on a target (which computer or server to protect) and an action to take. After both have submitted their intentions, the human is presented with the agent's intended action. If the agent selected a pre-approved action, it executes automatically. If the agent selected a non-pre-approved action, the human must validate or modify the action before it can be executed. If both team members selected the same target, the human must resolve this overlap by modifying either their own or the agent's intention.

This design creates three key interaction scenarios that reveal team dynamics:(1) Overlap: When human and agent select the same target, requiring coordination, (2) Supervision: When the agent requires approval for non-pre-approved actions, (3) Backup: When multiple hosts are compromised, requiring strategic allocation of team resources.

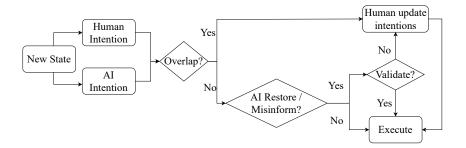


Figure 5.2: Semi-supervisory framework showing pre-approved and non-pre-approved actions

### 5.3.1 Autonomous Agent Types

The study compared three types of autonomous agents as teammates:

- 1. **Cognitive Agent**: An IBLT-based agent that learns from experience to make decisions based on similarity to past situations. The agent accumulates instances containing state, action, and utility information, retrieving these based on similarity to current situations. This agent adapts to the environment and to the human's behavior through experience.
- 2. **Heuristic Agent**: Follows rule-based strategies derived from expert knowledge about optimal network defense. This agent is highly competent but not adaptive, using fixed decision rules based on the status of network hosts and the history of attacks.
- 3. **Random Agent**: Makes decisions randomly, serving as a baseline for comparison. This agent has no strategic knowledge or learning capability.

Both the Cognitive and Heuristic agents were designed to achieve comparable performance levels when operating independently, allowing the study to isolate the effects of cognitive mechanisms rather than simple competence differences.

# 5.4 Experiment: Human-Autonomy Cyber Defense Team Performance

The experiment employed a between-subjects design with 156 participants randomly assigned to work with one of the three agent types: Cognitive (n=42), Heuristic (n=48), or Random (n=66). Each participant completed 7 episodes of the Team Defense Game, with each episode consisting of 25 steps protecting a network against an adversary.

Measurements included: (1) Team Performance: Total loss (points lost due to successful attacks and defensive actions) and recovery time (steps required to restore compromised hosts), (2) Collaborative Process Metrics: Frequency and handling of overlaps, supervision situations, and backup requirements, (3) Human Effort and Efficiency: Proportion of active versus passive actions and efficiency (loss reduction per effort expended), (4) Human Perception: Post-experiment ratings of agent trustworthiness and cooperativeness.

### 5.4.1 Key Findings

### **Team Performance**

Teams with Cognitive agents achieved significantly better performance than those with other agent types, as shown in Figure 5.3. HATs with Cognitive agents experienced lower average loss (M=-52.85, SD=27.42) compared to teams with Heuristic agents (M=-59.69, SD=28.54) and Random agents (M=-79.69, SD=49.11). Similarly, teams with Cognitive agents demonstrated faster recovery from compromised hosts (M=1.45, SD=1.89) than those with Heuristic (M=2.15, SD=2.09) or Random agents (M=3.92, SD=3.89). A two-way mixed measures ANOVA confirmed a significant main effect for agent type on team loss, F(2,152)=11.037,  $p_1.05$ , with post-hoc tests using Tukey's HSD indicating that teams with Cognitive agents achieved significantly lower losses than those with Heuristic or Random agents.

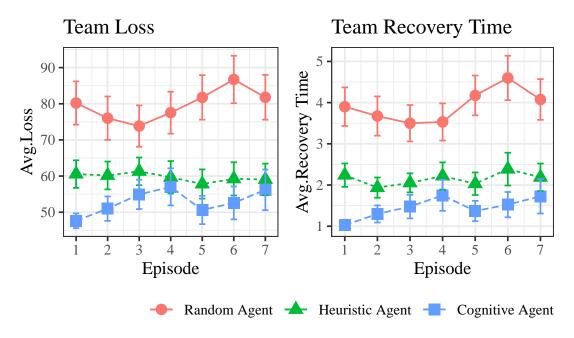


Figure 5.3: Team performance by agent type across episodes: team loss (left) and recovery time (right)

#### Collaborative Process Metrics

Analysis of collaborative processes revealed important differences in how teams functioned with different agent types: **Overlap Resolution**: Overlaps (both team members selecting the same target) occurred most frequently with Random agents. When resolving overlaps, humans were more likely to adjust their own actions (rather than the agent's) when working with Heuristic agents (29%) compared to Cognitive (21%) or Random agents (14%). This suggests greater trust in the Heuristic agent's predictable decision-making. **Supervision Dynamics**: Random agents required significantly more supervision (49% of actions) than Cognitive (37%) or Heuristic agents (33%). More importantly, humans agreed with agent recommendations at much higher rates when

working with competent agents (73-74% for Cognitive and Heuristic) compared to Random agents (35%). Figure 5.4 shows how this agreement proportion evolved over time, with humans rapidly learning to trust competent agents while decreasing trust in Random agents. **Backup Behavior**: Multiple breaches requiring team coordination occurred much more frequently with Random agents (41%) than with Cognitive or Heuristic agents (14%). However, humans provided backup less frequently when paired with Cognitive agents (11%) compared to Heuristic (17%) or Random agents (20%). This suggests potential over-reliance on the Cognitive agent's capabilities, with humans less inclined to provide assistance even when needed.

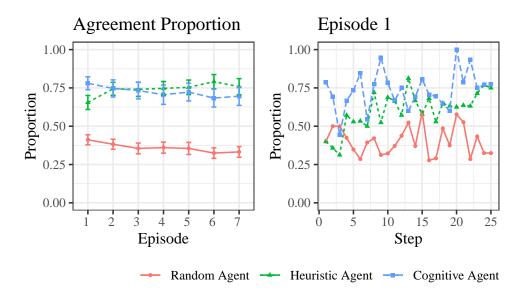


Figure 5.4: Evolution of human agreement with agent recommendations during the first episode (right) and across all episodes (left)

### Human Effort and Efficiency

Human effort (frequency of taking active actions vs. passive monitoring) decreased over episodes across all conditions, but efficiency varied significantly by agent type. As shown in Figure 5.5, participants working with Cognitive agents achieved the highest efficiency (M=39.40, SD=46.53), followed by those with Heuristic agents (M=29.82, SD=38.32) and Random agents (M=16.80, SD=27.02).

A two-way mixed measures ANOVA confirmed significant main effects for agent type  $(F(2,152)=7.949, p_i.001)$ , episode  $(F(3.46,630.77)=5.322, p_i.001)$ , and their interaction (F(8.30,630.77)=2.435, p=.012) on human efficiency. This indicates that Cognitive agents enabled humans to be more efficient with their actions, achieving better results with similar effort levels.

### **Human Perception of Agents**

Post-experiment questionnaires revealed that participants rated Cognitive and Heuristic agents significantly higher than Random agents on both cooperativeness and trustworthiness dimensions,

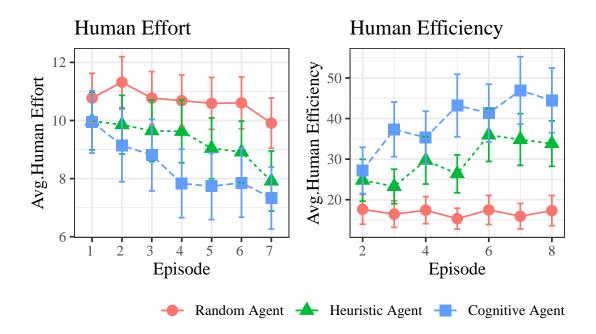


Figure 5.5: Evolution of human effort (left) and efficiency (right) across episodes by agent type

as shown in Figure 5.6. Approximately 52% of participants in the Cognitive condition and 50% in the Heuristic condition agreed or strongly agreed that their agent was cooperative and trustworthy, compared to only 33% in the Random condition.

Open-ended feedback revealed interesting differences in how participants conceptualized their relationship with different agent types. Those working with Cognitive agents often noted inconsistency but adaptability, while those with Heuristic agents appreciated predictability but sometimes felt the agent was inflexible. Several participants with Random agents reported high confidence in the agent despite its poor performance, suggesting potential overreliance based on the agent's perceived authority rather than demonstrated capability.

### 5.5 Discussion and Conclusion

The findings suggest three major implications for the design of autonomous agents in cybersecurity teams.

Human-like Cognition Benefits Team Performance: Cognitive agents that emulate human learning processes significantly enhanced team performance compared to both heuristic and random agents. The ability to learn from experience and adapt to the individual play styles of human teammates appears to be particularly valuable in the dynamic cybersecurity context. However, the inconsistency and unpredictability of cognitive agents sometimes reduced human trust compared to more predictable heuristic agents. This suggests that cognitive agents should be designed to maintain adaptability while providing explanations for changing behaviors to maintain human trust.

Competence Affects Trust and Reliance: Both cognitive and heuristic agents were perceived

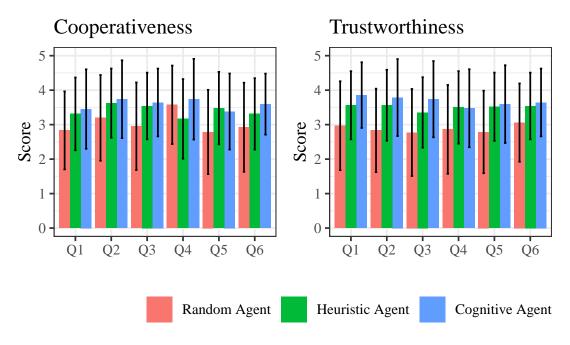


Figure 5.6: Human ratings of agent cooperativeness (left) and trustworthiness (right) by agent type

as more trustworthy and cooperative than random agents, demonstrating that baseline competence is a prerequisite for effective human-AI teaming. However, competent agents also required less human effort, potentially leading to over-reliance. This highlights the importance of calibrating trust appropriately by designing agents that actively signal when they require human assistance or have uncertainty in their decision-making processes.

**Trust Calibration Remains Challenging:** The experiment revealed that human trust in agents developed rapidly within the first episode and was strongly influenced by perceived competence. However, this trust was not always well-calibrated to actual agent performance. Some participants developed excessive trust in random agents despite poor performance, while others maintained skepticism toward cognitive agents despite their superior performance. This suggests that trust development in HAT contexts is influenced by factors beyond objective performance metrics, including consistency, predictability, and transparency.

These findings contribute to our understanding of human-AI teaming in cybersecurity by demonstrating how cognitive mechanisms influence collaboration effectiveness and trust dynamics. The results highlight the potential of cognitive agents to enhance team performance while also revealing important challenges in building appropriate trust and reliance.

### 5.5.1 Limitations and Future Work

While this study provides valuable insights into HAT dynamics in cybersecurity, several limitations and future directions should be acknowledged. First, participants were not cybersecurity professionals, potentially limiting the generalizability of findings to operational contexts. Future work should validate these findings with expert participants from security operations centers.

Second, the experimental task, while capturing essential elements of cyber defense decision-making, simplified many aspects of real-world security operations. Future research should explore more complex scenarios with higher fidelity to operational environments, including variable threat levels, resource constraints, and mission priorities.

Finally, the study focused primarily on performance and process metrics with limited investigation of how the agents' decision-making processes were perceived and understood by human teammates. Future work should incorporate explainable AI techniques to make cognitive agent reasoning more transparent to human teammates, potentially enhancing trust calibration and team coordination.

Despite these limitations, this research represents an important step toward understanding how cognitive mechanisms influence human-AI collaboration in cybersecurity contexts. By demonstrating the potential of cognitively inspired agents to enhance team performance, the findings provide guidance for designing more effective human-autonomy teams for cyber defense.

# Chapter 6

# Conclusion & Future Directions

### 6.1 Conclusion

This dissertation has addressed fundamental challenges in cybersecurity by investigating the integration of human cognitive processes with artificial intelligence capabilities through a multiagent modeling perspective. The research presented here moves beyond traditional approaches to cybersecurity by systematically studying how cognitive mechanisms influence both offensive and defensive behaviors in complex, dynamic environments. The empirical and computational studies conducted across the three research dimensions reveal several profound insights about the nature of decision-making in cybersecurity contexts. Perhaps most significantly, the findings demonstrate that realistic modeling of human cognitive processes—including bounded rationality, experiential learning, and categorical reasoning—is not merely an academic exercise but a practical necessity for developing effective cybersecurity systems.

In examining human-like adversaries, this work challenges prevailing assumptions about attacker modeling. The observation that even highly skilled human defenders struggle against cognitive attackers while performing well against deterministic strategies reveals a critical vulnerability in current training approaches. This discovery suggests a fundamental reconsideration of how we prepare cybersecurity professionals, shifting from static pattern recognition to adaptive response against dynamic threats. The cognitive attacker implementation demonstrates that Instance-Based Learning Theory provides a viable framework for creating realistic adversary emulations that capture the adaptive, learning behaviors characteristic of sophisticated human attackers.

The investigation of multi-defender interactions revealed unexpected patterns in cooperative behavior that contradict theoretical predictions about information and incentives. The inverse U-shaped relationship between information availability and cooperation—with experiential information producing higher cooperation rates than both minimal and descriptive information—challenges fundamental assumptions about strategic decision-making in group contexts. Furthermore, the emergence of selective cooperation strategies and third-player effects on dyadic relationships reveals that triadic interactions cannot be understood as simple aggregations of independent pairs. These findings have significant implications for the design of information sharing frameworks in cybersecurity, suggesting that experiential learning mechanisms may be

more effective than abstract descriptions of theoretical benefits. The cognitive model developed to explain these findings—integrating dynamic weighting, category learning, and contrast effects—provides a psychologically plausible account of how individuals navigate multiple cooperative relationships despite cognitive limitations. This model successfully reproduced human behavior patterns without parameter fitting, demonstrating its explanatory power. Beyond cybersecurity, this theoretical contribution advances our understanding of human social learning in complex environments with multiple interdependent relationships.

The exploration of human-AI teaming in cyber defense revealed both the potential and challenges of cognitive integration in operational contexts. Cognitive agents that learn from experience significantly enhanced team performance and human efficiency compared to both heuristic and random agents. However, the findings also uncovered important challenges in trust calibration and relationship development. The observation that humans sometimes developed excessive trust in random agents despite poor performance, while maintaining skepticism toward cognitive agents despite superior performance, highlights the complex psychological dimensions of human-AI collaboration. These insights extend beyond cybersecurity to inform human-AI teaming across domains where appropriate trust calibration is critical.

Methodologically, this dissertation has introduced novel experimental paradigms—including the Interactive Defense Game and Team Defense Game—that enable systematic investigation of complex phenomena while maintaining experimental control. These platforms provide valuable tools for future research at the intersection of cognitive science and cybersecurity, allowing for controlled manipulation of factors such as information availability, incentive structures, and team composition.

The overall findings from this research suggest a fundamental shift in how we conceptualize cybersecurity operations—moving from isolated technical solutions toward integrated sociotechnical systems that acknowledge the cognitive dimensions of both threats and defenses. By recognizing the distinctive cognitive processes that shape adversarial, cooperative, and collaborative behaviors, we can develop more robust approaches to cybersecurity that leverage complementary human and artificial intelligence capabilities. This integrated perspective has important implications for cybersecurity practice. For training and education, it suggests developing scenarios that incorporate adaptive adversaries rather than deterministic patterns. For information sharing frameworks, it indicates prioritizing experiential feedback about concrete outcomes over abstract payoff descriptions. For human-AI teaming, it emphasizes designing agents that can adapt to individual human working styles while providing appropriate transparency about their decision processes.

The cognitive modeling approaches developed in this dissertation provide a foundation for future research across multiple domains. As artificial intelligence continues to evolve, understanding how to integrate human and AI decision-making will become increasingly crucial not only for cybersecurity but for all domains characterized by complex, dynamic threats and resource constraints. By advancing our understanding of the cognitive mechanisms that enable effective decision-making in multiagent contexts, this research contributes to the broader goal of creating more secure and resilient systems through the thoughtful integration of human and artificial intelligence.

# 6.2 Future Work: Human-like Adversaries Modeling

In our recent work, we demonstrated that cognitive attackers capable of learning from experience and adapting their strategies present significantly greater challenges for defenders compared to deterministic attack strategies. Building on this foundation, I propose integrating large language models (LLMs) with human-like decision-making processes to create more realistic, adaptive cyber adversaries for security training and testing. This approach promises to better prepare defenders for the unpredictable and evolving tactics used by real human attackers in operational environments. The proposed research will develop autonomous attackers that incorporate key aspects of human cognition, including bounded rationality, risk assessment under uncertainty, learning from past interactions, and strategic adaptation to defensive countermeasures. Unlike traditional rule-based attack simulations that follow predetermined patterns, these cognitive adversaries will exhibit the flexibility, creativity, and unpredictability characteristic of human attackers, while maintaining the controllability needed for systematic training and evaluation.

Recent advances in LLMs have demonstrated promising capabilities for cybersecurity applications. Works like PenHeal have shown that LLMs can effectively model multistage penetration testing workflows, combining reconnaissance, vulnerability assessment, and exploitation in coherent attack chains [108]. Similarly, PENTESTGPT has demonstrated that LLMs can generate contextually appropriate attack strategies based on system descriptions and security configurations [49]. These capabilities, when combined with cognitive models of decision-making under uncertainty, create an opportunity to develop adversaries that not only execute technical attacks but do so with human-like strategic reasoning. The technical approach involves developing a multi-component system where an LLM handles higher-level attack planning and reasoning about network configurations, while a cognitive decision-making module manages experience-based learning, adaptation, and execution of attack sequences. This integration preserves the human-like learning patterns that made our previous cognitive attackers effective, while adding the sophisticated reasoning capabilities of modern language models.

Modeling Adversarial Biases A critical aspect of human attackers that distinguishes them from current automated systems is their susceptibility to cognitive biases and limitations. Even sophisticated attackers make non-optimal decisions due to information processing constraints, risk perception biases, and emotional factors. These biases create exploitable patterns that defenders can leverage if properly understood. Our proposed framework will incorporate these cognitive limitations, modeling how attackers prioritize targets, balance exploration versus exploitation, and respond to deceptive defensive measures. By systematically investigating how multiple biases interact in sequential attack decisions, we can identify network configurations and defensive strategies that maximize attacker inefficiency while minimizing defensive resources. The framework will support the exploration of adversarial biases in specific cyber scenarios. By modeling how attackers with different cognitive profiles navigate the same network environment, we can identify defensive configurations that exploit these biases most effectively.

**Evaluation and Applications** The system will be evaluated in simulated environments of increasing complexity, measuring both technical performance (attack success rates, time to ob-

jective) and behavioral fidelity (similarity to human attack patterns observed in controlled experiments and real-world incidents). Beyond serving as training tools, these human-like adversaries can also be used to evaluate defensive systems, identify potential vulnerabilities in security postures, and develop novel defensive strategies that specifically target human cognitive weaknesses. This approach moves beyond traditional static defenses toward adaptive systems that can reconfigure to present the most challenging scenarios based on observed attacker behavior. For example, certain network topologies or deceptive elements might be particularly effective against attackers exhibiting specific cognitive biases like confirmation bias or availability heuristics. Recent work by Singer et al. demonstrates the feasibility of using LLMs to execute multistage network attacks [220], while Li and Zhu have explored symbiotic game models for cyber deception operations [141]. The research extends beyond technical implementation to include fundamental questions about adversarial cognition: How do attackers build and update mental models of target networks? How do they allocate attention across multiple potential attack vectors? How do emotional factors like frustration or overconfidence affect their strategic decisions? Answering these questions through computational modeling will not only improve our ability to create realistic training scenarios but also enhance our theoretical understanding of adversarial behavior in cybersecurity contexts.

## 6.3 Future Work: Toward Autonomous Intelligent Cyber Defense

The increasing complexity and scale of cyber threats necessitate advanced approaches to network defense that can operate effectively in dynamic environments without constant human supervision. As cyber attacks grow in sophistication, traditional defensive approaches that rely on static rules or signatures become increasingly inadequate. My research aims to address this challenge by developing autonomous intelligent cyber defense systems that can adapt to evolving threats, anticipate attacker behaviors, and coordinate defensive actions across complex networks.

Advanced cyber defense requires intelligent systems that can demonstrate human-like reasoning capabilities while operating at machine speed. These systems must develop a sense of causality that discovers relationships between objects and events, allowing incorporation of temporal and spatial information into reasoning processes. They must also balance potentially conflicting objectives while operating safely in poorly understood environments, requiring advances in risk-aware online planning [131]. As these systems grow more complex, they must transition from isolated defensive tools to coordinated teams of defensive agents, working together with human operators to protect critical infrastructure.

I have begun exploring this research direction through preliminary work on reinforcement learning for adaptive cyber deception [56]. In this work, I demonstrated how a defender trained through deep reinforcement learning could strategically deploy deceptive elements throughout an attack graph to significantly delay attackers compared to static or heuristic approaches. By learning through self-play, the defender developed strategies that dynamically responded to attacker progression, maximizing the effectiveness of limited defensive resources. This initial exploration proved the viability of reinforcement learning for cyber defense but also revealed significant challenges that must be addressed for real-world deployment.

Over the next five years, my research will focus on two critical directions that build upon

this foundation: developing sample-efficient learning methods for cyber defense and advancing multiagent reinforcement learning frameworks for cooperative defense. These directions address key limitations of current approaches while pushing toward truly autonomous cyber defense systems.

### 6.3.1 Sample-Efficient Learning for Complex Defense Environments

A fundamental challenge in applying reinforcement learning to cybersecurity problems is the sample complexity of current approaches. Model-free reinforcement learning algorithms typically require millions of interactions with the environment to learn effective policies—a requirement that is impractical for operational networks. My research will address this challenge through three complementary approaches to sample-efficient learning: model-based reinforcement learning, incorporation of domain knowledge, and reinforcement learning from human feedback (RLHF).

Model-based reinforcement learning offers a promising path forward by enabling agents to learn an internal model of environment dynamics. This approach allows for planning and simulation-based learning that can drastically reduce the number of real-world interactions required during training [261]. In the context of cyber defense, I will develop techniques that enable defenders to construct and refine mental models of attacker behavior patterns, vulnerabilities, and network dynamics. The resulting models will allow defenders to anticipate potential attack vectors through simulated what-if scenarios, effectively learning from imagined experience rather than requiring extensive real-world interactions.

Research in understanding effective memory structure and processes will benefit from a collaboration with cognitive scientists to understand memory in biological systems [11]. New approaches are needed to address potential issues with memory systems such as catastrophic forgetting, limited storage capacity, and development of new methods to efficiently use external knowledge stores. My work will investigate how cognitive architectures can inform the design of memory systems for cyber defense agents, drawing inspiration from human memory processes to develop more robust and flexible learning mechanisms.

Domain knowledge incorporation represents another avenue for improving sample efficiency. Cyber defense involves well-established principles and heuristics developed through decades of operational experience. Rather than learning everything from scratch, my research will develop methods for encoding this domain knowledge as inductive biases that guide exploration and accelerate learning. This may include attack graph structures, common vulnerability patterns, or established defense strategies. By structuring the learning problem with appropriate priors, reinforcement learning agents can focus exploration on promising regions of the strategy space, significantly reducing training time while improving generalization to new threats [35].

Reinforcement learning from human feedback (RLHF) offers a third approach to sample efficiency by leveraging human expertise to guide agent learning. In cybersecurity contexts, experienced defenders possess tacit knowledge that may be difficult to formalize but can be expressed through demonstrations or feedback. My research will develop frameworks for capturing this expertise and incorporating it into reinforcement learning pipelines, allowing agents to learn from both simulation and human guidance. This approach will be particularly valuable for teaching agents about subtle signals or patterns that human defenders recognize but that might not be obvious in raw network data [39].

Together, these approaches to sample-efficient learning will enable reinforcement learning agents to develop effective defense strategies with substantially fewer training examples, making deployment in operational environments more practical. By combining model-based reasoning, domain knowledge, and human feedback, my research will produce cyber defense agents that learn more efficiently while maintaining the adaptability that makes reinforcement learning appealing for security applications.

### 6.3.2 Multiagent Reinforcement Learning for Cooperative Defense

As networks grow in scale and complexity, defense responsibilities are increasingly distributed across multiple agents with specialized roles and capabilities. Effective protection requires seamless coordination among these defensive components, along with adaptation to multiple concurrent threats. My second research direction will focus on multiagent reinforcement learning (MARL) frameworks that enable coordinated defensive actions across distributed agents.

Research in ad-hoc teamwork will enable entities (human and systems) to dynamically join together to address specific problems, then pursue separate tasks after the problem is solved [105]. In this type of teaming, there is no prior coordination between agents, and we cannot assume that the entities share the same types of learning algorithms or reward structures or that they have prior agreements regarding action coordination and information sharing. My research will address important problems within ad-hoc teaming, including ensuring that actions are understandable to fellow teammates, modeling the capabilities of team members, including humans in the ad-hoc teams, and dynamically modeling the performance of both the team and the individuals.

A key challenge in multiagent cyber defense is developing coordination mechanisms that balance local decision-making with team-level objectives. I will investigate techniques for enabling coordinated defensive actions across network zones without requiring full observation sharing or centralized control [91]. This approach is essential for large-scale networks where complete state information may be unavailable or prohibitively expensive to communicate. My work will develop efficient coordination mechanisms that enable defenders to perform complementary actions that amplify overall security outcomes while respecting communication constraints.

Another critical aspect of multiagent defense involves integrating diverse defensive techniques into a unified framework. While my preliminary work has focused primarily on deceptive deployments, effective cyber defense requires a broader arsenal including moving target defense, adaptive access control, and strategic resource allocation [145]. I will develop multiagent frameworks that can simultaneously reason about these diverse defensive options, learning when and how to deploy each technique for maximum combined impact. This integrated approach will overcome the limitations of current systems that treat different defensive mechanisms in isolation, failing to capture potential synergies between complementary tactics.

The presence of multiple attackers with diverse objectives and capabilities presents additional challenges for coordinated defense. Real-world networks face concurrent threats ranging from opportunistic attackers to targeted advanced persistent threats (APTs), each requiring different defensive responses. My research will develop models for recognizing and differentiating between multiple attacker profiles, enabling defenders to appropriately prioritize and respond to the most critical threats [67]. This will involve both technical advances in threat attribution and strategic

reasoning about optimal resource allocation across multiple simultaneous engagements.

The culmination of this research direction will be a framework for multiagent cyber defense that can scale to realistic network environments with thousands of nodes, heterogeneous services, and complex interdependencies. This will require novel approaches to state abstraction and hierarchical planning that can maintain computational tractability while capturing essential security dynamics [131]. By incorporating both topological and semantic features of network components, my methods will enable more nuanced defensive strategies that consider not just connectivity but also the business context and criticality of protected assets.

Over the next five years, this comprehensive research agenda will advance the state of the art in autonomous cyber defense, moving from theoretical models to practical deployable systems that can significantly enhance the security posture of real-world networks. By addressing key challenges in sample efficiency and multiagent coordination, my work will bridge the gap between current research prototypes and operational security solutions. The resulting advances will not only contribute to the academic understanding of reinforcement learning in adversarial contexts but also provide concrete tools and techniques that security practitioners can deploy to protect critical infrastructure against evolving cyber threats.

# 6.4 Future Work: Complementary Human-AI Teaming in Cyber Defense

My dissertation research has opened several promising research directions in human-AI teaming for cyber defense that I plan to pursue over the next five years. These directions build upon the foundations established in my work and address fundamental challenges in developing effective human-AI collaborative systems for cybersecurity operations.

## 6.4.1 Platform for Controlled Human-AI Team Experiments

The Team Defense Game (TDG) platform developed in this dissertation has demonstrated significant potential for studying human-AI collaboration in cyber defense scenarios. This work has revealed opportunities for platform advancements that would enable more sophisticated research into team dynamics and coordination mechanisms in cybersecurity contexts.

A natural extension of the current platform would be a parametric task generation framework that allows systematic control over the complexity dimensions of cyber defense scenarios. Such a framework would enable controlled experimentation by manipulating state-space complexity, strategic uncertainty, and interdependence patterns while maintaining ecological validity. As suggested by [135], multi-agent experimental environments benefit significantly from parameterized control over environmental complexity. The research direction involves creating structured multi-agent decision tasks that can be grounded in specific domain contexts through declarative schemas and language model prompting. This approach would facilitate the creation of realistic, variable, and interpretable team tasks for both experimental and training purposes.

Current research in human-AI teaming frequently overlooks the systematic control of interdependence structures between team members, focusing instead on varying the number of agents or

their observation capabilities. [157] have emphasized that emergent team cognition in human-AI settings depends critically on the structure of interdependence between team members. My work has revealed the importance of modeling both control and informational dependencies through directed graph representations. This representation enables the investigation of how different team structures—whether hierarchical, distributed, or hybrid—affect coordination dynamics and team performance. The ability to vary interdependence patterns while maintaining other experimental factors would provide unprecedented insights into optimal team configurations for different cyber defense contexts.

Communication infrastructure represents another critical dimension for advancement revealed by my dissertation work. While the current TDG implementation supports basic interaction between humans and agents, future platforms require a more comprehensive communication framework. Such a framework would encompass permission-based workflows for critical actions, explanation mechanisms for agent decision processes, and natural language channels for flexible dialogue. As [152] demonstrate, communication strategies play a pivotal role in team coordination, particularly in high-stakes environments like cybersecurity. This enhanced infrastructure would support investigations into how different communication modalities affect team performance, trust development, and coordination efficiency in cyber defense scenarios where rapid information exchange is crucial.

## 6.4.2 Human Behavior-Aware Agents as Team Members

Beyond platform enhancements, my research has opened significant opportunities for developing more sophisticated autonomous agents specifically designed for effective human-AI collaboration in cyber defense.

A primary direction involves designing agents with capabilities that complement rather than replicate human cognitive strengths. My dissertation has demonstrated that human-AI teams are most effective when their capabilities are complementary instead of redundant. [20] found that complementary team performance improved when AI systems were designed to address specific human cognitive limitations rather than mimic human expertise. Future work should develop a formal framework for capability characterization that identifies the relative strengths of humans and AI in different cyber defense sub-tasks. Such agents would excel at monitoring large volumes of network data without fatigue, detecting subtle correlations across disparate sources, maintaining comprehensive historical context, and rigorously quantifying confidence levels—all capabilities that complement known human cognitive limitations in cyber defense contexts.

Another promising direction involves developing agents capable of participating in ad-hoc teamwork scenarios where humans and systems dynamically join together to address specific cybersecurity incidents without prior coordination. Such scenarios are increasingly common in cybersecurity operations, where cross-organizational teams must rapidly form in response to major security incidents. [230] pioneered research in ad-hoc agent teaming, and extending this to human-AI contexts presents unique challenges. The research involves creating agents whose actions remain understandable to human teammates without extensive training, while dynamically modeling the capabilities of those teammates under limited observation. These agents must adapt their coordination strategies based on emerging team dynamics and balance exploration (learning about teammates) with exploitation (maximizing immediate team performance).

My work with cognitive models based on Instance-Based Learning Theory provides a foundation for agents with enhanced theory of mind capabilities. Such agents would model and predict human teammate behavior, including decisions about reliance and intervention. [175] demonstrated that cognitive models can successfully develop theory of mind capabilities through observation, enabling more accurate prediction of human decision processes. These models would incorporate factors such as trust dynamics, cognitive load, and expertise level to anticipate when humans are likely to appropriately rely on agent recommendations, unnecessarily override agent decisions, require additional explanation, or experience decision fatigue. By accurately modeling human teammates, agents can better adapt their behavior to complement human capabilities and compensate for predictable limitations or biases.

The development of multi-level adaptation mechanisms represents another significant research opportunity. My dissertation work suggests that effective agents must learn from interaction at multiple complementary time scales: rapid adaptation to immediate team dynamics through reinforcement learning, retention of team-specific strategies across multiple interactions with the same teammates, and acquisition of general principles for effective human-AI collaboration that transfer to new teammates and contexts. [196] argue that human learning occurs at multiple timescales, from immediate skill acquisition to long-term conceptual development, and AI systems that mirror this multi-level adaptation may collaborate more effectively with humans. These adaptation mechanisms must balance stability (maintaining consistent, predictable behavior) with flexibility (adjusting to changing team needs) to address a key challenge in human-AI teaming for cybersecurity operations.

## 6.4.3 Benchmarking AI Agents in Human-AI Teaming

Benchmarking autonomous agents for human-AI collaboration represents a foundational research direction that remains underdeveloped in current cybersecurity teaming frameworks. While many existing systems are evaluated against simplistic baselines such as random agents or static heuristics, a growing body of empirical work has demonstrated that the selection and design of benchmarking agents significantly influence not only objective team performance but also human trust calibration, reliance behavior, and overall teaming dynamics.

Recent studies in human-AI collaboration have revealed the limitations of self-play-optimized agents when placed in mixed human-AI teams. For example, in cooperative domains such as Overcooked and Hanabi, agents trained exclusively via reinforcement learning performed well with other agents but failed to coordinate effectively with human partners due to misaligned behavior models and unintuitive policies [33, 223]. In contrast, rule-based or behavior-cloned agents that more closely matched human expectations—though potentially suboptimal in isolated performance—were consistently preferred by human users. These findings underscore the need to benchmark agents not solely by technical metrics but also by their ability to function as teammates in realistic collaborative contexts.

A comprehensive benchmarking framework for human-AI teaming should incorporate multiple types of baseline agents beyond the random or naive comparator. These include: (1) extitscripted agents that reflect domain-specific heuristics or expert-defined protocols; (2) extitablated agents in which key architectural or behavioral components (e.g., explanation modules, human modeling layers) are systematically removed; (3) extitcapability-matched agents whose

task performance is tuned to be comparable with human counterparts to avoid dominance effects; and (4) extithuman-level proxies, such as behaviorally cloned agents trained on human data to provide a reference for naturalistic teaming behavior. The inclusion of these diverse baselines enables more diagnostic evaluation of proposed agents' contributions to team effectiveness.

Equally important is the need to benchmark along theoretically grounded design dimensions known to shape teaming dynamics. These include, but are not limited to: transparency and explanation generation [20, 260], adaptivity to partner behavior [33, 146], coordination strategy and timing [259], and levels of agent autonomy and initiative [205]. For example, slight changes in the wording of explanations—shifting from hedging to confident language—can significantly alter user acceptance and decision-making even when the underlying recommendation remains unchanged. Similarly, the timing of agent errors (e.g., early versus late in an interaction) has been shown to affect long-term trust trajectories. These findings highlight the necessity of rigorous control in benchmark agent design and careful documentation of seemingly minor implementation choices that may have outsized effects on human perception and teaming quality.

The research community has increasingly emphasized several best practices for benchmarking human-AI collaboration. These include using standardized testbeds with reproducible configurations [33], adopting within-subject or counterbalanced user studies to isolate treatment effects [20], and reporting multi-dimensional outcome measures encompassing not only task success but also human-centered metrics such as trust alignment, cognitive load, and perceived cooperativeness [260, 146]. In addition, recent work advocates for the open release of agent code, experimental protocols, and evaluation data to facilitate reproducibility and comparative analysis across research groups.

Cyber defense settings present unique challenges for agent benchmarking. Unlike tabletop games or general assistance tasks, cyber operations are characterized by adversarial dynamics, noisy and ambiguous data, and time-sensitive decision demands. Benchmarking agents in these environments must therefore account for additional factors such as alert fatigue, the interpretability of intrusion detection recommendations, and the robustness of trust calibration under uncertainty and deception. Moreover, the interaction between agent autonomy and human oversight becomes particularly critical when agents are authorized to perform disruptive actions, such as isolating network nodes or blocking communications. Evaluating agents across varied scenarios that simulate these conditions—while holding constant key contextual variables—will be essential to identify designs that generalize across threat landscapes and organizational structures.

Taken together, these considerations underscore that benchmarking agents for human-AI teaming is not merely a technical exercise but a complex methodological challenge. It involves principled design of baseline agents, theoretically motivated variation along critical interaction dimensions, and ecologically valid evaluation methodologies that reflect the demands of real-world cyber defense operations.

## 6.4.4 Comprehensive Team Evaluation

My research has demonstrated that traditional cybersecurity metrics focused solely on technical performance are insufficient for evaluating effective human-AI teams. This opens up research opportunities for developing more nuanced evaluation methodologies that capture the quality of collaboration alongside technical outcomes.

Trust calibration measurement represents a critical research direction. My work has shown that trust is not uniformly beneficial; rather, appropriate trust calibration—trusting the right agent for the right task at the right time—is crucial for team performance. [13] demonstrated that calibrated trust in AI systems significantly improves team performance in deferred decision tasks, particularly when humans appropriately adjust their reliance based on AI capabilities. Future research should develop and validate quantitative measures of trust calibration that assess the alignment between agent capability and human reliance, appropriate adjustment of trust levels based on observed performance, situation-specific modulation of trust based on task characteristics, and resistance to overtrust in high-confidence but incorrect agent recommendations. These measures would enable more sophisticated evaluation of human-AI team dynamics beyond simple trust or distrust dichotomies.

The study of team resilience under adversarial conditions represents another important direction emerged from my work. Cyber defense teams must function effectively even when under attack or facing resource limitations. [153] argue that team resilience in cybersecurity contexts involves not only technical robustness but also adaptive capacity in the face of unexpected challenges. Future research should develop experimental protocols for testing team resilience under adversarial conditions, including scenarios with communication disruption that limit information sharing between team members, resource constraints that force prioritization decisions under time pressure, deception attempts where adversaries try to manipulate team trust relationships, and recovery situations that test the team's ability to reestablish effective coordination after failures. These protocols would provide insight into the robustness of different team configurations and identify specific vulnerabilities that can be addressed through training or agent design.

Cognitive workload distribution assessment has also emerged as a promising research direction. My work suggests that effective teams distribute cognitive workload appropriately across members according to their capabilities and current capacity. [42] demonstrated that effective command-and-control teams dynamically redistribute cognitive burden based on evolving task demands and individual capacity. Future research should develop methods to assess cognitive load distribution in human-AI teams through both behavioral and physiological measures when possible. This research would determine whether human-AI teams successfully offload appropriate cognitive burdens to autonomous agents while keeping humans engaged in decisions that benefit from human judgment and expertise, and whether agents appropriately adjust their autonomy level based on detected human cognitive load.

The research directions outlined in this section represent a comprehensive agenda for advancing human-AI teaming in cyber defense. By simultaneously developing more sophisticated experimental platforms, designing team-aware cognitive agents, implementing nuanced evaluation methodologies, and addressing open-world decision challenges, this agenda would address fundamental questions about effective human-AI collaboration in complex cybersecurity environments. The long-term vision is to develop human-AI teams that achieve performance exceeding what either humans or AI could accomplish individually, while maintaining appropriate human oversight in critical decision contexts. This work would contribute not only to improved cybersecurity operations but also to broader understanding of complementary human-AI collaboration in high-stakes environments.

## Bibliography

- [1] M. A. Abbasi, C. Dovrolis, and K. B. Parag. "Toward Realistic Models of Network Security Games". In: *Journal of Information Security and Applications* 15.1 (2010), pp. 81–100.
- [2] Robert K Abercrombie, Bob G Schlicher, and Frederick T Sheldon. "Security analysis of selected AMI failure scenarios using agent based game theoretic simulation". In: 2014 47th Hawaii International Conference on System Sciences. IEEE. 2014, pp. 2015–2024.
- [3] Ibrahim Adedeji Adeniran et al. "Strategic Risk Management in Financial Institutions: Ensuring Robust Regulatory Compliance". In: *Finance & Accounting Research Journal* (2024). DOI: 10.51594/farj.v6i8.1508.
- [4] P. Aggarwal et al. "Designing Effective Masking Strategies for Cyberdefense through Human Experimentation and Cognitive Models". In: *Computers & Security* 117 (2022), p. 102671.
- [5] Palvi Aggarwal, Cleotilde Gonzalez, and Varun Dutt. "HackIt: A real-time simulation tool for studying real-world cyberattacks in the laboratory". In: *Handbook of Computer Networks and Cyber Security*. Springer, 2020, pp. 949–959.
- [6] Afrah Almansoori, Mostafa Al-Emran, and Khaled Shaalan. "Exploring the Frontiers of Cybersecurity Behavior: A Systematic Review of Studies and Theories". In: *Applied Sciences* (2023). DOI: 10.3390/app13095700.
- [7] T. Alpcan and T. Başar. *Network Security: A Decision and Game-Theoretic Approach*. Cambridge University Press, 2010.
- [8] M. Alsharnouby, F. Alaca, and S. Chiasson. "Why phishing still works: User strategies for combating phishing attacks". In: *International Journal of Human-Computer Studies* 82 (2015), pp. 69–82.
- [9] J. R. Anderson. "How Can the Human Mind Occur in the Physical Universe?" In: (2007).
- [10] John R Anderson and Lael J Schooler. "Reflections of the environment in memory". In: *Psychological Science* 2.6 (1991), pp. 396–408.
- [11] John R. Anderson. *How Can the Human Mind Occur in the Physical Universe?* Oxford Series on Cognitive Models and Architectures. New York: Oxford University Press, 2007.
- [12] James Andreoni and John H Miller. "Rational cooperation in the finitely repeated prisoner's dilemma: Experimental evidence". In: *The economic journal* 103.418 (1993), pp. 570–585.

- [13] Duri Long de Araujo, Finale Doshi-Velez, and Julie Shah. "Calibrated Trust: Enabling Effective Human-Autonomy Teaming in Task Planning". In: *ACM Transactions on Interactive Intelligent Systems* 11.3-4 (2021), pp. 1–29.
- [14] Travis Ashley et al. "Gamification of Cybersecurity for Workforce Development in Critical Infrastructure". In: *Ieee Access* (2022). DOI: 10.1109/access.2022.3216711.
- [15] Robert Axelrod. "The Evolution of Cooperation". In: *Basic Books* (1984).
- [16] Chris L. Baker, Rebecca Saxe, and Joshua B. Tenenbaum. "Bayesian theory of mind: Modeling joint belief-desire attribution". In: *Proceedings of the Annual Meeting of the Cognitive Science Society* 33.33 (2011).
- [17] Tyler Balon and Ibrahim Baggili. "Cybercompetitions: A Survey of Competitions, Tools, and Systems to Support Cybersecurity Education". In: *Education and Information Technologies* (2023). DOI: 10.1007/s10639-022-11451-4.
- [18] T. Ban et al. "Combat Security Alert Fatigue with AI-Assisted Techniques". In: *Proceedings of the 14th Cyber Security Experimentation and Test Workshop* (2021), pp. 9–16.
- [19] Gagan Bansal et al. "Beyond accuracy: The role of mental models in human-AI team performance". In: *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing* 7.1 (2019), pp. 2–11.
- [20] Gagan Bansal et al. "Does the whole exceed its parts? The effect of AI explanations on complementary team performance". In: *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 2021, pp. 1–16.
- [21] Mohan Baruwal Chhetri et al. "Towards Minimizing Alert Fatigue: An XAI-based Decision Support Tool for SOC Analysts". In: *arXiv preprint arXiv:2208.03026* (2022).
- [22] Mohan Baruwal Chhetri et al. "Towards reducing cognitive load of SOC analysts via human-centric alert management". In: *Digital Threats: Research and Practice* 4.2 (2023), pp. 1–18.
- [23] Zisis Batzos et al. "Gamification and Serious Games for Cybersecurity Awareness and First Responders Training: An Overview". In: (2023). DOI: 10.36227/techrxiv. 22650952.
- [24] N. Ben-Asher and C. Gonzalez. "Effects of Cyber Security Knowledge on Attack Detection". In: *Computers in Human Behavior* 48 (2015), pp. 51–61.
- [25] Terry Benzel. "The science of cyber security experimentation: The DETER project". In: *Proceedings of the 27th Annual Computer Security Applications Conference*. 2010, pp. 137–148.
- [26] Pieter van den Berg and Franz J. Weissing. "Cooperative behavior in social environments". In: *Current Opinion in Behavioral Sciences* 34 (2020), pp. 127–134.
- [27] Oscar Bjurling et al. "Enabling Human-Autonomy Teaming in Aviation: A Framework to Address Human Factors in Digital Assistants Design". In: *Journal of Physics Conference Series* (2024). DOI: 10.1088/1742-6596/2716/1/012076.
- [28] N. Buchler et al. "Mission Command in the Age of Network-Enabled Operations: Social Network Analysis of Information Sharing and Situation Awareness". In: *Frontiers in Psychology* 7 (2016), p. 937.

- [29] Nico Büchler et al. "Learning in the context of real-time strategic interaction". In: *Nature Communications* 2 (2011), p. 209.
- [30] Sergio Caltagirone, Andrew Pendergast, and Christopher Betz. "The diamond model of intrusion analysis". In: *Center for Cyber Intelligence Analysis and Threat Research* (2023).
- [31] Valerio Capraro and Christoph Hauert. "A group fitness theory for cooperation in difficult tasks". In: *Scientific Reports* 3 (2013), p. 2013.
- [32] David Carmel and Shaul Markovitch. "Opponent modeling in multi-agent systems". In: *IJCAI*. Vol. 95. 1995, pp. 40–45.
- [33] Micah Carroll et al. "On the Utility of Learning About Humans for Human-AI Coordination". In: *Advances in Neural Information Processing Systems (NeurIPS)*. Vol. 32. 2019.
- [34] M. A. Champion et al. "Team-Based Cyber Defense Analysis". In: *IEEE International Multi-Disciplinary Conference on Cognitive Methods in Situation Awareness and Decision Support* (2012), pp. 218–221.
- [35] Marie Chapelle et al. "Adaptive Cybersecurity: Machine Learning-Based Intrusion Detection and Response Systems". In: *Applied Sciences* 13.2 (2023), p. 983.
- [36] Gabriele Chierchia et al. "Integrating social and nonsocial cognitive control". In: *Cognition* 163 (2017), pp. 42–60.
- [37] Erin K Chiou and John D Lee. "Negotiated and reciprocal exchange structures in humanagent cooperation". In: *Computers in Human Behavior* 90 (2019), pp. 288–297.
- [38] Noman H. Chowdhury, Marc T. P. Adam, and Geoffrey Skinner. "The impact of time pressure on cybersecurity behaviour: a systematic literature review". In: *Behaviour & Information Technology* 38.12 (2019), pp. 1290–1308. DOI: 10.1080/0144929X.2019. 1583769.
- [39] Paul F Christiano et al. "Deep reinforcement learning from human preferences". In: *Advances in Neural Information Processing Systems* 30 (2017).
- [40] Merijke Coenraad et al. "Experiencing Cybersecurity One Game at a Time: A Systematic Review of Cybersecurity Digital Games". In: *Simulation & Gaming* (2020). DOI: 10. 1177/1046878120933312.
- [41] Edward JM Colbert, Alexander Kott, and Lawrence P Knachel. "The game-theoretic model and experimental investigation of cyber wargaming". In: *The Journal of Defense Modeling and Simulation* 17.1 (2020), pp. 21–38.
- [42] Nancy J Cooke et al. "Team cognition in experienced command-and-control teams". In: *Journal of Experimental Psychology: Applied* 19.3 (2013), p. 233.
- [43] Jacob W Crandall et al. "Cooperating with machines". In: *Nature Communications* 9.1 (2018), pp. 1–12.
- [44] E. A. Cranford et al. "Modeling Cognitive Dynamics in (End)-user Response to Phishing Emails". In: *Proceedings of the 17th ICCM* (2019).
- [45] E. A. Cranford et al. "Toward Personalized Deceptive Signaling for Cyber Defense Using Cognitive Models". In: *Topics in Cognitive Science* 12.3 (2020), pp. 992–1011.

- [46] L. F. Cranor. "A Framework for Reasoning About the Human in the Loop". In: *Proceedings of the 1st Conference on Usability, Psychology, and Security.* 2008, pp. 1–15.
- [47] Sourya Joyee De and Daniel Le Métayer. "PRIAM: A Privacy Risk Analysis Methodology". In: (2016). DOI: 10.1007/978-3-319-47072-6\\_15.
- [48] Mustafa Demir, Nathan J. McNeese, and Nancy J. Cooke. "The Evolution of Human-Autonomy Teams in Remotely Piloted Aircraft Systems Operations". In: *Frontiers in Communication* (2019). DOI: 10.3389/fcomm.2019.00050.
- [49] Gelei Deng et al. "{PentestGPT}: Evaluating and harnessing large language models for automated penetration testing". In: *33rd USENIX Security Symposium (USENIX Security 24)*. 2024, pp. 847–864.
- [50] Mina Deng et al. "A Privacy Threat Analysis Framework: Supporting the Elicitation and Fulfillment of Privacy Requirements". In: *Requirements Engineering* (2010). DOI: 10.1007/s00766-010-0115-7.
- [51] A. R. Dennis and R. K. Minas. "Security on Autopilot: Why Current Security Theories Hijack Our Thinking and Lead Us Astray". In: *ACM SIGMIS Database: The DATABASE for Advances in Information Systems* 49.1 (2018), pp. 15–38. DOI: 10.1145/3184444. 3184448.
- [52] Fábio Martins Dias et al. "Risk Management Focusing on the Best Practices of Data Security Systems for Healthcare". In: *International Journal of Innovation* (2021). DOI: 10.5585/iji.v9i1.18246.
- [53] Qing Dong and Ce Pang. "Risk-Based Non-Myopic Sensor Scheduling in Target Threat Level Assessment". In: *Ieee Access* (2021). DOI: 10.1109/access.2021.3078830.
- [54] Yinuo Du, Baptiste Prébot, and Cleotilde González. "Turing-Like Experiment in a Cyber Defense Game". In: *AAAI Spring Symposium Series* (2024). DOI: 10.1609/aaaiss. v3i1.31271.
- [55] Yinuo Du et al. "A Cyber-War Between Bots: Cognitive Attackers are More Challenging for Defenders than Strategic Attackers". In: *ACM Transactions on Social Computing* 8.3-4 (2025), pp. 1–22.
- [56] Yinuo Du et al. "Learning to play an adaptive cyber deception game". In: *Proc. of the 21st International Conference on Autonomous Agents and Multiagent Systems. Auckland, New Zealand.* Vol. 6. 2022.
- [57] Yinuo Du et al. "Towards Autonomous Cyber Defense: Predictions From a Cognitive Model". In: *Proceedings of the Human Factors and Ergonomics Society Annual Meeting* (2022). DOI: 10.1177/1071181322661504.
- [58] Yinuo Du et al. "Towards autonomous cyber defense: predictions from a cognitive model". In: *Proceedings of the Human Factors and Ergonomics Society Annual Meeting* (2022).
- [59] V. Dutt, Y.-S. Ahn, and C. Gonzalez. "Cyber Situation Awareness: Modeling Detection of Cyber Attacks With Instance-Based Learning Theory". In: *Human Factors* 55.3 (2013), pp. 605–618.
- [60] V. Dutt, Y.-S. Ahn, and C. Gonzalez. "Cyber Situation Awareness: Modeling the Security Analyst in a Cyber-attack Scenario through Instance-based Learning". In: *Data and Applications Security and Privacy XXV* (2011), pp. 280–292.

- [61] Varun Dutt, Young-Suk Ahn, and Cleotilde Gonzalez. "Cyber Situation Awareness: Modeling the security analyst in a cyber-Attack Scenario through Instance-Based Learning". In: *Data and Applications Security and Privacy XXV*. Ed. by Yingjiu Li. Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 280–292. ISBN: 978-3-642-22348-8.
- [62] J. Dykstra and C. L. Paul. "Cyber Operations Stress Survey (COSS): Studying Fatigue, Frustration, and Cognitive Workload in Cybersecurity Operations". In: 11th USENIX Workshop on Cyber Security Experimentation and Test (CSET 18) (2018).
- [63] Ido Erev and Alvin E. Roth. "Learning and equilibrium as useful approximations: Accuracy of prediction on randomly selected constant sum games". In: *Economic Theory* 33.1 (2007), pp. 29–51.
- [64] Ernst Fehr and Simon Gächter. "Cooperation and Punishment in Public Goods Experiments". In: *American Economic Review* 90.4 (2000), pp. 980–994. DOI: 10.1257/aer. 90.4.980.
- [65] Ernst Fehr and Ivo Schurtenberger. "Cooperative phenotypes predict reciprocal behavior". In: *Science* 365.6457 (2019).
- [66] Kimberly Ferguson-Walter et al. "Oppositional Human Factors in Cybersecurity: A Preliminary Analysis of Affective States". In: Nov. 2021. DOI: 10.1109/ASEW52652.2021. 00040.
- [67] Kimberly Ferguson-Walter et al. "Oppositional human factors in cybersecurity: A preliminary analysis of affective states". In: 2021 IEEE Aerospace Dependable and Secure Computing Conference (ADSCC). IEEE, 2021, pp. 1–9.
- [68] Kimberly Ferguson-Walter et al. "The Tularosa study: An experimental design and implementation to quantify the effectiveness of cyber deception". In: *Hawaii International Conference on System Sciences* (2018).
- [69] Kimberly Ferguson-Walter et al. *The Tularosa Study: An Experimental Design and Implementation to Quantify the Effectiveness of Cyber Deception*. Tech. rep. Sandia National Lab.(SNL-NM), Albuquerque, NM (United States), 2018.
- [70] Drew Fudenberg and Eric Maskin. "The folk theorem in repeated games with discounting or with incomplete information". In: *Econometrica* (1986), pp. 533–554.
- [71] R. Garrido-Pelaz, L. González-Manzano, and S. Pastrana. "Shall we collaborate? A model to analyse the benefits of information sharing". In: *arXiv preprint arXiv:1606.02539* (2016).
- [72] Oliver Genschow et al. "Psychological Distance Modulates Goal-Based Versus Movement-Based Imitation." In: *Journal of Experimental Psychology Human Perception & Performance* (2019). DOI: 10.1037/xhp0000654.
- [73] Ali Gholami et al. "Design and Implementation of the Advanced Cloud Privacy Threat Modeling". In: *International Journal of Network Security & Its Applications* (2016). DOI: 10.5121/ijnsa.2016.8207.
- [74] G. Gigerenzer and R. Selten. "Bounded Rationality: The Adaptive Toolbox". In: *MIT Press* (2001).
- [75] Gerd Gigerenzer and Wolfgang Gaissmaier. "Heuristic decision making". In: *Annual Review of Psychology* 62 (2011), pp. 451–482.

- [76] Ella Glikson and Anita Williams Woolley. "Human trust in artificial intelligence: Review of empirical research". In: *Academy of Management Annals* 14.2 (2020), pp. 627–660.
- [77] Piotr J Gmytrasiewicz and Prashant Doshi. "A framework for reasoning about others: The theory of mind approach". In: *Journal of Artificial Intelligence Research* 24 (2005), pp. 1–35.
- [78] C. Gonzalez. "Building Human-like Artificial Agents: A General Cognitive Algorithm for Emulating Human Decision-making in Dynamic Environments". In: *Perspectives on Psychological Science* 19.5 (2023), pp. 860–873.
- [79] C. Gonzalez et al. "Design of Dynamic and Personalized Deception: A Research Framework and New Insights". In: (2020), pp. 1825–1834.
- [80] Cleotilde Gonzalez. "Building Human-Like Artificial Agents: A General Cognitive Algorithm for Emulating Human Decision-Making in Dynamic Environments". In: *Perspectives on Psychological Science* (2023), p. 17456916231196766.
- [81] Cleotilde Gonzalez, Javier F Lerch, and Christian Lebiere. "Instance-based learning in dynamic decision making". In: *Cognitive Science* 27.4 (2003), pp. 591–635.
- [82] Cleotilde Gonzalez and Jolie M Martin. "Scaling up instance-based learning theory to account for social interactions". In: *Negotiation and Conflict Management Research* 4.2 (2011), pp. 110–128.
- [83] Cleotilde Gonzalez et al. "A cognitive model of dynamic cooperation with varied interdependency information". In: *Cognitive science* 39.3 (2015), pp. 457–495.
- [84] Cleotilde Gonzalez et al. "Cognition and technology". In: *Cyber Defense and Situational Awareness*. Ed. by Alexander Kott, Cliff Wang, and Robert F. Erbacher. Cham: Springer International Publishing, 2014, pp. 93–117. ISBN: 978-3-319-11391-3. DOI: 10.1007/978-3-319-11391-3\_6. URL: 10.1007/978-3-319-11391-3\_6.
- [85] M. Granåsen and D. Andersson. "Measuring Team Effectiveness in Cyber-Defense Exercises: A Cross-Disciplinary Case Study". In: *Cognition, Technology & Work* 19.2-3 (2017), pp. 329–345.
- [86] Jelena Grujić and Tom Lenaerts. "Social dynamics within decomposed games: giving teammates a helping hand in public goods games". In: *Physical Review E* 89.2 (2014), p. 023811.
- [87] Oliver Guidetti, Craig Speelman, and Peter Bouhlas. "A Review of Cyber Vigilance Tasks for Network Defense". In: *Frontiers in Neuroergonomics* (2023). DOI: 10.3389/fnrgo. 2023.1104873.
- [88] Oliver Guidetti, Craig Speelman, and Peter Bouhlas. "The WACDT, a Modern Vigilance Task for Network Defense". In: *Frontiers in Neuroergonomics* (2023). DOI: 10.3389/fnrgo.2023.1215497.
- [89] T. Gundu. "Understanding Cyber Security Skills in a Dynamic Systems Environment". In: SSRN Electronic Journal (2020).
- [90] Tapiwa Gundu. "Learn, Unlearn and Relearn: Adaptive Cybersecurity Culture Model". In: *International Conference on Cyber Warfare and Security* (2024). DOI: 10.34190/iccws.19.1.2177.

- [91] Jayesh K. Gupta, Maxim Egorov, and Mykel J. Kochenderfer. "Cooperative Multi-agent Control Using Deep Reinforcement Learning". In: *Autonomous Agents and Multiagent Systems. AAMAS 2017. Workshops, Best Papers, Revised Selected Papers.* Ed. by Gita Sukthankar and Juan A. Rodriguez-Aguilar. Vol. 10642. Lecture Notes in Computer Science. Springer, 2017, pp. 66–83. DOI: 10.1007/978-3-319-71682-4\_5.
- [92] R. S. Gutzwiller et al. "Oh, Look, A Butterfly! A Framework for Distracting Attackers to Improve Cyber Defense". In: *Human Factors and Ergonomics Society Annual Meeting Proceedings* 64.1 (2020), pp. 437–441.
- [93] Robert S Gutzwiller, Benjamin A Clegg, and C.A.P. Smith. "Part Task Training in the Context of Automation: Current and Future Directions". In: *ARL Technical Report* (2013).
- [94] Anderson Kevin Gwenhure and Flourensia Sapty Rahayu. "Gamification of Cybersecurity Awareness for Non-It Professionals: A Systematic Literature Review". In: *International Journal of Serious Games* (2024). DOI: 10.17083/ijsg.v11i1.719.
- [95] Lee Hadlington. "Employees attitude towards cyber security and risky online behaviours: An empirical assessment in the United Kingdom". In: *International Journal of Cyber Criminology* 12.1 (2018), pp. 269–283. DOI: 10.5281/zenodo.1467909.
- [96] Samuel N Hamilton and Wendy L Hamilton. "Adversary modeling and simulation in cyber warfare". In: *IFIP International Information Security Conference*. Springer. 2008, pp. 461–475.
- [97] S. T. Hamman et al. "Represented Model of a Concrete Cyber Security Training Game". In: *International Journal of Cyber Warfare and Terrorism* 7.3 (2017), pp. 42–56.
- [98] B. P. Hámornik and C. Krasznay. "A team-level perspective of human factors in cyber security: Security operations centers". In: *International Conference on Applied Human Factors and Ergonomics*. 2017, pp. 224–236.
- [99] Balázs Péter Hámornik and Csaba Krasznay. "A team-level perspective of human factors in cyber security: security operations centers". In: *International Conference on Applied Human Factors and Ergonomics*. Springer. 2017, pp. 224–236.
- [100] Jochim Hansen, Hans Alves, and Yaacov Trope. "Psychological Distance Reduces Literal Imitation: Evidence From an Imitation-Learning Paradigm." In: *Journal of Experimental Psychology Human Perception & Performance* (2016). DOI: 10.1037/xhp0000150.
- [101] Marc Harper et al. "Reinforcement learning produces dominant strategies for the Iterated Prisoner's Dilemma". In: *PloS One* 12.12 (2017), e0188046.
- [102] Jayde Harris et al. "Social schema for human-autonomy teams: Impact of perceived warmth and competence on trust and automation use". In: *Human Factors* 65.7 (2023), pp. 1301–1317.
- [103] Fritz Heider. "Attitudes and cognitive organization". In: *The Journal of Psychology* 21.1 (1946), pp. 107–112. DOI: 10.1080/00223980.1946.9917275.
- [104] M. Hendrix, A. Al-Sherbaz, and V. Bloom. "Game Based Cyber Security Training: Are Serious Games Suitable for Cyber Security Training?" In: *International Journal of Serious Games* 3.1 (2016), pp. 53–61.
- [105] Pablo Hernandez-Leal et al. "A survey of learning in multiagent environments: Dealing with non-stationarity". In: *arXiv preprint arXiv:1709.02779* (2019).

- [106] Ralph Hertwig and Ido Erev. "Decisions from experience and statistical probabilities: Why they trigger different choices than a priori probabilities". In: *Journal of Behavioral Decision Making* 22.1 (2009), pp. 1–35.
- [107] Cecilia Heyes. "What's Social About Social Learning?" In: *Journal of Comparative Psychology* (2012). DOI: 10.1037/a0025180.
- [108] Junjie Huang and Quanyan Zhu. "Penheal: A two-stage llm framework for automated pentesting and optimal remediation". In: *Proceedings of the Workshop on Autonomous Cybersecurity*. 2023, pp. 11–22.
- [109] Lisa M. Huang and Jeffrey W. Sherman. "Attentional Processes in Social Perception". In: (2018). DOI: 10.1016/bs.aesp.2018.03.002.
- [110] S.-W. Huang et al. "From Non-discrimination to Diversity: A Study on Anti-discrimination Policies in Online Marketplaces". In: *Proceedings of the ACM on Human-Computer Interaction* 2 (2018), pp. 1–28.
- [111] Wei Huang, Sarah Goldstein, and Ehab Al-Shaer. "Beyond Automation: Human Creativity in Exploit Chain Development". In: *Proceedings of the 2022 IEEE Symposium on Security and Privacy (SP)*. San Francisco, CA, USA: IEEE, 2022, pp. 278–295. DOI: 10.1109/SP46214.2022.9833647.
- [112] Mark I. Hwang. "Decision Making Under Time Pressure: A Model for Information Systems Research". In: *Information & Management* 27.4 (1994), pp. 197–203. DOI: 10.1016/0378-7206(94)90048-5.
- [113] C. E. Irvine, M. F. Thompson, and K. Allen. "CyberCIEGE: Gaming for Information Assurance". In: *IEEE Security & Privacy* 3.3 (2005), pp. 61–64.
- [114] Ge Jin et al. "Evaluation of Game-Based Learning in Cybersecurity Education for High School Students". In: *Journal of Education and Learning (Edulearn)* (2018). DOI: 10. 11591/edulearn.v12i1.7736.
- [115] Chelsea K Johnson et al. "Decision-Making Biases and Cyber Attackers". In: 2021 36th IEEE/ACM International Conference on Automated Software Engineering Workshops (ASEW). IEEE. 2021, pp. 140–144.
- [116] David A. Kalkstein et al. "Social Learning Across Psychological Distance." In: *Journal of Personality and Social Psychology* (2016). DOI: 10.1037/pspa0000042.
- [117] Hamdi Kavak et al. "Simulation for cybersecurity: state of the art and future directions". In: *Journal of Cybersecurity* 7.1 (2021), tyab005.
- [118] Harold H Kelley and Anthony J Stahelski. "Social interaction basis of cooperators' and competitors' beliefs about others." In: *Journal of personality and social psychology* 16.1 (1970), p. 66.
- [119] T. D. Kelley, E. Avery, and L. N. Long. "A Hybrid Cognitive Architecture with Primal Affect and Physiology". In: *IEEE Transactions on Autonomous Mental Development* 5.2 (2013), pp. 109–118.
- [120] F. Kheiri, C. Gonzalez, and C. Lebiere. "Learning Collaborative Strategies in Cyber Defense through Cognitive Modeling". In: *Journal of Cybersecurity* 8.1 (2022), pp. 1–15.

- [121] Elmar Kiesling et al. "Simulation-based optimization of information security controls: An adversary-centric approach". In: *2013 Winter Simulations Conference (WSC)*. IEEE. 2013, pp. 2054–2065.
- [122] Oliver Kirchkamp et al. "Behavioral spillovers and cooperation". In: *Journal of Economic Behavior & Organization* 130 (2016), pp. 160–175.
- [123] Max Kleiman-Weiner et al. "Coordinate to cooperate or compete: Abstract goals and joint intentions in social interaction". In: *Topics in Cognitive Science* 8.2 (2016), pp. 424–435.
- [124] Max Kleiman-Weiner et al. "Learning to cooperate: The evolution of social rewards in repeated interactions". In: *Topics in Cognitive Science* 8.4 (2016), pp. 860–877.
- [125] G. Klein. "Naturalistic Decision Making". In: *Human Factors* 50.3 (2008), pp. 456–460.
- [126] Gary Klein. *Sources of Power: How People Make Decisions*. Cambridge, MA: MIT Press, 1998. ISBN: 978-0262611466.
- [127] Gary Klein. The Power of Intuition: How to Use Your Gut Feelings to Make Better Decisions at Work. New York, NY: Currency/Doubleday, 2004. ISBN: 978-0385502894.
- [128] Gary A. Klein et al., eds. *Decision Making in Action: Models and Methods*. Norwood, NJ: Ablex Publishing Corporation, 1993. ISBN: 978-0893917432.
- [129] Oleksandr Korystin and Nataliia Svyrydiuk. "Activities of Illegal Weapons Criminal Component of Hybrid Threats". In: (2021). DOI: 10.2991/aebmr.k.210320.016.
- [130] R. Kosmider, J. C. Gibbens, and Rachelle Avigad. "Identification, Assessment and Management of New and Re-emerging Animal-related Risks: UK Perspective". In: *Veterinary Record* (2017). DOI: 10.1136/vr.104258.
- [131] Alexander Kott et al. "Autonomous Intelligent Agents for Cybersecurity: Opportunities and Challenges". In: *Computer* 56.5 (2023), pp. 81–94.
- [132] Sara Kraemer, Pascale Carayon, and John Clem. "Human and organizational factors in computer and information security: Pathways to vulnerabilities". In: *Computers & Security* 28.7 (2009), pp. 509–520.
- [133] Mona Kriesten, Mamello Thinyane, and David Ormrod. "Leveraging Gamification for Cyber Threat Intelligence for Resilience in Satellite Cyber Supply Chains". In: *European Conference on Cyber Warfare and Security* (2024). DOI: 10.34190/eccws.23.1.2203.
- [134] Halima Ibrahim Kure, Shareeful Islam, and Mohammad A. Razzaque. "An Integrated Cyber Security Risk Management Approach for a Cyber-Physical System". In: *Applied Sciences* (2018). DOI: 10.3390/app8060898.
- [135] Marc Lanctot et al. "A unified game-theoretic approach to multiagent reinforcement learning". In: *Advances in Neural Information Processing Systems*. 2017, pp. 4193–4206.
- [136] Pat Langley. "The computational gauntlet of human-like learning". In: *Proceedings of the aaai conference on artificial intelligence*. Vol. 36. 11. 2022, pp. 12268–12273.
- [137] R. Langner. "Stuxnet: Dissecting a Cyberwarfare Weapon". In: *IEEE Security & Privacy* 9.3 (2011), pp. 49–51.
- [138] C. Lebiere et al. "A Functional Model of Sensemaking in a Neurocognitive Architecture". In: *Computational Intelligence and Neuroscience* (2013).
- [139] Chulgoo Lee, Byungho Kim, and Myongjin Kim. "Mitigating the Impact of Work Overload on Cybersecurity Behavior: The Moderating Influence of Corporate Ethics—A

- Mediated Moderation Analysis". In: *Sustainability* 15.19 (2023), p. 14327. doi: 10. 3390/su151914327.
- [140] Joel Z Leibo et al. "Multi-agent Reinforcement Learning in Sequential Social Dilemmas". In: *Proceedings of AAMAS* (2017).
- [141] Tao Li and Quanyan Zhu. "Symbiotic game and foundation models for cyber deception operations in strategic cyber warfare". In: *arXiv preprint arXiv:2403.10570* (2024).
- [142] Yong Li et al. "Research on New Power Business Threat Modeling Method Based on Bayesian Networks". In: (2024). DOI: 10.1117/12.3033224.
- [143] Shihui Lim et al. "Opponent modeling in deep reinforcement learning". In: *International Conference on Machine Learning* (2016).
- [144] Marianne Lindroth et al. "Advancing Cybersecurity Through Civic Skills". In: *International Journal of Information Security and Cybercrime* (2024). DOI: 10.19107/ijisc. 2024.01.01.
- [145] Igor Linkov, Alexander Kott, and Stephen Cauffman. "Toward a More Adaptive Cyber Defense: Clarifying Terms". In: *IEEE Security & Privacy* 21.1 (2023), pp. 76–80.
- [146] Ziyue Liu et al. "A Hierarchical Approach to Population Training for Human-AI Collaboration". In: *Proceedings of the 32nd International Joint Conference on Artificial Intelligence (IJCAI)*. 2023.
- [147] Andrew J. Lohn et al. "Autonomous Cyber Defense". In: (2023). DOI: 10.51593/2022ca007.
- [148] Michele Loi and Markus Christen. "Ethical Frameworks for Cybersecurity". In: (2020). DOI: 10.1007/978-3-030-29053-5\\_4.
- [149] Ryan Lowe et al. "Multi-agent actor-critic for mixed cooperative-competitive environments". In: *Advances in Neural Information Processing Systems*. 2017, pp. 6379–6390.
- [150] Craig Lyn, Alan R Chappell, and Jeremy McEver. "Opportunities for Autonomy in Tactical Cyber Operations". In: *MITRE Technical Report* (2019).
- [151] Joseph B. Lyons et al. "Human–Autonomy Teaming: Definitions, Debates, and Directions". In: *Frontiers in Psychology* (2021). DOI: 10.3389/fpsyg.2021.589585.
- [152] Jean MacMillan, Elliot E Entin, and Daniel Serfaty. "Communication overhead: The hidden cost of team cognition". In: *Team cognition: Understanding the factors that drive process and performance* (2004), pp. 61–82.
- [153] Janet L Manaois and Douglas L Van Bossuyt. "Resilience Engineering in Secure and Dependable Cyber-Physical Systems". In: *Systems* 10.4 (2022), p. 136.
- [154] V. F. Mancuso, G. J. Funke, and A. J. Strang. "Seeing the Big Picture: Investigating Antecedents of Situational Awareness in Cyber Defense". In: *Proceedings of the Human Factors and Ergonomics Society Annual Meeting* 58.1 (2014), pp. 280–284.
- [155] V. F. Mancuso et al. "Human Factors of Cyber Attacks: A Framework for Human-Centered Research". In: *Proceedings of the Human Factors and Ergonomics Society Annual Meeting* 56.1 (2012), pp. 399–403.
- [156] Jolie M Martin et al. "A description–experience gap in social interactions: Information about interdependence and its effects on cooperation". In: *Journal of Behavioral Decision Making* 27.4 (2014), pp. 349–362.

- [157] Nathan J McNeese et al. "Coactive emergence theory: A foundation for understanding human-autonomy team cognition". In: *Proceedings of the Human Factors and Ergonomics Society Annual Meeting* 65.1 (2021), pp. 427–431.
- [158] Nathan J McNeese et al. "Teaming with a synthetic teammate: Insights into human-autonomy teaming". In: *Human Factors* 60.2 (2018), pp. 262–273.
- [159] Nathan J. McNeese et al. "Understanding the Role of Trust in Human-Autonomy Teaming". In: (2019). DOI: 10.24251/hicss.2019.032.
- [160] James J Meadows and Samuel Sambasivam. "Mandatory Gamified Security Awareness Training Impacts on Texas Public Middle School Students: A Qualitative Study". In: *Issues in Informing Science and Information Technology* (2023). DOI: 10.28945/5129.
- [161] A. Mermoud et al. "To Share or Not to Share: A Behavioral Perspective on Human Participation in Security Information Sharing". In: *Journal of Cybersecurity* 5.1 (2019), pp. 1–13.
- [162] Alain Mermoud et al. "To share or not to share: a behavioral perspective on human participation in security information sharing". In: *Journal of Cybersecurity* 5.1 (2019), tyz006.
- [163] Andreas Michael, Laila Nockur, and Thomas Schlösser. "Prosocial behavior in the time of COVID-19: The effect of private and public role models". In: *Behavioral and Experimental Economics* 89 (2020), pp. 101–124.
- [164] George A Miller. "The magical number seven, plus or minus two: Some limits on our capacity for processing information". In: *Psychological Review* 63.2 (1956), pp. 81–97.
- [165] F. Moisan et al. "Not All Prisoner's Dilemma Games Are Equal: Incentives, Social Preferences, and Cooperation". In: *Decision* 5.4 (2018), pp. 306–322.
- [166] Frédéric Moisan and Cleotilde Gonzalez. "Security under uncertainty: Adaptive attackers are more challenging to human defenders than random attackers". In: *Frontiers in Psychology* 8 (2017), p. 982.
- [167] Frédéric Moisan and Cleotilde Gonzalez. "Security under uncertainty: adaptive attackers are more challenging to human defenders than random attackers". In: *Frontiers in psychology* 8 (2017), p. 982.
- [168] Frédéric Moisan et al. "Not all Prisoner's Dilemma games are equal: Incentives, social preferences, and cooperation." In: *Decision* 5.4 (2018), p. 306.
- [169] S. Monleon, C. Gonzalez, and C. Wiese. "Modeling the Dynamics of Information Sharing in Cyber Defense". In: *Computational and Mathematical Organization Theory* 29.1 (2023), pp. 75–98.
- [170] José AG Moreira et al. "Social dilemmas and cooperation". In: *Physics of Life Reviews* 10.4 (2013), pp. 208–243.
- [171] A. Nagarajan et al. "Exploring Game Design for Cybersecurity Training". In: *IEEE International Conference on Cyber Technology in Automation, Control, and Intelligent Systems* (2012), pp. 256–262.
- [172] Catherine Neubauer et al. "Developing a New Human-Autonomy Team Cohesion Scale". In: *Proceedings of the Human Factors and Ergonomics Society Annual Meeting* (2021). DOI: 10.1177/1071181321651324.

- [173] A. Newell. *Unified Theories of Cognition*. Harvard University Press, 1994.
- [174] T. N. Nguyen and C. Gonzalez. "Theory of Mind from Observation in Cognitive Models and Humans". In: *Topics in Cognitive Science* 14.4 (2022), pp. 665–686.
- [175] Thuy Ngoc Nguyen and Cleotilde Gonzalez. "Theory of mind from observation in cognitive models and humans". In: *Topics in cognitive science* 14.3 (2022), pp. 665–686.
- [176] C. Nobles. "Stress, Burnout, and Security Fatigue in Cybersecurity: A Human Factors Problem". In: *HOLISTICA–Journal of Business and Public Administration* 13 (2022), pp. 49–72.
- [177] Calvin Nobles et al. "Stress, Fatigue, and Cognitive Load in Cybersecurity Operations". In: *Computers & Security* (2022), p. 102706.
- [178] Martin A Nowak. "Five rules for the evolution of cooperation". In: *Science* 314.5805 (2006), pp. 1560–1563.
- [179] Martin A. Nowak. "The evolution of cooperation". In: *Scientific American* 307.1 (2010), pp. 34–39.
- [180] Tom O'Neill et al. "Human–Autonomy Teaming: A Review and Analysis of the Empirical Literature". In: *Human Factors the Journal of the Human Factors and Ergonomics Society* (2020). DOI: 10.1177/0018720820960865.
- [181] S. Y. Oh et al. "Best Practices for Implementing Adaptive Cyber Defense". In: *Proceedings of the 2019 IEEE International Conference on Big Data*. 2019, pp. 3529–3538.
- [182] Sang Ho Oh et al. "Applying Reinforcement Learning for Enhanced Cybersecurity Against Adversarial Simulation". In: *Sensors* (2023). DOI: 10.3390/s23063000.
- [183] Sang Ho Oh et al. "Employing Deep Reinforcement Learning to Cyber-Attack Simulation for Enhancing Cybersecurity". In: *Electronics* (2024). DOI: 10.3390/electronics13030555.
- [184] Serghei Ohrimenco and Valeriu Cernei. "Cybersecurity Risk". In: (2024). DOI: 10. 53486/escst2023.17.
- [185] Raja Parasuraman, Mustapha Mouloua, and Robert Molloy. "Performance consequences of automation-induced complacency". In: *The International Journal of Aviation Psychology* 3.1 (1993), pp. 1–23.
- [186] Christopher Paul and Kathleen Whitley. "Human-computer interaction in cyber security operations". In: *Advances in Information Security* 5 (2014), pp. 117–130.
- [187] Corina Pelau, Dan-Cristian Dabija, and Ionut Nica. "What makes a chatbot human-like? A comparative analysis of human-like characteristics in AI interfaces". In: *Computers in Human Behavior* 124 (2021), p. 106871.
- [188] Matjaž Perc et al. "Statistical physics of human cooperation". In: *Physics Reports* 687 (2017), pp. 1–51.
- [189] Alexander Peysakhovich and Adam Lerer. "Consequentialist learning in repeated interactions". In: *AAMAS* (2018).
- [190] Shari Lawrence Pfleeger and Deanna D. Caputo. "Leveraging behavioral science to mitigate cyber security risk". In: *Computers & Security* 31.4 (2012), pp. 597–611. DOI: 10.1016/j.cose.2011.12.010.

- [191] Dhanya Pramod. "Gamification in Cybersecurity Education; A State of the Art Review and Research Agenda". In: *Journal of Applied Research in Higher Education* (2024). DOI: 10.1108/jarhe-02-2024-0072.
- [192] Baptiste Prebot, Yinuo Du, and Cleotilde Gonzalez. "Learning about simulated adversaries from human defenders using interactive cyber-defense games". In: *Journal of Cybersecurity* 9.1 (2023).
- [193] William H Press and Freeman J Dyson. "Iterated Prisoner's Dilemma contains strategies that dominate any evolutionary opponent". In: *PNAS* 109.26 (2012).
- [194] Charis Psaltis, Gerard Duveen, and Anne-Nelly Perret-Clermont. "The Social and the Psychological: Structure and Context in Intellectual Development". In: *Human Development* (2009). DOI: 10.1159/000233261.
- [195] P. Pusey, M. Gondree, and Z. Peterson. "The Outcomes of Cybersecurity Competitions and Implications for Underrepresented Populations". In: *IEEE Security & Privacy* 14.6 (2016), pp. 90–95.
- [196] Maithra Raghu et al. "Can deep reinforcement learning solve Erdos-Selfridge-Spencer games?" In: *International Conference on Machine Learning* (2017), pp. 2954–2963.
- [197] Fiza Abdul Rahim et al. "Cybersecurity Vulnerabilities in Smart Grids With Solar Photovoltaic: A Threat Modelling and Risk Assessment Approach". In: *International Journal of Sustainable Construction Engineering Technology* (2023). DOI: 10.30880/ijscet. 2023.14.03.018.
- [198] P. Rajivan and C. Gonzalez. "Creative Persuasion: A Study on Adversarial Behaviors and Strategies in Phishing Attacks". In: *Frontiers in Psychology* 9 (2018), p. 135.
- [199] P. Rajivan, M. A. Janssen, and N. J. Cooke. "Agent-based Model of a Cyber Security Defense Analyst Team". In: *Proceedings of the Human Factors and Ergonomics Society Annual Meeting* 57.1 (2013), pp. 314–318.
- [200] Prashanth Rajivan and Cleotilde Gonzalez. "Creative persuasion: a study on adversarial behaviors and strategies in phishing attacks". In: *Frontiers in psychology* 9 (2018), p. 135.
- [201] S. Rajtmajer et al. "Constrained Social-Energy Minimization for Multi-Party Sharing in Online Social Networks". In: *Autonomous Agents and Multi-Agent Systems* 31.4 (2017), pp. 867–896.
- [202] David G Rand et al. "Cluster analysis reveals a binary effect of storage on Drosophila courtship conditioning". In: *Nature Communications* 11.1 (2020), pp. 1–12.
- [203] Anatol Rapoport. "A note on the index of cooperation for Prisoner's Dilemma". In: *Journal of Conflict Resolution* 11.1 (1967), pp. 100–103.
- [204] S. Roy et al. "A Survey of Game Theory as Applied to Network Security". In: 43rd Hawaii International Conference on System Sciences (2010), pp. 1–10.
- [205] Vildan Salikutluk et al. "An Evaluation of Situational Autonomy for Human-AI Collaboration in a Shared Workspace Setting". In: *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. ACM. 2024.
- [206] D. D. Salvucci and N. A. Taatgen. "Threaded Cognition: An Integrated Theory of Concurrent Multitasking". In: *Psychological Review* 115.1 (2008), pp. 101–130.

- [207] D. Sanchez, S. Martinez, and C. Gonzalez. "Cognitive Models for Adaptive Cyber Defense". In: *Journal of Defense Modeling and Simulation* 19.4 (2022), pp. 289–304.
- [208] Francisco C Santos et al. "Social networks and evolutionary games". In: *Nature* 456.7219 (2008), pp. 943–946.
- [209] Iqbal H Sarker. "AI-based modeling: Techniques, applications and research issues towards automation, intelligent and smart systems". In: *SN Computer Science* 4.1 (2023), p. 74.
- [210] M. A. Sasse, S. Brostoff, and D. Weirich. "Transforming the 'Weakest Link': A Human/Computer Interaction Approach to Usable and Effective Security". In: *BT Technology Journal* 19.3 (2015), pp. 122–131.
- [211] Beau G. Schelble et al. "Designing Human-Autonomy Teaming Experiments Through Reinforcement Learning". In: *Proceedings of the Human Factors and Ergonomics Society Annual Meeting* (2020). DOI: 10.1177/1071181320641340.
- [212] Bruce Schneier. Secrets and Lies: Digital Security in a Networked World. John Wiley & Sons, 2000.
- [213] Vikram Shashank, Aisha Patel, and Marco Rodriguez. "Polymorphic Malware Evolution: A Study of Creative Evasion Techniques". In: *Journal of Cybersecurity* 15.3 (2023), pp. 412–429. DOI: 10.1109/JCYBER.2023.0053.
- [214] Jeffrey W. Sherman et al. "Attentional Processes in Stereotype Formation: A Common Model for Category Accentuation and Illusory Correlation." In: *Journal of Personality and Social Psychology* (2009). DOI: 10.1037/a0013778.
- [215] O. Shoetan, K.-K. R. Choo, and A. Hassanzadeh. "Analyzing Cyber Attacker's Risk Tolerance and Decision Making". In: *Journal of Cyber Security and Mobility* 11.1 (2022), pp. 127–154.
- [216] Philip Olaseni Shoetan et al. "Synthesizing Ai's Impact on Cybersecurity in Telecommunications: A Conceptual Framework". In: *Computer Science & It Research Journal* (2024). DOI: 10.51594/csitrj.v5i3.908.
- [217] H. A. Simon. "A Behavioral Model of Rational Choice". In: *The Quarterly Journal of Economics* 69.1 (1955), pp. 99–118.
- [218] H. A. Simon. "Models of Man: Social and Rational". In: (1957).
- [219] Herbert A Simon. *Bounded rationality*. MIT Press, 1990.
- [220] Brian Singer et al. "On the Feasibility of Using LLMs to Execute Multistage Network Attacks". In: *arXiv preprint arXiv:2501.16466* (2025).
- [221] Laurens Sion et al. "An Architectural View for Data Protection by Design". In: (2019). DOI: 10.1109/icsa.2019.00010.
- [222] Laurens Sion et al. "SPARTA: Security & Amp; Privacy Architecture Through Risk-Driven Threat Assessment". In: (2018). DOI: 10.1109/icsa-c.2018.00032.
- [223] Ho Chit Siu et al. "Evaluation of Human-AI Teams for Learned and Rule-Based Agents in Hanabi". In: *Advances in Neural Information Processing Systems (NeurIPS)*. Vol. 34. 2021, pp. 16183–16195.
- [224] Florian Skopik, Timea Pahi, and Markus Leitner. *Cyber Situational Awareness in Public-Private-Partnerships: Organisationsübergreifende Cyber-Sicherheitsvorfälle effektiv bewältigen*. 1st. Berlin, Germany: Springer Vieweg, 2018, p. 347. ISBN: 978-3-662-56083-9.

- [225] Tommy van Steen and Julia R.A. Deeleman. "Successful Gamification of Cybersecurity Training". In: *Cyberpsychology Behavior and Social Networking* (2021). DOI: 10.1089/cyber.2020.0526.
- [226] J. Steinke et al. "Improving Cybersecurity Incident Response Team Effectiveness Using Teams-Based Research". In: *IEEE Security & Privacy* 13.4 (2015), pp. 20–29.
- [227] Jeffrey R Stevens et al. "Cognitive constraints in strategic decision making". In: *Nature Human Behaviour* 2.3 (2018), pp. 213–221.
- [228] Alexander J Stewart and Joshua B Plotkin. "From extortion to generosity, evolution in the Iterated Prisoner's Dilemma". In: *PNAS* 110.38 (2013).
- [229] Kevin Stine et al. "Integrating Cybersecurity and Enterprise Risk Management (ERM)". In: (2020). DOI: 10.6028/nist.ir.8286-draft2.
- [230] Peter Stone et al. "Ad hoc autonomous agent teams: Collaboration without pre-coordination". In: *Twenty-Fourth AAAI Conference on Artificial Intelligence*. 2010.
- [231] Barry Strauch. "Ironies of automation: Still unresolved after all these years". In: *IEEE Transactions on Human-Machine Systems* 48.5 (2017), pp. 419–433.
- [232] Valdemar Švábenský et al. "Enhancing Cybersecurity Skills by Creating Serious Games". In: (2018). DOI: 10.1145/3197091.3197123.
- [233] Attila Szolnoki and Matjaž Perc. "Modelling the evolution of human fairness in structured populations". In: *Scientific Reports* 9.1 (2019), pp. 1–9.
- [234] Hamed Taherdoost. "Towards an Innovative Model for Cybersecurity Awareness Training". In: *Information* (2024). DOI: 10.3390/info15090512.
- [235] Paul Théron and Alexander Kott. "When Autonomous Intelligent Goodware Will Fight Autonomous Intelligent Malware: A Possible Future of Cyber Defense". In: (2019). DOI: 10.1109/milcom47813.2019.9021038.
- [236] D. R. Thomas and R. Sule. "Security evaluation of the risk tolerance of individuals to cybercrime". In: *IEEE Security & Privacy* 19.5 (2021), pp. 42–50.
- [237] Godwin Thomas and Mary-Jane Sule. "A Service Lens on Cybersecurity Continuity and Management For organizations' Subsistence and Growth". In: *Organizational Cybersecurity Journal Practice Process and People* (2022). DOI: 10.1108/ocj-09-2021-0025.
- [238] J. N. Tioh, M. Mina, and D. W. Jacobson. "Cyber Security Training Using Interactive Game Mechanisms". In: *Journal of Cyber Security Technology* 1.2 (2017), pp. 86–105.
- [239] P. M. Todd and G. Gigerenzer. "Ecological Rationality: Intelligence in the World". In: *Oxford University Press* (2012).
- [240] Peter M Todd et al. "Ecological rationality: Intelligence in the world". In: *Oxford University Press* (2012).
- [241] Deepak K. Tosh et al. "An evolutionary game-theoretic framework for cyber-threat information sharing". In: *IEEE International Conference on Military Communications* (2015), pp. 585–590.
- [242] Ulubilge Ulusoy and Garrett E. Reisman. "Human Factors Respect in Human Autonomy Teams". In: (2024). DOI: 10.36227/techrxiv.171502760.01527479/v1.
- [243] Jay J. Van Bavel et al. "Contextual Sensitivity in Scientific Reproducibility". In: *Proceedings of the National Academy of Sciences* (2016). DOI: 10.1073/pnas.1521897113.

- [244] V. D. Veksler et al. "Cognitive Models in Cybersecurity: Learning from Expert Analysts and Predicting Attacker Behavior". In: *Frontiers in Psychology* 11 (2020), p. 1049.
- [245] A. Vishwanath. "Mobile Device Affordance: Explicating How Smartphones Influence the Outcome of Phishing Attacks". In: *Computers in Human Behavior* 63 (2016), pp. 198–207.
- [246] A. Vishwanath et al. "Why Do People Get Phished? Testing Individual Differences in Phishing Vulnerability within an Integrated, Information Processing Model". In: *Decision Support Systems* 51.3 (2011), pp. 576–586.
- [247] Joel H. K. Vuolevi and Paul A. M. Van Lange. "On the boundaries of social exchange: Vague systems and social interdependence". In: *Journal of Experimental Social Psychology* 48.1 (2012), pp. 100–111.
- [248] J. Wang, R. Chellappa, and P. J. Phillips. "Safety in Cyberspace: A Cognitive Science Perspective". In: *IEEE Systems, Man, and Cybernetics Magazine* 1.3 (2015), pp. 27–34.
- [249] Zhihui Wang et al. "Evolving strategies for the Iterated Prisoner's Dilemma". In: *Scientific Reports* 8.1 (2018), pp. 1–8.
- [250] Robert L West and Christian Lebiere. "Simple games as dynamic, coupled systems: Randomness and other emergent properties". In: *Cognitive Systems Research* 1.4 (2001), pp. 221–239.
- [251] R.W. Wohleber, K. Stowers, and Y. Lin. "Understanding and designing human-autonomy team trust: A hierarchical cognitive task analysis approach in urban search and rescue". In: *Applied Ergonomics* 109 (2023), p. 107866.
- [252] Junhui Wu et al. "Too much or too little? A meta-analysis of the social comparison effects on cooperation". In: *Psychological Bulletin* 146.7 (2020), pp. 651–679.
- [253] Z. Yan et al. "Information Security Knowledge Sharing Behavior Analysis in Knowledge Network". In: *Journal of Industrial Information Integration* 16 (2019), p. 100100.
- [254] Wako Yoshida, Raymond J. Dolan, and Karl J. Friston. "Game Theory of Mind". In: *PLOS Computational Biology* 4.12 (2008), e1000254. DOI: 10.1371/journal.pcbi. 1000254.
- [255] Jessica Young et al. "Attention and strategic decision making". In: *Cognitive Science* 43.4 (2019), e12721.
- [256] Dandan Zhang et al. "The dynamics of belief updating in human cooperation: findings from inter-brain ERP hyperscanning". In: *NeuroImage* 198 (2019), pp. 1–12.
- [257] J. Zhang, J. Zhuang, and V. R. R. Jose. "The Role of Risk Preferences in a Multi-target Defender-attacker Resource Allocation Game". In: *Reliability Engineering & System Safety* 169 (2018), pp. 95–104.
- [258] Jing Zhang, Jun Zhuang, and Victor Richmond R Jose. "The role of risk preferences in a multi-target defender-attacker resource allocation game". In: *Reliability Engineering & System Safety* 169 (2018), pp. 95–104.
- [259] Qian Zhang et al. "Investigating AI Teammate Communication Strategies to Support Human-AI Collaboration in Decision-Making". In: *Proceedings of the ACM on Human-Computer Interaction (CSCW)*. Vol. 7. CSCW2. 2023, pp. 1–30.

- [260] Qian Zhang et al. "You Complete Me: Human-AI Teams and Complementary Expertise". In: *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. ACM. 2022, pp. 1–19.
- [261] [Author name] Zhao. "Model-based reinforcement learning: A survey". In: [Journal name] (2023).
- [262] Qinian Zhong and Qian Kang. "Ransomware Detection With Opcode Analysis and GAN-Based Unsupervised Learning". In: (2024). DOI: 10.21203/rs.3.rs-3819158/v1.
- [263] V. Zimmermann and K. Renaud. "Moving from a 'Human-as-Problem" to a 'Human-as-Solution" Cybersecurity Mindset". In: *International Journal of Human-Computer Studies* 131 (2019), pp. 169–187.



## A Cyber-War Between Bots: Cognitive Attackers are More Challenging for Defenders than Strategic Attackers

YINUO DU, Software and Social Systems Department, Carnegie Mellon University, Pittsburgh, United States

BAPTISTE PREBOT, Department of Social and Decision Sciences, Carnegie Mellon University, Pittsburgh, United States

TYLER MALLOY, Department of Social and Decision Sciences, Carnegie Mellon University, Pittsburgh, United States

CLEOTILDE GONZALEZ, Department of Social and Decision Sciences, Carnegie Mellon University, Pittsburgh, United States

Adversary emulation is commonly used to test cyber-defense performance against known threats to organizations. However, many adversary emulation methods often rely on automated planning and underplay the role of human cognition. Consequently, defenders are often underprepared for human attackers who can think creatively and adapt their strategies. In this article, we propose the design of adversarial cognitive agents that are dynamic, adaptable, and able to learn from experience. These cognitive agents are built based on the theoretical principles of Instance-Based Learning Theory (IBLT) of experiential choice in dynamic tasks, making them more challenging than strategically optimal adversaries for human defenders. Our research offers three main contributions. First, in a simulation experiment, we demonstrate how IBL attacker agents can learn from experience and become as efficient as optimal strategic algorithms against a strategic defender. In a second simulation experiment, the IBL attackers are pitted against an IBL defender, showing that the IBL attacker can be a more challenging adversary for the IBL defender, while the IBL defender can learn to counter carefully crafted optimal attack strategies. To test these observations, we conducted a third experiment, where humans played the role of defenders against both strategic and IBL attackers in an interactive task. The results confirm the predictions of the second simulation experiment: Cognitive attackers are more challenging for human defenders than strategic attackers. These insights contribute to informing future adversary emulation efforts and training of cyber defenders.

CCS Concepts: • Computing methodologies  $\rightarrow$  Modeling methodologies; Model verification and validation; Cognitive science; • Security and privacy  $\rightarrow$  Human and societal aspects of security and privacy; • Human-centered computing  $\rightarrow$  Empirical studies in interaction design;

Additional Key Words and Phrases: Adversary emulation, cognitive modeling, cyber security game

This research was sponsored by the Army Research Office and accomplished under Australia-USA MURI Grant Number W911NF-20-S-000 and by the Army Research Laboratory under Cooperative Agreement Number W911NF-13-2-0045 (ARL Cyber Security CRA).

Authors' Contact Information: Yinuo Du, Software and Social Systems Department, Carnegie Mellon University, Pittsburgh, Pennsylvania, United States; e-mail: yinuod@andrew.cmu.edu; Baptiste Prebot, Department of Social and Decision Sciences, Carnegie Mellon University, Pittsburgh, Pennsylvania, United States; e-mail: baptiste.prebot@ensc.fr; Tyler Malloy, Department of Social and Decision Sciences, Carnegie Mellon University, Pittsburgh, Pennsylvania, United States; e-mail: tylermal@andrew.cmu.edu; Cleotilde Gonzalez, Department of Social and Decision Sciences, Carnegie Mellon University, Pittsburgh, Pennsylvania, United States; e-mail: coty@cmu.edu.



This work is licensed under a Creative Commons Attribution 4.0 International License.

© 2025 Copyright held by the owner/author(s). ACM 2469-7818/2025/03-ART7 https://doi.org/10.1145/3712672

7:2 Y. Du et al.

#### **ACM Reference Format:**

Yinuo Du, Baptiste Prebot, Tyler Malloy, and Cleotilde Gonzalez. 2025. A Cyber-War Between Bots: Cognitive Attackers are More Challenging for Defenders than Strategic Attackers. *ACM Trans. Soc. Comput.* 8, 3-4, Article 7 (March 2025), 22 pages. https://doi.org/10.1145/3712672

#### 1 Introduction

Cyber systems have gradually populated all the personal and collective layers of society. From banks to hospitals, from electric grids to industrial facilities, the interconnectivity of systems has created new opportunities for criminals. Cyber security is a domain of great complexity, defined by uncertainty, lack of visibility, extreme speeds, and partial information. In this adversarial context, defenders and attackers confront each other using digital weapons that are beyond the limits of human capabilities for perception and assessment. Defenders need extensive experience to effectively defend against dynamic and distributed attacks.

Cyber wargaming and adversary emulation (i.e., Red teams) are common practices in organizations to train defenders (i.e., Blue teams) and to develop appropriate defense algorithms [10, 15]. However, the design of emulated adversaries can be expensive and time-consuming, especially for scaled networks with a large attack surface and rich defense arsenals. Autonomous agents have been developed to mitigate this problem [5, 47, 51]. For example, Kotenko [30] modeled a DDoS attack, and Razak et al. [43] simulated network intrusions. However, these simulations do not specifically portray the attacker. The early models contained static patterns prescribed for the attacker agents to follow [22]. These models eventually gave way to graph-based [29] and state-based [1] attack simulation methods, which provide a useful characterization of the attacker's profile, such as goals, starting points, and available time. This group of simulation methods models and stores generic attack patterns with preconditions and postconditions in a knowledge base. Additional attack pattern attributes include the cost of attempts, execution time, base success probability, and maximum attempts.

Despite technical fidelity, most automated adversary simulation methods ignored the social context and lacked a dynamic behavior component [28]. Human attackers have varying levels of risk tolerance, which might affect their choice of target and attack methods [55]. Human attackers can also learn from their experiences [31], dynamically adapt to defenses they encounter, and modify their strategies accordingly, making them more dangerous over time as they become more adept at evading detection and exploiting vulnerabilities. Thus, to improve the training of defenders, the emulated adversaries need to exhibit behavior similar to that of the human attackers and have the capability to learn and adapt to the defender's actions.

Cognitive architectures and theories of human decision-making have made significant progress in emulating human-like behavior in dynamic environments. Unlike typical computational algorithms that aim to make optimal decisions, cognitive architectures adhere to human constraints such as forgetting, limited attention, and bounded rationality [16]. Cognitive models based on **instance-based learning theory (IBLT)** [20] have been implemented in the context of cybersecurity to model human cognitive processes. Dutt et al. [13] proposed an IBL model that represents the awareness of cyber situation of a human analyst and is capable of making concrete predictions about the recognition and comprehension processes of a security expert in a cyber attack. A more recent model from Du et al. [12] uses a cyber security scenario in which the IBL defender learns to defeat the most aggressive, optimal, but deterministic attack strategy. Cognitive models have also been used as embedded computational agents to simulate human interactions with software and networks [18, 35, 52]. Previous work has focused independently on

understanding defense behaviors and developing cognitive models of blue agents [12, 13] or the attack preferences and biases of the attacker [11]. However, past work on cognitive modeling in cyber security systems has rarely considered the real-time social interactions of attackers and defenders together.

The attacker and the defender can influence each other in cyber adversarial scenarios [53]. Such dynamics between attackers and defenders can make defenders more vulnerable to adversarial actions compared to even random attackers [36]. For example, in a simple and abstract game, humans were found to handle random attacks more effectively than adaptive attacks [36]. This suggests that commonly used random attack security algorithms may be less effective than human-inspired adaptive attack strategies in training human defenders. The characteristics of human attackers have been studied in a phishing experiment. For example Rajivan and Gonzalez [42] found that individual creativity is a predictor of an adversary's ability to evade detection. Cognitive biases and emotions are also believed to affect attacker behavior and decision-making [14, 26].

Given the current evidence on cyber wargaming and adversary emulation, we hypothesize that cognitive agents that are capable of emulating *human* adversaries will be more challenging for cyber defenders than deterministic attacker strategies. If this hypothesis is correct, then training cyber defenders against cognitive adversaries will result in better-prepared defenders than the current procedure for training cyber defenders against strategic attackers. This article aims to test this hypothesis.

The contributions of this article are as follows: First, we propose a cognitive model of a red agent that uses the theoretical principles of IBLT (that is,  $IBL_{Red}$ ) [20]. In a simulation experiment, we compare the performance of the  $IBL_{Red}$  agent with that of a deterministic, highly accurate, and targeted attack strategy ( $Beeline_{Red}$ ) and with a strategy that explores the network without prior knowledge about the location of targets ( $Meander_{Red}$ ). In Experiment 1, we demonstrate that  $IBL_{Red}$  is capable of learning and achieving performance similar to the optimal agent  $Beeline_{Red}$ . In Experiment 2, we tested the three red agents ( $IBL_{Red}$ ,  $Beeline_{Red}$ , and  $Meander_{Red}$ ) against a cognitive defender  $IBL_{Blue}$  to demonstrate that  $IBL_{Red}$  is the most challenging attacker for  $IBL_{Blue}$ . In Experiment 3, we validate the simulation findings of Experiment 2 in an experiment in which human participants play the role of defender. The results of the experiment confirm that cognitive adversaries are more challenging to human defenders than strategic adversaries.

#### 2 Instance-based Learning Theory

IBLT is a cognitive theory of decision-making. It is based on the idea that decisions are made by recognizing similar past experiences, integrating them into generating the expected utility of decision alternatives, and selecting the alternative with the maximum expected utility [20]. The development of cognitive models for cyber defense is based on a large body of work on applying cognitive science to cyber security (e.g., Reference [17]).

Although both the process and the mechanisms of IBLT have been published, we repeat the mathematical formulations of the theory here for completeness. The central element of IBLT is the "instance." It represents a unit of memory resulting from evaluating potential choice alternatives. Each decision is stored in an instance, structured with three elements that are built over time: A situation state s that is composed of a set of characteristics f; a decision or action a taken corresponding to an alternative in state s; and an expected utility or experienced result x of the action taken in a state. Concretely, for an IBL agent, an option k = (s, a) is defined by action a in state s. At time t, assume that  $n_{kt}$  different instances  $(k_i, x_{ik_it})$  for  $i = 1, \ldots, n_{kt}$ , associated with k. Each instance i in memory has an activation value, which represents the ease of retrieving this information from memory [4]. Here, we consider a simplified version of the activation equation,

7:4 Y. Du et al.

which only captures recency, frequency, and noise in memory:

$$\Lambda_{ik_{i}t} = \ln \left( \sum_{t' \in T_{ik_{i}t}} (t - t')^{-d} \right) + \sigma \ln \frac{1 - \xi_{ik_{i}t}}{\xi_{ik_{i}t}}, \tag{1}$$

where d and  $\sigma$  are the decay and noise parameters, respectively, and  $T_{ik_it} \subset \{0, \dots, t-1\}$  is the set of previous timestamps in which instance i was observed. The rightmost term represents noise to capture individual variation in activation, and  $\xi_{ik_it}$  is a random number drawn from a uniform distribution U(0,1) at each step and for each instance and option.

Activation of an instance i is used to determine the probability of retrieving an instance from memory. The probability of an instance where a soft-max function defines i:

$$P_{ik_{i}t} = \frac{e^{\Lambda_{ik_{i}t}/\tau}}{\sum_{j=1}^{n_{kt}} e^{\Lambda_{jk_{j}t}/\tau}},$$
(2)

where  $\tau$  is the Boltzmann constant (i.e., the "temperature") in the Boltzmann distribution. For simplicity,  $\tau$  is often defined as a function of the same  $\sigma$  used in the activation equation  $\tau = \sigma \sqrt{2}$ .

The expected utility of option k is calculated based on *Blending* as specified in discrete choice tasks [19]:

$$V_{kt} = \sum_{i=1}^{n_{kt}} P_{ik_i t} x_{ik_i t}. {3}$$

The choice rule is to select the option corresponding to the maximum blended value. When the agent receives delayed results, the agent updates the expected utilities using a credit assignment mechanism [38].

#### 3 Cyber Security Scenario

Testing attacker and defender agents requires a simulation or training platform that encapsulates cyber elements in an integrated environment. On such a platform, defense agents can confront attack agents in cyber scenarios and network simulations. Here, we use the interactive defense game, based on CybORG AI gym [7, 8, 50] with adversarial cyber-operation scenarios to allow users to train agents in a simple but realistic environment. We adopt the CAGE cyber defense scenario [49] to perform experimental simulations using IBL agents as cyber defenders and cyber attackers. This framework was also presented in References [12, 41], and in the following, we outline its main structural elements and the particularities of the cyber defense scenario.

The attacker (hereafter the Red agent) interacts with the environment through high-level actions that aim to progress and impact the network; the defender (hereafter the Blue agent) aims to stop the progression of the attacker and remove it from the network. Each combat between an attacker and a defender is an episode, and the duration of episode was set to 25 steps to ensure that the Blue agent could fully observe the attack strategies.

Figure 1 illustrates the topology of the network chosen for this scenario. The network is divided into three subnets: Subnet 1 consists of user hosts that are not critical, Subnet 2 consists of enterprise servers designed to support the user activities on Subnet 1, and Subnet 3 contains the critical operational server and operational hosts. Attackers usually start with user subnets and establish their entry point through social engineering or spear phishing. In this scenario, the entry point is host *User0*.

Figure 2 summarizes the phases of a targeted attack led by the Red agent (red arrows) and countermeasures for the Blue agent to stop it (blue arrows). The Red agent starts by searching

<sup>&</sup>lt;sup>1</sup>http://janus.hss.cmu.edu:8084/

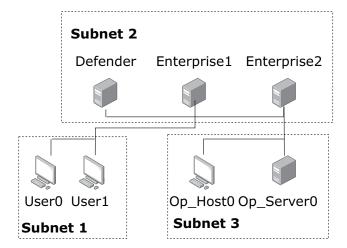


Fig. 1. Adaptation of the cage challenge network.

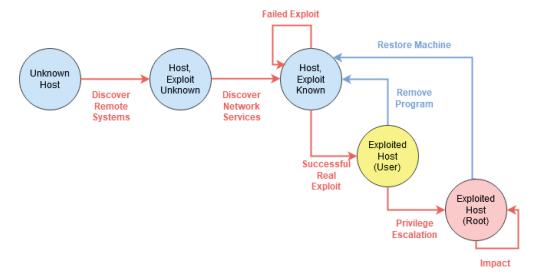


Fig. 2. Effect of actions on the host state (diagram from Reference [49]).

for hosts on the network with *DiscoverRemoteSystems*. To identify vulnerabilities on a target host, the next step is to *DiscoverNetworkServices*. A successful *ExploitRemoteService* on target can obtain *User* level access for the Red agent, which can be escalated to a more privileged *Root* level by *PrivilegeEscalate*. The Blue agent can *Remove* its adversary at the *User* level and use *Restore* if the Red agent has escalated. It can also *Analyse* the activities for additional information or passively *Monitor* the network.

#### 3.1 Red Agents

We used three types of red agents: Two strategic attackers: (1) a highly efficient deterministic agent,  $Beeline_{Red}$ , and (2) a stochastic agent,  $Meander_{Red}$ ; and (3) a dynamic cognitive agent,  $IBL_{Red}$ .

Advanced attackers well-funded by organizations or governments are well planned and highly organized to increase the probability of the success of the attack [3]. This type of attackers extensively research their target, collecting the necessary information and intelligence on the assets of the organizations. Before launching attacks, the attackers can acquire details of the

7:6 Y. Du et al.

network layout, such as the types of switches, routers, anti-virus tools, firewalls, Web servers used, and ports open. The attackers can then build attack plans using well-known vulnerability databases such as the **Common Vulnerabilities and Exposures List (CVE)** and the NIST **National Vulnerability Database (NVD)** [34]. These plans allow attackers not only to establish a foothold, but also to penetrate deeper into the target's network.  $Beeline_{Red}$  assumes that the attacker has prior knowledge of the network topology and moves directly to the operational server following the red path ( $User0 \rightarrow User1 \rightarrow Enterprise1 \rightarrow Enterprise2 \rightarrow Op\_Server0$ ) (see Figure 1) in a predictive and deterministic way.

Novice attackers (also known as "Script Kiddies") who rely on pre-made exploit programs and files ("scripts") are usually not dedicated enough to their hacking [21]. Instead of making careful plans and collecting the necessary tools beforehand, they tend to scan random IP blocks on the Internet for weaknesses and exploit them as they are found. With the increasing amount of tools and scripts available on the Internet for free, novice attackers can also do a lot of damage to well-protected systems [9].  $Meander_{Red}$  assumes no prior knowledge about the network structure and behaves in a stochastic manner by choosing a random target to move forward.

In contrast, the cognitive agent,  $IBL_{Red}$ , is a novel contribution to this research, and it is a dynamic agent that learns from experience, as described in the IBLT section above.  $IBL_{Red}$  intends to represent cognitive memory-based decisions that can adapt their actions dynamically according to the conditions of the environment and the actions of the blue agent. The instances represent each decision made and are structured with the following three elements:

State,  $s_a$ : The state of the instances of the  $IBL_{Red}$  agent is composed of features, f, constructed using the concept of Attack Models and Attack Graphs introduced by Sheyner et al. [48] to model the security vulnerabilities of a network and their exploitation from the perspective of an attacker. Specifically, contextual characteristics include the success status of the previous action of the  $IBL_{Red}$  agent and the resources it occupied. A slot is dedicated to each type of resource in various states, as shown in Figure 2.

Specifically, a subnet can be newly *Detected* or already *Scanned*, while hosts are classified as *Detected*, *Scanned*, *Exploited* (*User*), *Exploited* (*Root*), *Impacted*.

The starting status denotes when the  $IBL_{Red}$  agent has just successfully established its foothold on the network on User0. At that point, only the User subnet is detected in addition to its entry point User0, while the rest of the slots are empty. The most successful final state for the  $IBL_{Red}$  agent is where all hosts and servers are exploited at the Root level and when critical  $Op\_Server0$  is impacted.

Action Space,  $a_a$ : The action space for the  $IBL_{Red}$  agent is dynamically constructed at each step based on the status of each host on the network. Each action consists of a target host and an applicable command. As shown in Figure 2,  $IBL_{Red}$  can choose to collect more information about hosts in the network or advance the attack status of known hosts.

Utility,  $z_a$ : A reward is calculated at each step, based on the attack status, as shown in Table 1. Higher rewards are assigned when the  $IBL_{Red}$  agent is able to access more significant systems. Only root access to the systems and successful impact on the operational server are rewarded. The agent  $IBL_{Red}$  receives a reward of 0 for any other action.

3.1.1 Metrics for the Red Agent. The performance of the Red agent was evaluated for each episode using the following metrics: (1) Reward: the cumulative rewards received during the execution of the scenario; (2) Impact duration: the average number of steps per episode that the Red agent successfully impacts the operational server; (3) Progress: the average number of steps

<b>Event or Action</b>	Reward
Administrator access on a Host	0.1
Administrator access on a Server	1
Successfully <i>Impact Op_Server0</i>	10

Table 1. Utility in the IBL Model: Events and Actions Costs

per episode that the Red agent took to penetrate *Enterprise subnet* and *Operational subnet*; and **(4) Action frequency:** the average proportion of command usage at each step in an episode.

#### 3.2 Blue Agents

We used two types of simulated blue agents: (1) a deterministic agent  $Sleepy_{Blue}$  and (2) a dynamic cognitive agent  $IBL_{Blue}$ .  $Sleepy_{Blue}$  does not attempt to stop the Red agent strategically and only takes the *Monitor* action to passively observe the state of the network. This agent helps simulate the situation where the defender fails to detect the existence of the stealthy attacker and thus does not employ defense measures.  $IBL_{Blue}$  is a dynamic agent proposed and tested in Du et al. [12]. It is also based on IBLT and has been demonstrated to resemble human-like defense decisions in empirical studies [40]. The instances represent each decision made and are structured with the following three elements:

State,  $s_d$ :  $IBL_{Blue}$  instance states are constructed to resemble the information that would be presented to a human defender in the scenario. Specifically, there are two slots for each host or server, representing the observed activity and the known compromised status of that host at a certain step in an episode. The order of (Activity, Compromised Status) pairs for each host is fixed to encode the identity of each host, i.e., the host name, IP address, and Subnet. The step index slot is included to resemble the step counter within each episode. The IBL agent has to choose the host to protect and the tool with which to protect it. Each action consists of a host and a command in the format of cmd host.

Action Space,  $a_d$ : At each step, based on the observed state of the network and the consequences of the attacker's previous actions,  $IBL_{Blue}$  selects a host or server to act on and one of four possible actions: Analyze is used to collect information on the level of compromise of the selected host; Remove is used to remove a suspected malicious agent from the host or server; if the malicious agent cannot be removed, then the blue agent can Restore a host or server to a previous stable state; and Monitor to just continue observing the system, which has essentially no effect on the state of the network.

*Utility*,  $z_d$ : The utility of the blue agent is the negative of the utility of the red agent, as the game is zero-sum.

3.2.1 Metrics for the Blue Agent. The performance of the Blue agent was also evaluated in terms of: (1) Action frequency: proportions of command usages at each step; and (2) Number of options: the average number of defense choices available to the blue agent. Each option is a tuple that contains the command and the target host.

#### 4 Simulation Experiment 1: Red Agents against a Deterministic Defender, Sleepy<sub>Blue</sub>

How well do the two strategic attackers ( $Beeline_{Red}$ ,  $Meander_{Red}$ ) compared to the cognitive red agent ( $IBL_{Red}$ ) perform against a passive defender ( $Sleepy_{Blue}$ )? The answer to this question provides a baseline for comparing the capabilities of the red agents.

7:8 Y. Du et al.

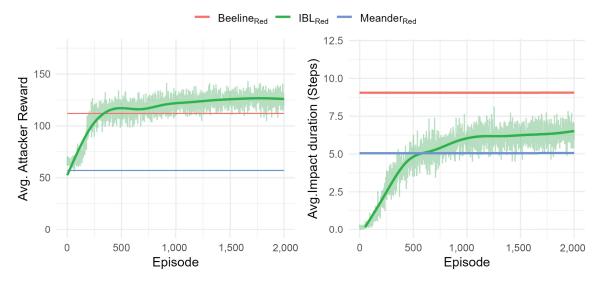


Fig. 3. Red agents performance when confronted by  $Sleepy_{Blue}$ : average reward per episode (left) and average impact duration (right).

Given that  $Beeline_{Red}$  is the best-known strategy in our scenario, we expect that this agent would perform better compared to the other two attackers ( $Meander_{Red}$  and  $IBL_{Red}$ ). The  $Beeline_{Red}$  agent will consistently receive the highest reward and maximum impact on the operational server, since  $Beeline_{Red}$  represents a highly effective but deterministic attacker. We also expect that the  $Meander_{Red}$  agent will consistently perform poorly against  $Sleepy_{Blue}$ , since this red agent behaves stochastically, without strategic knowledge of the path to move forward against the defender. In contrast, we expect that our newly proposed  $IBL_{Red}$  agent will initially perform poorly and similarly to  $Meander_{Red}$  when confronted with the  $Sleepy_{Blue}$  defender, but the  $IBL_{Red}$  agent would **learn** to take advantage of the ineffective  $Sleepy_{Blue}$  defender with practice and achieve a level of performance comparable to the optimal and strategic  $Beeline_{Red}$  agent.

#### 4.1 Methods

For each type of red agent (i.e.,  $Beeline_{Red}$ ,  $Meander_{Red}$ , and  $IBL_{Red}$ ), we performed 40 runs, each with 2,000 episodes. The  $IBL_{Red}$  agents were run with default decay d=0.5 and noise  $\sigma=0.25$  parameters. This means that the results presented here are all a priori predictions of the  $IBL_{Red}$  agent against the  $Sleepy_{Blue}$  defender.

#### 4.2 Results

4.2.1 Dynamic Cognitive Agents Learn Effective Attack Strategies against Passive Defenders. The left panel of Figure 3 shows the reward obtained by the Red agents against  $Sleepy_{Blue}$  defender. We observe that the reward of the  $Beeline_{Red}$  agents was consistently larger than the reward of the  $Meander_{Red}$  agents. Also, the reward of the  $IBL_{Red}$  agents was initially lower than the reward of the  $Beeline_{Red}$  agents; but the  $IBL_{Red}$  agents learned over the course of the episodes, approaching the rewards of the  $Beeline_{Red}$  agents after 2,000 episodes.

These observations are tested with one-way analysis of variance between subjects using *type of attack* as the main factor and *reward of attacker* as the dependent variable, aggregating for the first 500 and the last 500 episodes. As expected, the  $IBL_{Red}$  agents (M = 59.09, SD = 43.78) performed significantly worse than the  $Beeline_{Red}$  agents (M = 112.8, SD = 0) against  $Sleepy_{Blue}$  in the first 500 episodes [F(1,39998) = 30105, p < .001,  $\eta^2 = 0.43$ ]. In contrast, the  $IBL_{Red}$  agents (M = 104.64,

SD = 38.68) were able to achieve an average performance comparable to that of the  $Beeline_{Red}$  agents (M = 112.8, SD = 0) in the last 500 episodes. At the end of the 2,000th episode, 55% of the agents  $IBL_{Red}$  received a higher reward than  $Beeline_{Red}$ , which requires them to quickly penetrate the network to impact  $Op\_Server0$ , and at the same time fully exploit the remaining valuable systems when the opportunity arises. Most importantly, the agent  $IBL_{Red}$  learned such a complex and efficient strategy purely from experience according to IBLT [20] and without any explicit encoding of any strategy.

4.2.2 Impact Duration. The main goal of the attacker is to maintain constant Impact over  $Op\_Server0$ . The right panel of Figure 3 shows the number of successive impacts performed by the red agent on the  $Op\_Server$  after the first impact is achieved. As observed, the  $Beeline_{Red}$  agents consistently maintain a longer impact duration than the  $Meander_{Red}$  agents. The  $IBL_{Red}$  agents start poorly compared to the  $Beeline_{Red}$  and  $Meander_{Red}$  agents, but the  $IBL_{Red}$  agents are able to learn quickly and achieve a large number of impacts over the  $Op\_Server0$  with more task practice.

To test these observations, we performed a one-way ANOVA between subjects using *attacker* type as the main factor and *impact duration* as the dependent variable, aggregating for the first 500 and the last 500 episodes. As expected, the  $IBL_{Red}$  agents are capable of achieving a duration of impact similar to that of the network (M = 6.4, SD = 1) as the  $Beeline_{Red}$  agent (M = 9.0, SD = 3.37) when faced with  $Sleepy_{Blue}$  in the last 500 episodes [F(1,39998) = 901.4, p < .001,  $\eta^2 = 0.31$ ].

#### 5 Simulation Experiment 2: Red Agents against a Cognitive Defender, IBL<sub>Blue</sub>

After performing a baseline analysis of our three red agents against a passive defender, we compared the three red agents against an adaptive cognitive defender, the  $IBL_{Blue}$  agent. In Experiment 2, the strategies of the three agents,  $Beeline_{Red}$ ,  $Meander_{Red}$ , and  $IBL_{Red}^{Trained}$ , were kept the same as in Experiment 1. This means that the agent  $IBL_{Red}^{Trained}$  was trained against the agent  $Sleepy_{Blue}$  for 2,000 episodes and tested against the cognitive adaptive agent  $IBL_{Blue}$ .

Since the  $IBL_{Blue}$  agent can adapt to the strategies of the attackers, we expect that it will be able to learn effective defense strategies against the two static attackers: the  $Beeline_{Red}$  and  $Meander_{Red}$  agents. However, predictions of the performance of the  $IBL_{Blue}$  agent in defending against  $IBL_{Red}^{Trained}$  are less clear.

We expect that  $Beeline_{Red}$  will initially achieve a higher reward and a longer impact duration than the agents  $IBL_{Red}^{Trained}$  and  $Meander_{Red}$ . However, we expect that the determinism and static nature of  $Beeline_{Red}$  will be exploited by the learning agent  $IBL_{Blue}$ , resulting in a worse attack performance of  $Beeline_{Red}$  than  $IBL_{Red}^{Trained}$  with extended practice. We also expect that  $Meander_{Red}$  will start with a lower reward and a shorter impact duration than  $IBL_{Red}^{Trained}$ . Furthermore, although  $Meander_{Red}$  is stochastic, this attacker is not adaptive and is not a learning agent. For these reasons, we expect that, ultimately,  $Meander_{Red}$  will be exploited and effectively stopped by  $IBL_{Blue}$ .

#### 5.1 Methods

Similarly as in Experiment 1, we run 40  $IBL_{Blue}$  runs for each type of red agent (that is,  $Beeline_{Red}$ ,  $Meander_{Red}$ , and  $IBL_{Red}$ ). The IBL agents (red and blue) were configured with default decay d = 0.5 and noise  $\sigma = 0.25$  parameters.

#### 5.2 Results: Attacker Behavior

5.2.1 Dynamic Cognitive Agents Learn Effective Defense Strategies against Static Attackers, but Not Dynamic Ones.

7:10 Y. Du et al.

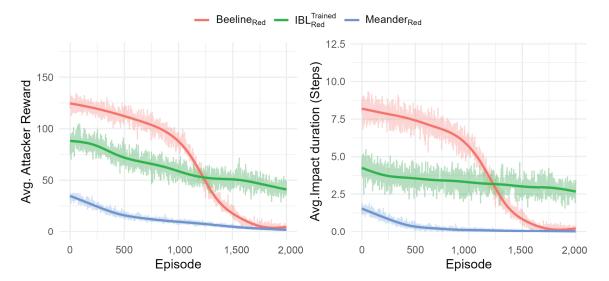


Fig. 4. Red agents performance when confronted by a cognitive defender,  $IBL_{Blue}$ . The average reward per episode (left) and average impact duration (right).

Attacker Reward. The left panel of Figure 4 shows the reward obtained by the Red agents against the  $IBL_{Blue}$  defender. We observe that the reward for the agents  $Beeline_{Red}$  was initially higher than the reward for the agents  $IBL_{Red}$ . However, over time, as the agent  $IBL_{Blue}$  learns to defend the network more effectively, the reward of the agents  $Beeline_{Red}$  drops significantly and approaches the reward of the agents  $Meander_{Red}$ . Meanwhile, the adaptive nature of the agent  $IBL_{Red}$  makes it difficult to learn a defense strategy within the training period investigated here.

Given that the agent  $IBL_{Blue}$  also starts naively learning from experience, the agents  $Beeline_{Red}$  performed better (M=112.8, SD=30.39) than the agents  $IBL_{Red}^{Trained}$  (M=80.34, SD=45.02) in the first 500 episodes [ $F(1,39998)=20579, p<.001, \eta^2=0.34$ ]. However,  $Beeline_{Red}$  agents (M=5.15, SD=14.68) performed significantly worse than  $IBL_{Red}^{Trained}$  agents (M=54.60, SD=34.31) in the last 500 episodes [ $F(1,39998)=31185, p<.001, \eta^2=0.44$ ]. The  $Beander_{Red}$  agents ( $Beander_{Red}$  agent) ( $Beander_{Red}$  agent) ( $Beander_{Red}$  agent) (Bean

Impact duration. As shown in the right panel of Figure 4, the duration of the impact shows trends similar to the reward obtained by the red agents against  $IBL_{Blue}$ .  $IBL_{Red}^{Trained}$  achieves shorter impact duration (M=4.08, SD=3.77) than  $Beeline_{Red}$  (M=8.93, SD=2.34) in the first 500 episodes [ $F(1,39998)=17240, p<.001, \eta^2=0.95$ ]. This relative disadvantage reversed in the last 500 episodes, where  $IBL_{Red}^{Trained}$  had a longer impact duration (M=3.25, SD=2.82) than  $Beeline_{Red}$ : (M=0.18, SD=1.11) [ $F(1,39998)=14699, p<.001, \eta^2=0.94$ ]. The agents  $Meander_{Red}$  start with a worse performance (M=0.8711282, SD=1.118891) than  $IBL_{Red}$  (M=3.25, SD=2.82) [ $F(1,39998)=19122, p<.001, \eta^2=0.32$ ] and are unable to cause an impact in the last 500 episodes.

Attacker Progress. To further explore the behavior of red agents, we analyzed the number of steps the attacker takes to reach a subnet (Enterprise and Operational). Taking into account the layered network structure shown in Figure 1, the progress of the Red agents can be marked by two milestones: penetration of the Enterprise subnet and the Operational subnet. The higher the number of steps required to reach a specific subnet, the more effective the defense agent is at protecting the network from attack.

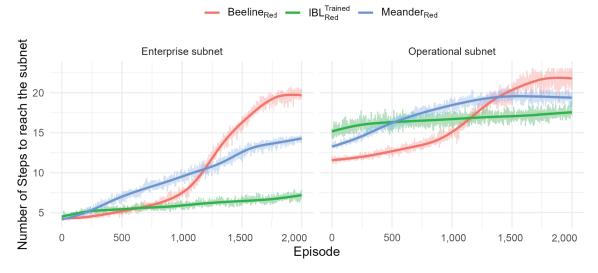


Fig. 5. Red agents progression capability when confronted by  $IBL_{Blue}$ . Number of steps taken to reach the Enterprise subnet (left) and Operational subnet (right).

We can observe the increasingly delayed and impeded forward progress of  $Beeline_{Red}$  in Figure 5. In general, each red agent takes longer to access both the Enterprise and Operational subnets at the end of the training period compared to earlier in the training, demonstrating that the  $IBL_{Blue}$  agent effectively learns to defend. At the end of the training period,  $Beeline_{Red}$  takes, on average, 15 more steps to enter the Enterprise subnet and 10 more steps to enter the Operational subnet compared to earlier in the training. However,  $IBL_{Red}^{Trained}$  shows relatively stable rates of the number of steps taken to reach these two subnets over the course of 2,000 episodes.

In the first 500 episodes,  $IBL_{Red}^{Trained}$  takes longer time to penetrate the Enterprise subnet (M=4.77, SD=0.33) and the Operational subnet (M=15.67, D=0.61) than  $Beeline_{Red}$  (Enterprise (M=4.59, SD=0.33), Operational (M=12.04, D=0.38)) [Enterprise:  $F(1,39998)=66.02, p<.001, \eta^2=0.06$ ] [Operational:  $F(1,39998)=12489, p<.001, \eta^2=0..93$ ]. But this relative disadvantage reversed in the last 500 episodes, where the  $IBL_{Red}^{Trained}$  propagated faster into the Enterprise subnet (M=6.12, SD=0.40) and the Operational subnet (M=16.90, D=0.52) than the  $Beeline_{Red}$ : Enterprise (M=19.20, SD=0.99) [ $F(1,39998)=74265, p<.001, \eta^2=0.99$ ], Operational  $(M=21.27, D=0.97), F(1,39998)=7940, p<.001, \eta^2=0.89$ ].

The faster progress speed of  $IBL_{Red}^{Trained}$  emerges from its learning capability against  $IBL_{Blue}$  in addition to the stochasticity of its actions. Although  $Meander_{Red}$  acts stochastically, it is significantly impeded by a lack of adaptation and takes, on average, 10 more steps to achieve each milestone than  $IBL_{Red}^{Trained}$ .

5.2.2 Dynamic Cognitive Agents Learn Defense Strategies That Disrupt the Behavior of Static Agents, but Not Dynamic Ones. Figure 6 compares the average distribution of the use of attack commands at each step of the first 500 episodes (left panels) versus the last 500 episodes (right panels).  $IBL_{Red}^{Trained}$  presents proportions of actions similar to  $Beeline_{Red}$  at the beginning, with higher Monitor proportion for  $IBL_{Red}^{Trained}$ . In the final episodes,  $Beeline_{Red}$  and  $Meander_{Red}$  are stuck in a loop of ExploitRemoteService and PrivilegeEscalate, while  $IBL_{Red}^{Trained}$  maintained a consistent distribution. This comparison constitutes further evidence of the inefficacy of strategic attack agents. The disappearance of Monitor, Analyse, and Impact actions can help explain the reason for the rapid drop in reward within the episodes.

7:12 Y. Du et al.

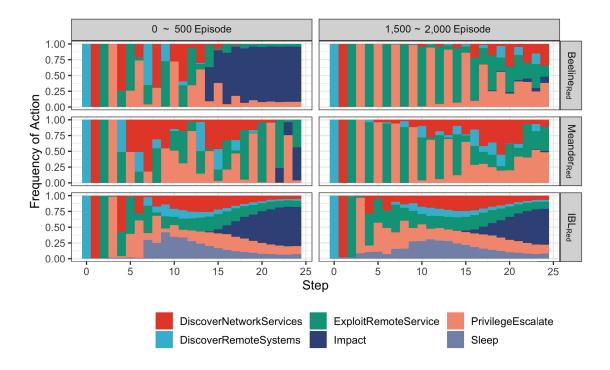


Fig. 6. Red agents evolution of action frequency in the first 500 episodes (left) and the last 500 episodes (right).

#### 5.3 Results: Defender Agent Behavior IBL<sub>Blue</sub>

Since the adversarial scenario is a zero sum game (i.e., the loss of the defender corresponds to the reward of the attacker), the performance results for  $IBL_{Blue}$  can be derived from the performance results of the red agents presented above. In particular,  $IBL_{Blue}$  performed better against  $Meander_{Red}$  and  $Beeline_{Red}$  than against  $IBL_{Red}^{Trained}$ . It was able to learn to defend effectively against  $Meander_{Red}$  and  $Beeline_{Red}$ , eventually achieving near-zero losses in the last 500 episodes. However,  $IBL_{Blue}$  was only able to reduce the loss of attacks  $IBL_{Red}^{Trained}$  by half.

To provide a deeper understanding of these results, we focus this section on the exploration of the defender's behavior against the three attackers.

5.3.1 Dynamic Cognitive Agents Adapt Defense Actions Taken against Static Agents, but Not Dynamic Ones. As presented in Figure 7, the dynamics of the use of defensive commands by the agent  $IBL_{Blue}$  shows a difference when confronting the agents  $Beeline_{Red}$  and  $Meander_{Red}$  in contrast to the agent  $IBL_{Red}^{Trained}$ .  $IBL_{Blue}$  agents faced with a strategic attacker are able to minimize the proportion of costly Restore action and stop the attacker with Remove in an earlier state of the cyber attack chain.

When the  $IBL_{Blue}$  against are defending against an  $IBL_{Red}^{Trained}$  attacker, they do not adjust the actions they take in a way that is able to respond appropriately to the attacker. This makes sense, given the constantly adapting nature of the  $IBL_{Red}^{Trained}$  attacker's behavior. These results indicate that the dynamics of this environment make it difficult for an adaptive defense strategy to learn effectively against an attacker that is also adapting their strategy.

5.3.2 Dynamic Cognitive Attackers Force Defenders to Choose from More Action Options Compared to Static Attackers. Psychology and behavioral research has shown that a large number

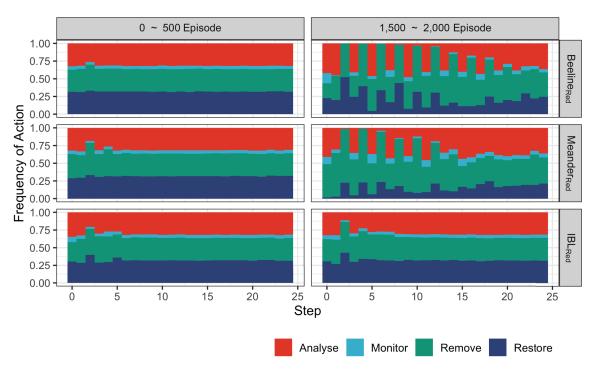


Fig. 7. IBL<sub>Blue</sub> evolution of action frequency in the first 500 episodes (left) and the last 500 episodes (right).

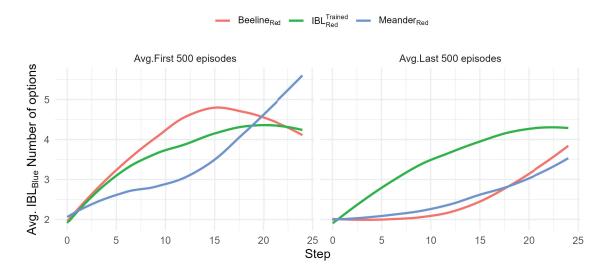


Fig. 8. Size of  $IBL_{Blue}$ 's option space in the first 500 episodes (left) and the last 500 episodes (right).

of options hinder learning speed [44, 46]. Defenders face this challenge of information overload. Figure 8 analyzes the number of options available to the  $IBL_{Blue}$  agent during the 25 steps of the episodes. As shown in the left panel, when  $IBL_{Blue}$  fought against  $Beeline_{Red}$  or  $Meander_{Red}$ , it was able to reduce the option space in the final 500 episodes compared to the first 500 episodes. In contrast, the option space of  $IBL_{Blue}$  stays about the same size from the first 500 episodes to the last 500 episodes when facing the adaptive agent  $IBL_{Red}^{Trained}$ . That is, agent  $IBL_{Blue}$  was unable to simplify its option space by impeding its progress and minimizing the number of attacked hosts.

7:14 Y. Du et al.

In summary,  $IBL_{Blue}$  is able to win the battle against  $Beeline_{Red}$  and  $Meander_{Red}$  after 2,000 episodes of repetitive interactions. It is able to block the opponent's progress, reduce the cognitive load of himself, and eliminate loss.  $IBL_{Red}^{Trained}$  appears to be much more persistent, and  $IBL_{Blue}$  can only slightly alleviate its impact.

#### 6 Experiment 3: Red Agents against Human Defenders

Although IBL models have shown successful replication of human behaviors in a variety of tasks, including cyberdefense [17], empirical verification of the findings of Experiment 2 with human defenders is needed. The first goal of this experiment is to compare the performance of human defenders faced with the three types of attackers,  $Beeline_{Red}$ ,  $Meander_{Red}$ , and  $IBL_{Red}^{Trained}$ . The second goal of this experiment is to validate the simulation results of predicted human behavior when paired against these three attackers, by comparing human performance with simulated  $IBL_{Blue}$  defenders. The results of this experimentation will provide support for the theoretical contribution of this work in motivating the use of cognitive agents to train cyber defenders.

#### 6.1 Experimental Design

Human participants completed the same cyber defense task and scenario as  $IBL_{Blue}$ , presented in Section 3. They performed the task using **Interactive Defense Game (IDG)** [40, 41], which provides an interactive decision game in the cyber task and environment.

#### 6.2 Participants

Participants were recruited through Amazon Mechanical Turk to participate in a cybersecurity study. The study was advertised to last between 30 and 45 minutes. The time it took between participants was  $M = 47.29 \pm 16.36$  minutes. Participants received a base compensation of \$4.50, and up to \$5.60 in bonus payment ( $M = 3.50 \pm 1.39$ ) based on their final score. <sup>2</sup> 186 participants (124 men, 61 women, 1 N / A) aged 21 to 65 years ( $M = 37.12 \pm 10.15$ ) completed the study. Seventeen of the 186 participants (9%) had more than five years of experience in the network operation and security area and at least a Master's degree in a related field. Each participant was randomly assigned to face one of the three red agents:  $Beeline_{Red}$ ,  $Meander_{Red}$ , or  $IBL_{Red}^{Trained}$ .

#### 6.3 Procedure

After giving their informed consent and completing a demographic questionnaire, the participants received instructions for the task followed by a short quiz to verify their basic understanding of the instructions for the task. The participants had to correctly answer all the questions before moving on to the next step of the experiment. The participants received feedback on the precision of their responses. There was no limit in the number of attempts the participants had to answer the questions correctly. However, we recorded the score of their first attempt and the number of times they tried to answer the questions. The participants then watched a video introduction to the IDG explaining the interface, the game controls, and the dynamics of an episode.

The participants then performed the task consisting of two phases: (1) a practice session and (2) a main task. The practice session consisted of two short episodes (i.e., games) of 10 steps each. The practice episodes were intended to familiarize participants with the interface and game controls. To do so, the participants successively faced Beeline and Meander; however, since these two deterministic attack strategies do not differ significantly during the first 10 steps, the participants did not have enough information to discriminate between them during the practice session.

<sup>&</sup>lt;sup>2</sup>As the score used in this experiment is negative (loss), the bonus payment was calculated by using the difference to the maximum possible loss and attributing \$0.005 per point: bonus=(total loss+1120)\*0.005.

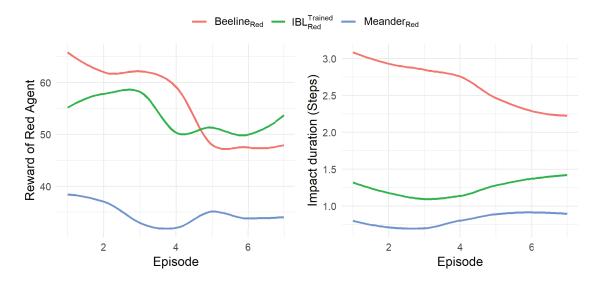


Fig. 9. Red agents performance when confronted by a human defender. The average reward per episode (left) and average impact duration (right).

Following the practice session, the participants performed the main task consisting of 7 episodes of 25 steps against the same type of adversary. No time restrictions were imposed. The initial state of the network was the same for all participants and for each of the episodes.

Subsequently, participants completed a post-experiment survey composed of two parts: (1) feedback on their performance and perceived strategy and (2) their experience in computer science and cyber defense. Finally, the participants received their final score and were dismissed. Experimental instructions, quizzes, and surveys, along with data and analysis scripts, can be accessed at <a href="https://osf.io/8vxej/?view\_only=c42691c2b5bb4c31a72b1ada00e38428">https://osf.io/8vxej/?view\_only=c42691c2b5bb4c31a72b1ada00e38428</a>.

#### 6.4 Results: Attacker Behavior against Human Defenders

6.4.1 Human Defenders Perform Similarly to Dynamic Cognitive Defenders against Both Static and Dynamic Attackers. The performance of three types of attackers (Beeline<sub>Red</sub>, Meander<sub>Red</sub>,  $IBL_{Red}^{Trained}$ ) against human defenders shown in the left panel of Figure 9 is in alignment with the simulation predictions shown in Figure 4. As human participants learn from experience, the  $Beeline_{Red}$  agents performed better (M = 66.25, SD = 5.39) than the  $IBL_{Red}^{Trained}$  agents (M = 55.51, SD = 4.70) in the first episode [F(1, 124) = 2.187, P = 0.142, P = 0.017].

However,  $IBL_{Red}^{Trained}$  agents posed a more persistent threat to the human defender than  $Beeline_{Red}$  agents. The performance of the  $Beeline_{Red}$  agents deteriorates rapidly. The  $Beeline_{Red}$  agents (M=47.33, SD=6.42) performed worse than the  $IBL_{Red}^{Trained}$  agents (M=54.141, SD=6.191) in the last episode [F(1,124)=0.577, p=0.449,  $\eta^2=0.005$ ]. The  $Meander_{Red}$  agents (M=5.26, SD=4.42) perform significantly worse [F(1,778)=46.762, p<0.001,  $\eta^2=0.057$ ] than  $IBL_{Red}^{Trained}$  in 7 episodes.

Consistent with the simulation prediction shown in Figure 4,  $IBL_{Red}^{Trained}$  demonstrates the persistent duration of the impact and is superior to  $Meander_{Red}$  (M=0.817, SD=1.488). In addition,  $IBL_{Red}^{Trained}$  achieves a shorter impact duration (M=1.29, SD=0.18) than  $Beeline_{Red}$  (M=3.15, SD=0.39) in the first episodes [ $F(1,872)=53.1, p<.001, \eta^2=0.95$ ]. This gap decreased significantly at the end of the seventh episode [ $Beeline_{Red}$ : (M=2.23, SD=0.41),  $IBL_{Red}^{Trained}$ : (M=1.44, SD=0.35)].

7:16 Y. Du et al.

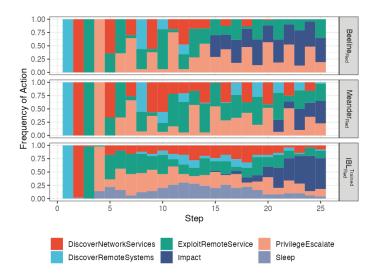


Fig. 10. Distribution of attack commands against human defenders.

These results show that the ordering of the highest and lowest attacker rewards and impact duration on early trials match between human and cognitive defenders. Additionally, the trends in adjusting the performance of attackers throughout the learning of the defender are similar. Overall, these results indicate a similar performance against all attack strategies between human and dynamic cognitive model defenders.

- 6.4.2 Distribution of Attack Commands. As demonstrated in Figure 10,  $IBL_{Red}^{Trained}$  used the Impact command the most often, while  $Beeline_{Red}$  and  $Meander_{Red}$  are unable to exert a consistent impact on the operational server and resort to ExploitRemoteService and PrivilegeEscalate.
- 6.4.3 Human Participants and Cognitive Defenders Behave Similarly against All Attacker Strategies. Consistent with the predictions of the simulated human performance, there was a significant difference in the mean reward by episode based on the type of Red agent (F(2,)=28.56, p<1e-10,  $\eta^2$ =0.5).
- 6.4.4 Action Frequency. Similar to  $IBL_{Blue}$ , the use of defensive commands by the human defender shows a difference when confronting the agent  $Beeline_{Red}$  and  $Meander_{Red}$  in contrast to the agent  $IBL_{Red}^{Trained}$ . Human participants are more passive and take more Analyse, Monitor actions than Remove, Restore actions. However, as shown in Figure 11, human participants have a consistent preference for the choice of action throughout the course of 7 episodes.
- 6.4.5 Size of Option Space. Figure 12 presents the number of options available to the human during the 25 steps of the episodes. Similarly to  $IBL_{Blue}$ , human participants can alleviate their cognitive load by narrowing down the option space when paired with  $Beeline_{Red}$ . The size of the option space remains approximately the same in the two stochastic conditions, that is,  $Meander_{Red}$  and  $IBL_{Red}^{Trained}$ .
- 6.4.6 Dynamic Cognitive Attackers Perform Best against the Most Efficient Human Defenders. To further investigate the human defenders against the various type of attacker models, we split the human defenders according to the distribution of attacker reward for each type of red agent. Efficient Defenders were those that resulted in attacker reward lower than the mean attack reward, and Inefficient Defenders were those that were equal or above the mean attack reward.

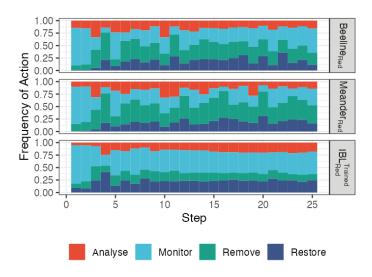


Fig. 11. Average action frequency of human defender.

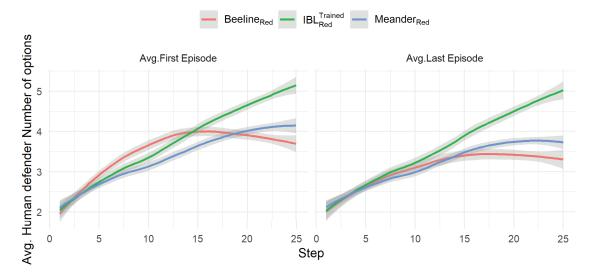


Fig. 12. Average size of the human defender's option space in the first (left) and the last episode (right).

Comparing performance in this way allowed us to determine if there was a clear difference in how red agents performed versus human participants who were better or worse able to learn attacker strategies. Furthermore, to compare performance once human defenders had enough experience to learn the attacker strategy adequately, we limited the statistical analysis to later (>4) episodes of the trial.

Figure 13 shows that  $IBL_{Red}^{Trained}$  had a significantly higher reward against efficient Defenders in later trials (mean: 33.19  $\pm$  33.31 (SD); Tukey's HSD p=0.040). Meanwhile,  $Beeline_{Red}$  had a higher reward against Inefficient Defenders in later trials (mean: 92.18  $\pm$  53.36 (SD); Tukey's HSD p=0.041). These results demonstrate that the  $IBL_{Red}^{Trained}$  strategy remained a challenge for efficient Defenders throughout the experiment. Furthermore, the difficult but deterministic nature of the  $Beeline_{Red}$  strategy was more difficult for inefficient Defenders.

7:18 Y. Du et al.

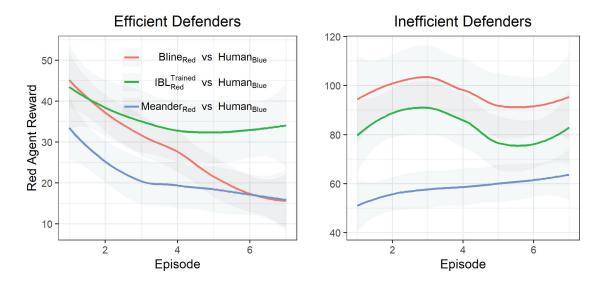


Fig. 13. Average red agent reward by episode, split between efficient human defenders and inefficient human defenders.

An explanation for the poorer performance of efficient Defenders against  $IBL_{Red}^{Trained}$  and inefficient Defenders against  $Beeline_{Red}$  is the defense style used by those groups. The four actions taken in the task can be described as passive (monitor and analyze) or active (remove and restore) [40]. Since  $Beeline_{Red}$  quickly attacks the operational server, it could be that Inefficient Defenders perform worse against  $Beeline_{Red}$ . Similarly, because  $IBL_{Red}^{Trained}$  is both stochastic and adaptive to defender behavior, it could be more difficult for efficient Defenders.

To test this explanation, we compared the proportion of active-type actions performed by efficient Defenders paired with  $IBL_{Red}^{Trained}$ ,  $Beeline_{Red}$ , and  $Meander_{Red}$ . This comparison demonstrated a higher rate of active actions in efficient Defenders paired with  $IBL_{Red}^{Trained}$  than both  $Beeline_{Red}$  (p=0.007) and  $Meander_{Red}$  (p<0.001), but no significant difference between  $Beeline_{Red}$  and  $Meander_{Red}$  (p=0.457). This shows that achieving proficient performance against  $IBL_{Red}^{Trained}$  required more active strategic actions.

#### 7 Discussion

Adversary emulation strategies can be used to train cyber defenders, develop intelligent systems for cyber defense, and test cyber defense capabilities. However, the process of developing effective adversary emulations can be expensive, and their evaluation is often subjective [45, 54]. Existing automated attacker strategies determine how an attacker could achieve a specific goal by assembling vulnerabilities in a graph or by executing a sequence of actions (for example, reconnaissance, escalate, exfiltrate, and lateral movement) in order [23]. Most of them are optimal, but also deterministic and predictable, and therefore can be easily thrown off by human defenders [2]. The first contribution of this article is to demonstrate that it is possible to evaluate intelligent cyber defense systems using cognitive models, aimed at emulating adaptive human decision-making [18]. We show that cognitive models that emulate human adversaries can be better test cases for cyber defenders and for technological capabilities than strategic attack agents that execute a sequence of actions in a fixed order.

Second, we present a cognitive model of an attacker based on IBLT [20],  $IBL_{Red}$ .  $IBL_{Red}$  is first trained against a static and inactive defender,  $Sleep_{Blue}$ . The main feature of  $IBL_{Red}$  is that it can learn from interactive feedback on the task, and we show that it can reach the same level

of effectiveness as the best adversarial strategy in this scenario. This result suggests that an IBL cognitive agent can be an effective dynamic and adaptive emulator of attack behavior. Importantly, the IBL attacker can adapt and learn according to the dynamics of the cyber defense environment. In recent years, **reinforcement learning (RL)** has also made rapid progress as an approach to building adaptive agents [33]. The difference is that RL algorithms are mostly focused on optimally solving computational problems, while cognitive agents focus more on replicating the way humans actually learn [16]. With advances in tool-supported RL agents [24] and human-autonomy teaming [39], defenders are likely to face both optimal adaptive agents and a human attacker. Thus, it is important to train against both types of adaptive attackers.

A third contribution of this article is to demonstrate the performance of an IBL model of the defender IBL<sub>Blue</sub> when paired with three different types of emulated attackers: the optimal attack strategy  $Beeline_{Red}$  the stochastic attack strategy  $Meander_{Red}$  and a cognitive attacker  $IBL_{Red}^{Trained}$ is more difficult for the  $IBL_{Blue}$  defender to learn than an optimal but stable optimal attack strategy. The explanation is that using a cognitive model to emulate attackers is more effective than using deterministic strategies. Cognitive models are dynamic and adaptive to the defender's actions, while the Beeline strategy is static and consistent. The agent  $IBL_{Blue}$  was able to learn the Beeline strategy and eventually take advantage of it, while it did not effectively hinder the progress of the agent  $IBL_{Red}$ . Our analyses show that it takes significantly more steps in time for  $Beeline_{Red}$  to reach the Enterprise subnet and, ultimately, more steps to reach the Operational server. The IBL<sub>Blue</sub> learns over time to prevent these actions from this Beeline<sub>Red</sub> strategy. However, it is significantly more difficult to prevent  $IBL_{Red}$  from reaching the Enterprise and Operational servers. We further verify that there is important learning that occurs from the first to the last episodes in terms of the actions taken by the attacker. For example, the number of impact actions is significantly reduced from the first to the last 500 episodes when the  $IBL_{Blue}$  agent confronts the  $Beeline_{Red}$  strategy, but the reduction in impact actions is minimal when the IBL<sub>Blue</sub> agent confronts the cognitive agent  $IBL_{Red}$ . Exploring the actions taken by the agent  $IBL_{Blue}$  suggests that the agent learns to decrease the restore actions when confronted with the agent Beeline<sub>Red</sub>, while maintaining a more consistent distribution of actions when confronted with the agent  $IBL_{Red}$ . When analyzing the options with which the agent  $IBL_{Blue}$  is confronted at each particular time, we observed an interesting effect: The IBL<sub>Blue</sub> agent learned to reduce its decision option space against Beeline<sub>Red</sub>, while the option space of the  $IBL_{Blue}$  agent against  $IBL_{Red}$  did not decrease substantially.

Finally, we corroborate the simulation predictions in an experiment involving human defenders facing  $IBL_{red}$  and  $Beeline_{Red}$  and  $Meander_{Red}$ . Consistent with the simulation results in  $IBL_{Blue}$ , we observe that the cognitive  $IBL_{Red}$  attacker poses a greater challenge to human defenders than the deterministic  $Beeline_{Red}$  and  $Meander_{Red}$  strategies. In "Interactive Computer-based" cybersecurity skill training, trainees are faced with an adversary and must take a proper course of action or face severe penalties [37]. In existing game training platforms such as CyberCIEGE [25] and CyberNEXS [6], the adversary is played by experienced hackers or emulated by automated attackers. The cognitive attacker agent can serve as another type of training partner in such interactive gaming platforms and mitigate the scarcity of human experts. Combined with techniques that trace the learning progress of trainees, the cognitive attacker agent can provide a personalized training experience for future cyber professionals.

#### 7.1 Conclusion and Limitations

In conclusion, we provide important steps towards establishing emulated adversaries that can effectively train cyber defenders and support the development of autonomous cyber defenders. We demonstrate that it is possible to use cognitive agents to produce adversaries that are adaptive to defenders' actions. These models can ultimately be more effective in learning cyber defense

7:20 Y. Du et al.

strategies than static and deterministic adversaries. In future work, we will also evaluate other types of adaptive agents (e.g., reinforcement learning–based agents) to further investigate the role of attacker adaptivity in human defender learning.

However, the cognitive agent in this work does not fully reflect the decision-making of the human attacker. In addition to learning through experience, human attackers also use various forms of fast and slow reasoning, including heuristics. Such reasoning abilities can improve attackers' efficiency but also make them vulnerable to biased decision-making. Defenders must prepare against experienced attackers who can make fast decisions with heuristics. On the flip side, there is also the potential to exploit these aspects of malicious actors, induce decision-making errors, and reduce the impacts or success of an attack. In future work, we will imbue the cognitive attacker agent with heuristic reasoning abilities. Through training against more human-like attacker agents, defenders can learn to disrupt attacker cognition.

Another limitation of this work is that the participants recruited from MTurk are not necessarily experts in cyber defense. Consistent with previous studies [27, 32], we found that defenders with different expertise have a different learning behavior. In future research, we will control the skill level of defenders by recruiting expert/novices from security operation centers and hacker communities.

Finally, demonstrating the benefits of using cognitive models in real-world cybersecurity environments remains a research challenge. Extending the scenario to the size of real-world networks can exponentially expand the state space in the cognitive model, and research on partially observable states for the defender will be required to account for imperfect network monitoring infrastructures. Future work will aim to improve the game model to be more representative of real-world environments; in particular, we will address the development of a collaborative defense environment to further explore human–AI collaboration in cyber defense.

#### References

- [1] Robert K. Abercrombie, Bob G. Schlicher, and Frederick T. Sheldon. 2014. Security analysis of selected AMI failure scenarios using agent based game theoretic simulation. In 47th Hawaii International Conference on System Sciences. IEEE, 2015–2024.
- [2] Ron Alford, Lukas Chrpa, Mauro Vallati, and Andy Applebaum. 2022. Knowledge reformulation and deception as a defense against automated cyber adversaries. In *The International FLAIRS Conference Proceedings*, Vol. 35. The Florida AI Research Society. https://journals.flvc.org/FLAIRS/article/view/130675
- [3] Adel Alshamrani, Sowmya Myneni, Ankur Chowdhary, and Dijiang Huang. 2019. A survey on advanced persistent threats: Techniques, solutions, challenges, and research opportunities. *IEEE Commun. Surv. Tutor.* 21, 2 (2019), 1851–1877.
- [4] John R. Anderson and Christian J. Lebiere. 1998. *The Atomic Components of Thought*. Psychology Press, New York. 504 pages. DOI: https://doi.org/10.4324/9781315805696
- [5] Andy Applebaum, Doug Miller, Blake Strom, Chris Korban, and Ross Wolf. 2016. Intelligent, automated red team emulation. In 32nd Annual Conference on Computer Security Applications. 363–373.
- [6] Tolulope Awojana and Te-Shun Chou. 2019. Overview of learning cybersecurity through game based systems. In Conference for Industry and Education Collaboration (CIEC'19).
- [7] Callum Baillie, Maxwell Standen, Jonathon Schwartz, Michael Docking, David Bowman, and Junae Kim. 2020. Cyborg: An autonomous cyber operations research gym. *arXiv:2002.10667* (2020).
- [8] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. 2016. OpenAI Gym. *ArXiv* abs/1606.01540 (2016).
- [9] Samuel Chng, Han Yu Lu, Ayush Kumar, and David Yau. 2022. Hacker types, motivations and strategies: A comprehensive framework. *Comput. Hum. Behav. Rep.* 5 (2022), 100167.
- [10] Edward J. M. Colbert, Alexander Kott, and Lawrence P. Knachel. 2020. The game-theoretic model and experimental investigation of cyber wargaming. *J. Defense Model. Simul.* 17, 1 (2020), 21–38.
- [11] Edward A. Cranford, Cleotilde Gonzalez, Palvi Aggarwal, Sarah Cooney, Milind Tambe, and Christian Lebiere. 2020. Toward personalized deceptive signaling for cyber defense using cognitive models. *Topics Cogn. Sci.* 12, 3 (2020), 992–1011.

- [12] Yinuo Du, Baptiste Prébot, Xiaoli Xi, and Cleotilde Gonzalez. 2022. Towards autonomous cyber defense: Predictions from a cognitive model. In *Human Factors and Ergonomics Society Annual Meeting*.
- [13] Varun Dutt, Young-Suk Ahn, and Cleotilde Gonzalez. 2011. Cyber situation awareness: Modeling the security analyst in a cyber-attack scenario through instance-based learning. In *Data and Applications Security and Privacy XXV*, Yingjiu Li (Ed.). Springer Berlin, 280–292.
- [14] Kimberly Ferguson-Walter, Robert Gutzwiller, Dakota Scott, and Craig Johnson. 2021. Oppositional human factors in cybersecurity: A preliminary analysis of affective states. DOI: https://doi.org/10.1109/ASEW52652.2021.00040
- [15] Kimberly Ferguson-Walter, Temmie Shade, Andrew Rogers, Michael Christopher Stefan Trumbo, Kevin S. Nauer, Kristin Marie Divis, Aaron Jones, Angela Combs, and Robert G. Abbott. 2018. *The Tularosa Study: An Experimental Design and Implementation to Quantify the Effectiveness of Cyber Deception*. Technical Report. Sandia National Lab, Albuquerque, NM (United States).
- [16] Cleotilde Gonzalez. 2024. Building human-like artificial agents: A general cognitive algorithm for emulating human decision-making in dynamic environments. *Perspectives on Psychological Science* 19, 5 (2024), 860–873.
- [17] Cleotilde Gonzalez, Palvi Aggarwal, Christian Lebiere, and Edward Cranford. 2020. Design of dynamic and personalized deception: A research framework and new insights. In 53rd Hawaii International Conference on System Sciences. 1825–1834.
- [18] Cleotilde Gonzalez, Noam Ben-Asher, Alessandro Oltramari, and Christian Lebiere. 2014. *Cognition and Technology*. Springer International Publishing, Cham, 93–117. DOI: https://doi.org/10.1007/978-3-319-11391-3\_6
- [19] Cleotilde Gonzalez and Varun Dutt. 2011. Instance-based learning: Integrating sampling and repeated decisions from experience. Psychol. Rev. 118, 4 (2011), 523.
- [20] Cleotilde Gonzalez, Javier F. Lerch, and Christian Lebiere. 2003. Instance-based learning in dynamic decision making. *Cogn. Sci.* 27, 4 (2003), 591–635.
- [21] Sara L. N. Hald and Jens M. Pedersen. 2012. An updated taxonomy for characterizing hackers according to their threat properties. In 14th International Conference on Advanced Communication Technology (ICACT'12). IEEE, 81–86.
- [22] Samuel N. Hamilton and Wendy L. Hamilton. 2008. Adversary modeling and simulation in cyber warfare. In *IFIP International Information Security Conference*. Springer, 461–475.
- [23] Jörg Hoffmann. 2015. Simulated penetration testing: From "Dijkstra" to "Turing test++." In *International Conference on Automated Planning and Scheduling*. 364–372.
- [24] Matthias Hutsebaut-Buysse, Kevin Mets, and Steven Latré. 2022. Hierarchical reinforcement learning: A survey and open research challenges. *Mach. Learn. Knowl. Extract.* 4, 1 (2022), 172–221.
- [25] Cynthia E. Irvine, Michael F. Thompson, and Ken Allen. 2005. CyberCIEGE: Gaming for information assurance. *IEEE Secur. Privac.* 3, 3 (2005), 61–64.
- [26] Chelsea K. Johnson, Robert S. Gutzwiller, Joseph Gervais, and Kimberly J. Ferguson-Walter. 2021. Decision-making biases and cyber attackers. In 36th IEEE/ACM International Conference on Automated Software Engineering Workshops (ASEW'21). IEEE, 140–144.
- [27] Keith S. Jones, Akbar Siami Namin, and Miriam E. Armstrong. 2018. The core cyber-defense knowledge, skills, and abilities that cybersecurity students should learn in school: Results from interviews with cybersecurity professionals. *ACM Trans. Comput. Educ.* 18, 3 (2018), 1–12.
- [28] Hamdi Kavak, Jose J. Padilla, Daniele Vernon-Bido, Saikou Y. Diallo, Ross Gore, and Sachin Shetty. 2021. Simulation for cybersecurity: State of the art and future directions. *J. Cybersecur.* 7, 1 (2021), tyab005.
- [29] Elmar Kiesling, Christine Strauss, Andreas Ekelhart, Bernhard Grill, and Christian Stummer. 2013. Simulation-based optimization of information security controls: An adversary-centric approach. In Winter Simulations Conference (WSC'13). IEEE, 2054–2065.
- [30] Igor Kotenko. 2007. Multi-agent modelling and simulation of cyber-attacks and cyber-defense for homeland security. In 4th IEEE Workshop on Intelligent Data Acquisition and Advanced Computing Systems: Technology and Applications. IEEE, 614–619.
- [31] Pat Langley. 2022. The computational gauntlet of human-like learning. In AAAI Conference on Artificial Intelligence. 12268–12273.
- [32] Clemens M. Lechner, Daniel Danner, Fabian Flöck, Lisa Posch, Arnim Bleier, and Markus Strohmaier. 2019. Measuring motivations of crowdworkers: The multidimensional crowdworker motivation scale. *ACM Transactions on Social Computing* 2, 2 (2019), 1–34.
- [33] Yutaka Matsuo, Yann LeCun, Maneesh Sahani, Doina Precup, David Silver, Masashi Sugiyama, Eiji Uchibe, and Jun Morimoto. 2022. Deep learning, reinforcement learning, and world models. Neural Netw. 152 (2022), 267–275.
- [34] Peter Mell, Karen Scarfone, and Sasha Romanosky. 2006. Common vulnerability scoring system. *IEEE Secur. Privac.* 4, 6 (2006), 85–89.

7:22 Y. Du et al.

[35] Konstantinos Mitsopoulos, Sterling Somers, Joel Schooler, Christian Lebiere, Peter Pirolli, and Robert Thomson. 2021. Toward a psychology of deep reinforcement learning agents using a cognitive architecture. *Topics in Cognitive Science* 14, 4 (2021), 756–79.

- [36] Frédéric Moisan and Cleotilde Gonzalez. 2017. Security under uncertainty: Adaptive attackers are more challenging to human defenders than random attackers. *Front. Psychol.* 8 (2017), 982.
- [37] Ajay Nagarajan, Jan M. Allbeck, Arun Sood, and Terry L. Janssen. 2012. Exploring game design for cybersecurity training. In IEEE International Conference on Cyber Technology in Automation, Control, and Intelligent Systems (CYBER'12). IEEE, 256–262.
- [38] Thuy Ngoc Nguyen, Chase McDonald, and Cleotilde Gonzalez. 2021. *Credit Assignment: Challenges and Opportunities in Developing Human-like AI Agents*. Technical Report. Carnegie Mellon University.
- [39] Thomas O'Neill, Nathan McNeese, Amy Barron, and Beau Schelble. 2022. Human–autonomy teaming: A review and analysis of the empirical literature. *Hum. Factors* 64, 5 (2022), 904–938.
- [40] Baptiste Prebot, Yinuo Du, and Cleotilde Gonzalez. 2023. Learning about simulated adversaries from human defenders using interactive cyber-defense games. arXiv:2304.01142 [cs.CR]
- [41] Baptiste Prébot, Yinuo Du, Xiaoli Xi, and Cleotilde Gonzalez. 2022. Cognitive models of dynamic decision in autonomous intelligent cyber defense. In 2nd International Conference on Autonomous Intelligent Cyber-defence Agents (AICA'22).
- [42] Prashanth Rajivan and Cleotilde Gonzalez. 2018. Creative persuasion: A study on adversarial behaviors and strategies in phishing attacks. *Front. Psychol.* 9 (2018), 135.
- [43] Shabana Razak, Mian Zhou, and Sheau-Dong Lang. 2002. Network intrusion simulation using OPNET. In *OPNET-WORKS Conference*. Citeseer.
- [44] Derek D. Reed, Brent A. Kaplan, and Adam T. Brewer. 2012. Discounting the freedom to choose: Implications for the paradox of choice. *Behav. Process.* 90, 3 (2012), 424–427.
- [45] Lorenzo Russo, Francesco Binaschi, Alessio De Angelis, A. Armando, M. Henauer, and A. Rigoni. 2019. Cybersecurity exercises: Wargaming and red teaming. *Next Gen. CERTs* 54 (2019), 44.
- [46] Barry Schwartz and Andrew Ward. 2004. Doing better but feeling worse: The paradox of choice. Positive Psychology in Practice (2004), 86–104.
- [47] Shishir Kumar Shandilya, Saket Upadhyay, Ajit Kumar, and Atulya K. Nagar. 2022. AI-assisted computer network operations testbed for nature-inspired cyber security based adaptive defense simulation and analysis. Fut. Gen. Comput. Syst. 127 (2022), 297–308.
- [48] Oleg Sheyner, Joshua Haines, Somesh Jha, Richard Lippmann, and Jeannette M. Wing. 2002. Automated generation and analysis of attack graphs. In *IEEE Symposium on Security and Privacy (SP'02)*. IEEE, 273–284.
- [49] Maxwell Standen, Martin Lucas, David Bowman, Toby J. Richer, Junae Kim, and Damian Marriott. 2021. CAGE challenge 1. In IJCAI-21 1st International Workshop on Adaptive Cyber Defense.
- [50] Maxwell Standen, Martin Lucas, David Bowman, Toby J. Richer, Junae Kim, and Damian Marriott. 2021. CybORG: A gym for the development of autonomous cyber agents. arXiv:2108.09118v1
- [51] Paul Theron, Alexander Kott, Martin Drašar, Krzysztof Rzadca, Benoît LeBlanc, Mauno Pihelgas, Luigi Mancini, and Agostino Panico. 2018. Towards an active, autonomous and intelligent cyber defense of military systems: The NATO AICA reference architecture. In *International Conference on Military Communications and Information Systems (ICM-CIS'18)*. IEEE, 1–9.
- [52] Vladislav D. Veksler, Norbou Buchler, Claire G. LaFleur, Michael S. Yu, Christian Lebiere, and Cleotilde Gonzalez. 2020.
  Cognitive models in cybersecurity: Learning from expert analysts and predicting attacker behavior. Front. Psychol. 11 (2020), 1049.
- [53] Robert L. West and Christian Lebiere. 2001. Simple games as dynamic, coupled systems: Randomness and other emergent properties. *Cogn. Syst. Res.* 1, 4 (2001), 221–239.
- [54] Jeong Do Yoo, Eunji Park, Gyungmin Lee, Myung Kil Ahn, Donghwa Kim, Seongyun Seo, and Huy Kang Kim. 2020. Cyber attack and defense emulation agents. *Appl. Sci.* 10, 6 (2020), 2140.
- [55] Jing Zhang, Jun Zhuang, and Victor Richmond R. Jose. 2018. The role of risk preferences in a multi-target defender-attacker resource allocation game. *Reliab. Eng. Syst. Safety* 169 (2018), 95–104.

Received 31 January 2023; revised 26 August 2024; accepted 15 December 2024

1



### Research paper

# Learning about simulated adversaries from human defenders using interactive cyber-defense games

Baptiste Prebot , Yinuo Du, Cleotilde Gonzalez 6\*

Social and Decision Sciences Department, Carnegie Mellon University, Pittsburgh, PA 15213, United States

\*Corresponding author. Social and Decision Sciences Department, Carnegie Mellon University, Pittsburgh, PA 15213, United States. E-mail: coty@cmu.edu

Received 19 September 2022; revised 16 June 2023; accepted 11 September 2023

#### **Abstract**

Given the increase in cybercrime, cybersecurity analysts (i.e. defenders) are in high demand. Defenders must monitor an organization's network to evaluate threats and potential breaches into the network. Adversary simulation is commonly used to test defenders' performance against known threats to organizations. However, it is unclear how effective this training process is in preparing defenders for this highly demanding job. In this paper, we demonstrate how to use adversarial algorithms to investigate defenders' learning using interactive cyber-defense games. We created an Interactive Defense Game (IDG) that represents a cyber-defense scenario, which requires monitoring of incoming network alerts and allows a defender to analyze, remove, and restore services based on the events observed in a network. The participants in our study faced one of two types of simulated adversaries. A Beeline adversary is a fast, targeted, and informed attacker; and a Meander adversary is a slow attacker that wanders the network until it finds the right target to exploit. Our results suggest that although human defenders have more difficulty to stop the Beeline adversary initially, they were able to learn to stop this adversary by taking advantage of their attack strategy. Participants who played against the Beeline adversary learned to anticipate the adversary's actions and took more proactive actions, while decreasing their reactive actions. These findings have implications for understanding how to help cybersecurity analysts speed up their training.

Keywords: cyber defense; human behavior; cyber adversary; interactive games; training

#### Introduction

The rapidly evolving attack capabilities to deploy increasingly sophisticated cyberattacks of unprecedented speed and scale require well-trained cybersecurity experts (i.e. defenders, analysts) [1, 2]. Cyber analysts are responsible for protecting an organization's computer network and digital assets. The job of these defenders consists of a wide variety of network-dependent tasks, including the examination of a large number of alerts to identify intrusion activities and determine whether a network is under attack, the detection of flaws in the organization's security, the development of appropriate protections, and, of course, the mitigation of threats. These activities often include making time-sensitive decisions that may involve disrupting the organization's work in order to protect their information.

Typically, cyber wargaming and adversary simulation are used to evaluate defense algorithms and strategies and to train defenders against new threats [3, 4]. Wargaming exercises mimic a potential threat to an organization by using threat intelligence to define what actions and behaviors an adversary may use. Wargaming emulators build scenarios that capture certain aspects of tactics, techniques, and procedures, to help test the efficacy of defense and identify vulnerability of the network [5]. Also, human defenders are usually recruited to interact with adversarial-simulated scenarios to help them learn from such an interaction [6, 7].

Despite a growing interest in cyber-defense behaviors in recent years [8–12], our understanding of the cognitive demands faced by cyber analysts is still limited [13]. Many factors in adversarial

2 Prebot et al.

behavior may influence defense strategies. For example, the aggressor's personality traits are known to influence their cyberattack behaviors [14, 15]: Machiavellianism was found to be a predictor of stealthy attacks, while narcissism and psychopathy were associated with shorter and more aggressive attacks (i.e. "brute force").

Human-in-the-loop cyber defense laboratory research is required to study both defensive and offensive cyber operations and to develop training protocols tailored to different types of attack strategies [16]. However, conducting meaningful laboratory research with simulated adversaries to study defender behavior is challenging. Participants with the skills and knowledge required to test highly technical tasks and sophisticated adversaries are hard to find and are often too busy to provide their time to test simulated adversaries [9, 17]. The design of simulated adversaries of high fidelity in terms of techniques also requires extensive threat intelligence collected through long-term tracking and clustering of intrusion activities [18]. Given the continuous evolution of network environments and adversaries, it is also unrealistic to derive a future-proof defense strategy at the granularity of current techniques.

To help mitigate this challenge, researchers have been using simulation tools and simplified games [19]. In the context of cybersecurity, these simplified testbeds are used to study the offensive and defensive sides of cyber deception [11, 20], to understand how the general public classifies phishing emails [15, 21], to investigate how the cybersecurity knowledge of the attacker affects the identification of attacks [22], and to study the behavior of the attacker under different levels of uncertainty about the attacker's strategy [23]. In this work, we adopt the Intrusion kill chain model [24] to simplify sophisticated cyberattacks into three tactical phases Establish initial foothold, Propagate through network, and Act on objectives [25]. Consequently, countermeasures such as Monitor, Analyze, Remove, and *Restore* are adopted to disrupt each phase of the attack lifecycle. By pairing defenders with various adversarial strategies constructed with the above tactics, we can learn about the behaviors of human defenders and their processes to address different types of attackers and adapt to dynamic network environments.

However, there is a lack of research on the impact of different adversarial strategies on defense behaviors and the development of defense strategies. Most adversarial cybersecurity games rely on game-theoretic approaches to determine the best defense strategies. These methods often only consider a particular adversary and assume that opponents act "rationally" (i.e. exhibit optimization behavior). These techniques assume the availability of information to adversaries rather than uncertainty and provide individuals with an exact payoff matrix [26, 27]. This leads to a misrepresentation of the reality of the highly dynamic cyber environment, where analysts must work with incomplete and flawed information. While game-theoretic approaches can be useful in determining the optimal defense strategies against known attacks, they provide an unrealistic representation of the attacker's intentions [28-30]; leading to unrealistic representations that might ultimately perform poorly in dynamic cyber-defense environments against unfamiliar adversaries [30-32].

#### Goals and research method

In this research, we address the question of how human defenders behave against different attack strategies and how it affects the emergence of defense strategies. We defined two adversarial strategies in a particular but generic network setting. One adversarial strategy (i.e. Meander) was stealthy; and another one was direct and speedy (i.e. Beeline), reflecting two contrastic attack strategies.

In a recent experiment, ref. [33] proposed a cognitive model based on Instance-Based Learning (IBL) theory [34] that acted as a defender.

This model was paired with both, the Beeline and Meander adversarial strategies to provide predictions of the potential performance of human defenders. The simulation experiment captured the differences in attack strategies and their effect on defenders outcomes. Mainly, the Beeline strategy resulted in a worst performance for the model than the Meander strategy. But it showed that the IBL defender was able to learn over the course of repeated episodes of the defense task. While this is an interesting prediction, human data were not available to validate these observations.

We designed an Interactive Defense Game (IDG) in a cybersecurity scenario and conducted a laboratory study to test human defense behavior against the two adversarial strategies. Similarly to ref. [33], we expect participants who face a Beeline strategy to have more difficulty defending their network against intrusions than participants who face the Meander strategy; and we also expect that humans will learn over the course of repetitions of the defense task.

#### Interactive defense game

The IDG is a web-based interactive cyber-defense game developed to study how human defenders make decisions in a cybersecurity situation. The IDG does not require any installation and can be played remotely using a web browser (Demo of the game: http://janus.hss.cmu.edu:8084/). It provides human participants with a graphical interface to observe network events and analyze the information about a computer network, similar to the way Intrusion Detection Systems (IDS) present network events to human defenders. IDS are common tools to monitor the activities on a network and to help detect possible intrusions or attacks [13].

#### The task of a cyber defender in the IDG

In the IDG, participants play the role of cybersecurity analysts hired by a fictitious manufacturing company to protect their computer network from external malicious activity. The network we use is a simplified version of common corporate network topologies. It is composed of hosts, staff computers, and servers grouped in subnets. Attackers are trying to gain access to the Operational Server (Op\_Server0) to steal information and disrupt production. The easiest way for them to do so is to enter the network through one of the staff computers on the first subnet and progressively make their way up to the critical Op\_Server0 by gaining administrator access to every host on their way.

Each host on which an attacker got administrator-level access costs the defenders some points. The goal of the defender is to minimize the number of points lost.

To perform this task, the defenders use the IDG interface shown in Fig. 1. They must actively monitor the activity of the network to try to identify malicious activity and take actions to block the progression of the attacker. The hosts of the network are characterized by the subnet to which they belong, an IP address, and a host name. Additionally, the system provides the defenders with two dynamic piece of information about each host, the Compromise level and the Activity. When targeting a host, the attacker will first try to gain user-level access to the machine, then try a privilege escalation to gain administrator-level access, and progress to the next target in the network. The Compromise level indicates the status of infection of the host. The second dynamic element provides information about the last Activity detected by the system, like scans or exploitation attempts performed by the attacker on this host. However, not all attacker's activities can be detected by the system. More advanced actions, e.g. privilege escalation attempts and their consequences, are automatically detected. Thus, the defenders have to understand the

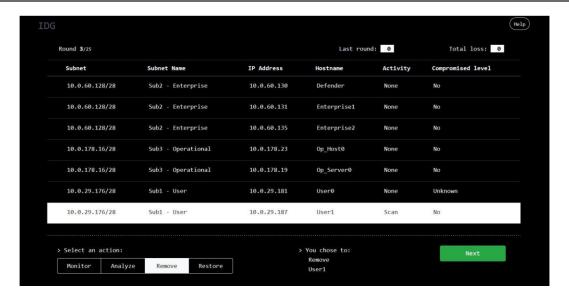


Figure 1. Illustration of the IDG user interface.

observable activity and compromised levels to anticipate future actions of the attackers.

Based on these observable elements, defenders can select among a set of actions represented in buttons on the bottom right of the screen: Monitor, Analyze, Remove, and Restore. Human defenders can select a host by clicking on its row in the table and then choose one of the four actions to perform on that particular host. Only the Monitor action does not require to select a target, it applies to the whole network.

Then, after clicking on the "next" button, the selected action takes effect, and the defender can see the result (i.e. amount of points lost) from the execution of that action in the "last round" value. A new and updated version of the environment is presented to the human defender, demonstrating the new state (activity and compromised levels) of the network elements. The "last round" outcome provides immediate feedback regarding the effectiveness of the past action, and the "total loss" presents the human defender with a cumulative account of the loss during the game. Each game lasts a fixed number of *steps*, each step representing one action.

#### Defense scenario and attack strategies

Human defenders in the IDG are asked to defend a computer network against a red agent. The specific network we used in this scenario is illustrated in Fig. 2.

The network is composed of seven hosts (four computer hosts and three servers) distributed across three subnets. Subnet 1 consists of user hosts that are not critical, subnet 2 consists of enterprise servers designed to support the user activities on subnet 1, and subnet 3 contains the critical operational server and an operational host.

Two types of attack strategy are implemented. They differ by the assumption of the attacker's prior knowledge and illustrate attack behaviors that may result from differences in the attacker's personality traits [14, 15]. In the *Beeline* strategy, attackers route directly through subnet nodes to the Operational Server. The *Meander* strategy does not assume any prior knowledge of the network from the attacker. Attackers following this strategy wonder through the network, trying to gain privileged access to every host in a subnet before advancing further into the network. As a consequence, the Beeline strategy is a direct, rapid, and targeted strategy that can reach the

Operational Server faster than an attacker following the Meander strategy.

The outcome at each step is calculated as shown in Table 1. If the attacker successfully gains administrator access to a user host, the defender loses 0.1 points, while losing administrator access to a server is penalized by -1.0 points. The loss is applied in each step as long as the attacker is not removed from that host or server by the defender. Defenders also receive a negative reward if they have to use the Restore action (-1), because of the important consequences of this action on the system availability. Finally, if the attacker successfully perform the Impact action on the Operational Server, the defender is penalized by -10 points. As Beeline can reach the Operational Server earlier than Meander, it can repeatedly Impact the Operational Server for longer (unless stopped by the defender). As a consequence, and because of the weight accorded to the Impact action, Beeline is potentially more harmful than Meander. For the defender, the implications are a higher theoretical maximum loss against Beeline (-160) than against Meander (-100) (These results are estimated using a completely passive defender. Attackers are able to perform their attack without being disturbed. Beeline then reaches the Operational\_server five steps earlier than Meander).

#### Methods

#### Experimental design

The goal of this experiment is to compare the behavior of human defenders faced with the two types of attack strategy discussed above: *Beeline* and *Meander*.

Given the characteristics of the Beeline strategy that can be faster and more damaging to defenders compared to the Meander strategy, we expected that defenders would initially perform worse against Beeline than against Meander. This hypothesis was preregistered with the Open Science Framework (https://osf.io/u3nfh).

#### **Participants**

Participants were recruited through Amazon Mechanical Turk to participate in a cybersecurity study. The study was advertised to last between 35 and 45 minutes. The time it took across participants was  $M=47.02\pm13.16$  minutes. Participants received a base compensation of \$4.5, and up to \$5.6 in bonus payment ( $M=3.96\pm1.39$ ) based

4 Prebot *et al.* 

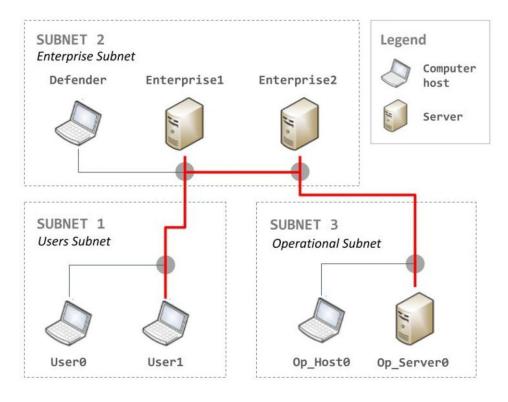


Figure 2. Topology of the network being defended in the IDG scenario. The red line represents the path any attacker needs to take to access the Operational Server.

Table 1. Cost table.

Event or action	Subnet	Point cost
Attacker has administrator access on a Host	Subnet 1, 2, 3	-0.1
Attacker has administrator access on a Server	Subnet 1, 2, 3	-1
Attacker runs IMPACT attack on Operational Server Defender restore an Host or Server	Subnet 3	$-10 \\ -1$

This table was provided to the participants during the instruction phase and was accessible anytime during the experiment through a "help" button.

on their final score (As the score used in this experiment is negative (loss), the bonus payment was calculated by using the difference to the maximum possible loss and attributing 0.005\$ per point: bonus = (total loss + 1120)\*0.005).

A total of 120 participants (89 male, 30 female, 1 N/A) aged 21–65 years ( $M=36.77\pm11.00$ ) completed the study. A total of 12 of the 120 participants (10%) had more than 5 years of experience in the network operation and security area and at least a Master's degree in a related field (In the follow-up survey, participants' expertise was assessed through two likert-scale questions concerning their highest degree in network operation and security, and the years of experience in this area. A one-way ANOVA on those two groups (experts and novices) reveals no effect of the experience on the total losses [F(1, 2.07) = 4.2693, P = .17,  $\eta^2 = .71$ )).

Each participant was randomly assigned to face one of the two adversarial strategies.

#### Procedure

After giving their informed consent and completing a demographic questionnaire, participants received instructions for the task followed by a short quiz to verify their basic understanding of the task instructions, including the network topology, attacker's goal, and the loss calculation process. Participants had to correctly answer all the

questions before moving on to the next step of the experiment. Participants received feedback on the accuracy of their responses and were allowed to modify their responses if they were incorrect. There was no limit in the number of attempts the participants had to answer the questions correctly. However, we recorded the score of their first attempt and the number of times they tried to answer the questions.

Next, participants watched a video introduction to the IDG, explaining the interface, the game controls, and the dynamics of an episode.

Then, participants performed the task consisting of two phases: (1) a practice session and (2) a main task. The practice session consisted of two short episodes (i.e. games) of 10 steps each. The practice episodes were intended to familiarize participants with the interface and game controls. Each of the practice episodes was associated with one of the attacker strategies; however, since the two attack strategies do not differ significantly during the first 10 steps, the participants did not have enough information to discriminate between the two adversarial strategies during the practice session.

Following the practice session, the participants performed the main task consisting of seven episodes of 25 steps each. No time restrictions were imposed. The experimental conditions were kept constant throughout the episodes, which means that each participant played seven episodes against the same adversarial strategy. The

**Table 2.** Descriptive statistics (mean  $\pm$  SD) regarding average loss, number of disruptions, recovery time, and success rate per episode.

	Beeline	Meander
Loss	$-56.12 \pm 50.73$	$-34.76 \pm 30.40$
Disruptions	$0.94 \pm 0.81$	$0.49 \pm 0.52$
Recovery time (steps)	$2.75 \pm 3.55$	$1.31 \pm 1.69$

For contextualization, the maximum loss per episode is -160 against Beeline and -100 against Meander.

initial state of the network was the same for all participants and for each of the episodes.

Subsequently, participants completed a postexperiment survey composed of two parts: (1) feedback on their performance and perceived strategy, and (2) their experience in computer science and cyber defense. Finally, the participants received their final score and were dismissed. The experimental instructions, quiz, and surveys, along with the data and analysis scripts, can be accessed at https://osf.io/u3nfh.

#### Outcome and process metrics

We measured the outcome of the defense performance in the IDG using three metrics:

- Loss: total number of points lost by the defender during the scenario. For reference, the maximum loss per episode resulting from Beeline actions is -160, while it is -100 against Meander.
- **Disruptions:** number of server disruptions that occur within each episode. One disruption represents a set of consecutive steps between a successful impact attack on the Operating Server and the successful recovery by the defender.
- Recovery time: the average number of steps per episode that the defender takes to remove the attacker from the Operational Server after it is disrupted.

We also measured defense process behaviors in addition to defender decisions (i.e. which action is chosen in each step). The attacker actions were also logged for each step and were used to analyze the human behaviors and strategies of defense:

- Proportion of defense actions: number of times that each of the four defense actions—Analyze, Monitor, Remove, and Restore is used by a participant within each episode, divided by the length of the episode (25 steps).
- Proportion of attacker's targets: number of times each host or subnet is being targeted by the attacker within each episode, divided by the length of the episode (25 steps). This is indicative of the attacker's path in the network.
- Proportion of defense strategy: the frequency with which each
  of three coded strategies of defense have been used (*Reactive*,
  Proactive, and Passive) within each episode. Details of calculations of these strategies are presented in the "Defense strategies"
  section below.

#### **Results**

#### Outcome metrics

Table 2 presents the average loss, the number of disruptions, and the recovery time of the participants who played against the Beeline attack strategy and those who faced the Meander attack strategy.

These observations corroborate some expected differences between the two attack strategies in each of the three metrics for outcome performance. In general, the participants lost more points against the Beeline strategy than against the Meander strategy. The average number of disruptions to the operational server within one episode was larger when playing against the Beeline than when playing against the Meander strategy. It also took more steps within an episode to remove the attacker from the operational server when disrupted by the Beeline than the Meander attacker.

We analyzed the outcome metrics over episodes to determine whether the defenders improve with practice against each of the two adversaries. Figure 3 shows the average of each of the three outcome metrics per episode. Generally, we observe more stability over episodes in the participants' outcomes against the Meander adversary than against the Beeline adversary. In other words, the initially poorer performance of participants against a Beeline adversary improves with more practice with this adversary, while the performance of participants against the Meander adversary does not improve much over episodes.

The participants' losses are lower and relatively more stable against the Meander adversary; however, the participants' losses are larger against the Beeline adversary, and they decrease with more practice against this adversary. In addition, the average number of server disruptions is initially higher for participants confronted with the Beeline adversary compared to those confronted with the Meander adversary. However, the number of disruptions decreases with more episodes against the Beeline adversary. A similar result is observed in the average recovery time per episode; where the time is longer for participants playing against the Beeline adversary compared to the Meander adversary, but it decreases with more episodes.

These observations were tested using mixed-effects analysis of variance (ANOVAs) that included the adversary as a betweensubjects factor, the episode as a within-subjects factor, and their interaction. The results for each of the three outcome metrics are reported in Table 3.

Statistical results indicate that the loss, disruptions, and recovery time of the defenders are significantly different when facing the Beeline or Meander adversary. With the exception of average recovery time, we also found consistent significant effects of the episode and the interactions between the adversary and the episode in the Loss and Disruptions.

*Post-hoc* one-way ANOVAs for each of the metrics confirm what we observed in the figure: loss and disruptions improved over the course of episodes *only* when participants confront the Beeline adversary, but not when paired against the Meander adversary. Losses were lower with more episodes only in the Beeline adversary [F(4.29, 278.7) = 7.69, P < .001,  $\eta^2 = .02$ ] but not in the Meander [F(4.12, 214.1) = 1.256, P = .29,  $\eta^2 = .01$ ]; and the number of disruptions decreased only in the Beeline adversary [F(4.93, 320.45) = 10.70, P < .001,  $\eta^2 = .08$ ] and not in the Meander [F(6, 312) = 1.95, P = .07,  $\eta^2 = .02$ ].

The analyses above demonstrate significant differences in defense outcomes when defenders confront Beeline or Meander adversary. The results suggest that Beeline is initially a significantly more damaging attack strategy than Meander. This makes sense by the definition of the strategy, where the Beeline adversary advances directly through the subnets to the operational sever. However, importantly, participants were able to learn the behavior of the Beeline adversary and improve their defense in a way that the loss and number of disruptions improved with more episodes in the task. Participants were more successful against the Meander strategy; however, they were unable to significantly improve their performance with more episodes.

In what follows, we further analyze the process by which participants behaved over the course of the episodes. We analyze the 6 Prebot et al.

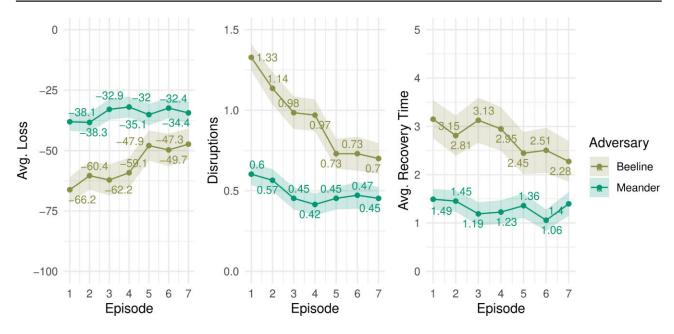


Figure 3. Outcome metrics over time with standard error of the mean. From left to right: loss; disruptions; and recovery time.

Table 3. Results of the mixed ANOVAs regarding the effect of adversary type and episodes on outcome metrics.

Metric		NumDF	DenDF	F-value	P	P significance	$\eta^2$
Loss							
	Adversary	1.00	117.00	8.44	.004	**	.06
	Episode	4.45	520.94	5.99	<.001	***	.01
	Adversary:Episode	4.45	520.94	3.54	.005	**	.01
Disruptions	, ,						
*	Adversary	1.0	117.00	24.24	<.001	***	.10
	Episode	5.1	596.38	10.08	<.001	***	.04
	Adversary:Episode	5.10	596.38	4.34	<.001	***	.02
Recovery time	, ,						
•	Adversary	1.0	117.00	8.87	.004	**	.06
	Episode	4.78	559.48	2.09	.068		.00
	Adversary:Episode	4.78	559.48	1.62	.157		.00

<sup>\*</sup>P <.05, \*\*P <.01, and \*\*\*P <.001.

Table 4. Descriptive statistics (mean  $\pm$  SD) regarding the average proportion of command usage per attacker type.

	Beeline	Meander
Analyze	.20 ±.14	.19 ±.11
Monitor	$.36 \pm .20$	$.30 \pm .19$
Remove	$.32 \pm .19$	$.39 \pm .22$
Restore	.19 ±.09	$.19 \pm .09$

participants proportion of actions, the dynamics of defense actions over time, and characterize their defense strategies. We also explore the individual differences of these behaviors.

#### Process metrics

#### Defense actions

We analyzed the defense actions taken by the participants while executing the task. Table 4 presents the overall average proportion of use of each of the four defense actions—*Analyze*, *Monitor*, *Remove* and *Restore*—in each of the two adversary strategies.

In general, the Monitor and Remove actions seem to be more popular compared to the Analyze and Restore actions among defenders, regardless of the strategy. ANOVAs performed for each adversary group revealed significant differences on the proportion of use of these actions when facing Beeline  $[F(3, 264) = 17.91, P < .001, \eta^2 = .17)$  and when facing Meander  $[F(3, 208) = 18.80, P < .001, \eta^2 = .21]$ . *Post-hoc* comparisons using Tukey's HSD corrections confirm that, regardless of the type of adversary, the proportion of use of Monitor and Analyze; Monitor and Restore; Remove and Analyze; and Remove and Restore were significantly different at P < .001.

Overall, participants in both conditions used Monitor and Remove actions significantly more often than Analyze and Restore (We noted a weak but significant positive correlation between the proportion of Analyze command used and the Cybersecurity background of participants (Spearman rank correlation:  $R_s = .23$ , P = .011). "Expert" subjects seemed to be overly focused on the Analyze action. However, the discussion of this result is beyond the scope of this paper).

To observe the dynamics of the use of these defense actions over the course of episodes, we analyzed the proportions of actions on two levels: (1) across episodes, to observe potential learning and progres-

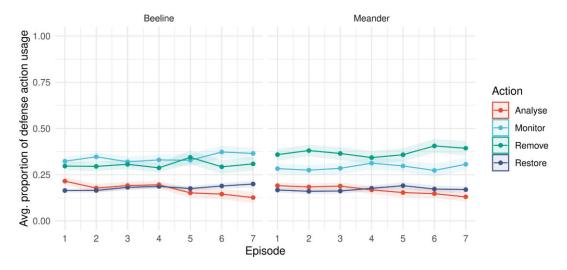


Figure 4. Average proportion of defense action usage over episodes with standard error of the mean.

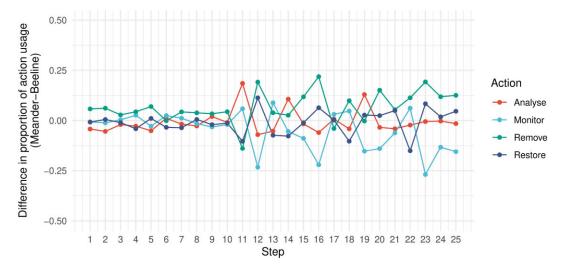


Figure 5. Difference in average proportion of action usage between Meander and Beeline conditions. A positive value indicates a higher proportion of the command in the Meander condition, and a negative one indicates a higher proportion in the Beeline condition.

sive establishment of a defense strategy, and (2) within episodes, aggregating all episodes and analyzing across the 25 steps of episodes.

Figure 4 shows the average proportion of actions over the course of the seven episodes. The defender's behavior appears to be very similar in both adversary strategies across episodes. The main differences observed are that the actions Monitor and Remove are more common than the actions Analyze and Restore. In addition, the action Remove is more common when the defender confronts the Meander than when confronting the Beeline adversary.

However, mixed-effect ANOVAs on the proportion of each of the action types only revealed a significant effect of the episode on the proportion of Analyze action [ $F(4.33, 506.54) = 8.318, P < .001, \eta^2 = .02$ ] when playing against the Beeline and also the Meander adversaries. No effects of the type of adversary were found for any of the actions.

We also analyzed the proportion of actions performed at each step over all episodes. To highlight the differences between the two adversaries, we calculated the difference between the proportion of actions taken by participants facing the Meander opponent and the proportion of actions taken by participants facing the Beeline opponent. Figure 5 presents this difference.

We observe a larger number of Remove actions initially in the Meander compared to the Beeline, and the larger number of Analyse actions in the Beeline compared to Meander in the first 10 steps. The difference in the proportion of actions is relatively consistent and stable during the first 10 steps. However, after step 10, we observe significant variability in this difference of the proportion of actions, noticing that the participants against the Beeline adversary engage in more Monitor actions than those playing against the Meander.

The proportion of actions against Beeline and Meander was tested for each type of action during steps 1–10, and then during steps 11–25. Table 5 indicates that the only significant difference is in the proportion of Monitor and Remove actions during steps 11–25. The proportion of Monitor actions for participants who confronted the Beeline strategy was higher than those who confronted the Meander strategy. Also, the proportion of Remove actions for participants who confronted the Meander strategy was higher than those who confronted the Beeline strategy.

To explain these defense behaviors within episodes, we analyzed the types of targets that each of the adversarial strategies attacked in each of the steps aggregated across all episodes. Figure 6 represents 8 Prebot et al.

Table 5. Results of the ANOVA regarding the effect of adversary type in groups of steps 1-10 and 11-25.

	Command	NumDF	DenDF	F-value	P	P significance	$\eta^2$
1–10							
	Analyze	1.00	686.40	3.53	.06		.08
	Monitor	1.00	670.47	0.08	.784		.03
	Remove	1.00	610.51	2.61	.107		.07
	Restore	1.00	685.28	0.27	.601		.04
11-25							
	Analyze	1.00	1014.13	0.08	.78		.03
	Monitor	1.00	1016.06	38.80	<.001	***	.23
	Remove	1.00	992.60	24.47	<.001	***	.20
	Restore	1.00	1025.17	1.72	.191		.05

<sup>\*\*\*</sup>P <.001.

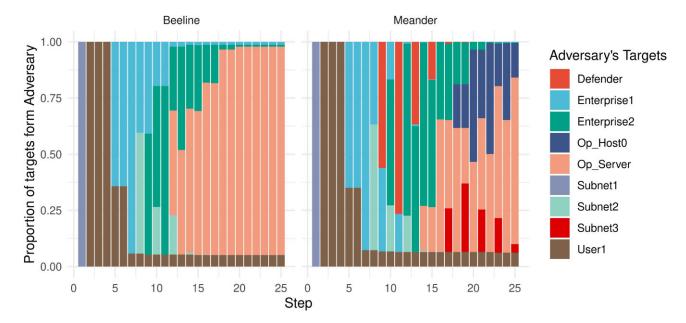


Figure 6. Evolution of the proportion of attack by target across steps.

the proportion of targets that each of the adversaries attacked on each step.

We observe that both adversaries start by attacking Subnet 1, then move to User 1, then to Enterprise 1, and then to Subnet 2. This similarity of adversarial actions appears during the first eight steps of the game. After these steps, Meander starts to target different hosts, such as "Defender," while Beeline moves on to Enterprise 2 and then directly to the Operational Server. This illustration explains the differences in the two attack strategies and explains why the human defenders' actions vary after step 10 and differs in the Monitoring and Removing actions during steps 11–25.

The analysis of defense actions provides evidence of an evolution in the dynamics of defender's decisions throughout the experiment. We propose that this behavior is the result of the participant's learning to defend against their opponent, which explains the performance improvement observed in Fig. 3. There are at least two possibilities to evaluate whether participants improved their understanding of the opponent's strategy and the optimality of their decisions.

Given the sequential nature of the game, the optimality of a decision at a specific step should be defined based on the effect that each particular action will have in *future* steps. Thus, it is possible but not trivial to calculate such "optimal" decision at each step. For each action taken by each individual participant, one would need to calculate the sequence of 25-n actions that would result in the lowest

loss by the end of the episode at each step n and for all the future steps. This is a computationally intensive model and not a trivial optimization algorithm that we considered but decided not to pursue.

Instead, we looked to characterize defense strategies and developed a set of defense heuristics, that may inform human behavior.

#### Defense strategies

To capture the level of understanding of the opponent's strategy and to identify defense actions that would be cognitively plausible, we developed a set of defense heuristics and classified the defense actions into three groups of strategies: *Reactive*, *Proactive*, and *Passive* strategies.

In the cyber literature, *proactive* and *reactive* strategies usually refer to the general approach institutions have for their cybersecurity, i.e. anticipating future threats versus patching security flaws that could expose them to known threats [35–38]. Here, as we focus on the operational level rather than the organizational one, we categorized each individual decision and action according to the following definition:

The passive strategy represents defense actions that have no direct effect on the state of the network or slowing or stopping the progress of the adversary in the network.

Table 6. Heuristics.

Behavior	Strategy
Recovering a compromised host at the user or administrator level	Reactive
Recovering the Operational Server when it is impacted	Reactive
Blocking an initial Impact attempt	Proactive
Preventing a host from being compromised	Proactive
Repeating a successful action	Proactive
Monitoring or Analyzing	Passive

**Table 7.** Descriptive statistics (mean  $\pm$  SD) regarding the average proportion of defense strategy per attacker type.

	Beeline	Meander
Reactive	.27 ±.15	$.26 \pm .16$
Proactive	$.19 \pm .19$	$.15 \pm .20$
Passive	$.48 \pm .22$	$.45 \pm .24$

- The *reactive* strategy represents actions that result in an improved state of the network, such as the recovery of infected hosts. These are actions that the defender takes after hosts have already been attacked by the adversary and defense points have been lost.
- The proactive strategy is characterized by preventive actions.
   These are actions that reflect an anticipation of the next adversarial move or a prediction of the intention of the adversary, in a way that the defender is able to block the progression of the attack.

Table 6 presents the set of high-level heuristics used to categorize defense actions into one of the three strategies. Using the defender action, the state of the network (e.g. is the defender targeting a host that is or has been attacked), and the effect of the defense action, we coded each of these heuristics. Using this coding scheme, 91% of all defender's actions were categorized.

In particular, we characterized proactive actions as a way to determine whether the defenders were ahead of the attacker by choosing the action that would prevent the attacker from doing damage to the network in the future. A repetition of proactive actions reflects an advanced understanding of the opponent's strategy, and explains the learning across episodes.

The overall proportion of reactive, proactive, and passive strategies coded from the defenders' actions when confronted with Beeline and Meander adversaries are presented in Table 7. The table indicates that passive strategies are more common than proactive strategies.

Figure 7 presents the proportion of these strategies per episode. This figure illustrates that passive strategies are most common, regardless of the type of adversary. The proportion of reactive strategies decreases over the course of episodes, while the proportion of proactive strategies increases. This pattern appears to be very similar for both adversaries, although the increase of proactive strategies appears to be faster against the Beeline adversary compared to the Meander adversary.

The mixed-ANOVA results shown in Table 8 indicates a significant effect of the episode on the proportion of each strategy in both types of adversaries. It also shows a significant interaction between the episode and the type of adversary for the proportion of *proactive* strategy.

*Post-hoc* one-way ANOVAs, and considering the Bonferroni adjusted *P*-value (*P*.adj), it can be seen that the simple main effect of Episode on the proportion of Proactive strategy was significant against Beeline [ $F(2.46, 159.66) = 9.152, P.adj < .001, \eta^2 = .04$ ] but not against Meander [ $F(3.11, 161.83) = 2.930, P.adj = .068, \eta^2 = .01$ ].

#### Individual differences

Figure 8 represents the proportion of each strategy fit per episode for each individual participant separately. Furthermore, these panels are organized according the overall loss of each of the participants, where the top-left panel represents the participant with the maximum loss and the bottom-right panel represents the participant with the minimum loss.

This figure immediately reveals the variability in the individual behaviors and the connections between the strategy that each participant used and the individual loss. Many unsuccessful defenders use passive strategies more often, while more successful defenders were more proactive.

#### Strategy and loss correlations

The association between the strategy and the total loss across both adversaries, was also analyzed through correlations. Scatter plots in Fig. 9 represent the relationship between each individual defender's total loss score and the proportion of each strategy.

Spearman's correlation tests indicate a strong significant positive correlation between the participant's loss and the proportion of proactive strategy (Spearman rank correlation:  $R_s = 0.66$ , P < .001). That is, generally, defenders with a higher proportion of proactive behaviors are more likely to lose fewer points, i.e. to protect the network better. Being proactive, such as performing a Remove action that prevents a host from being exploited, is an efficient way to prevent loses and being more successful in protecting the network.

Similarly, Spearman's correlation tests indicate a moderate significant negative correlation between the defender's loss and its proportion of passive strategy (Spearman rank correlation:  $R_s = -0.45$ , P < .001). Defenders with larger number of passive actions were more likely to lose more points since they are not taking any active defense action, i.e. they are not protecting the network.

Finally, the correlation between the defender's loss and the proportion of reactive strategy was not significant.

#### Discussion

We designed a simple cyber-defense game as a web-based application, to study human defense decisions against simulated adversaries. In this experiment, we measured the impact of two different deterministic attack strategies on defenders' behaviors. To do so, we analyzed their performance, their defense choices and behaviors, and their strategies.

As expected, the defenders performance reflects the difference in "aggressiveness" of the attack strategy in terms of Loss, Recovery Time, and number of Disruptions. Indeed, as an attacker following the Beeline strategy was quicker to reach the Operational Server than one following a Meander strategy, it resulted in significantly bigger Loss for the human defender, more Disruptions and longer Recovery Time. However, we have observed that, over the episodes and independently from the condition, participants have managed to improve their performance and lower their Loss. Two possible explanations can be investigated for the overall improvement: (1) the number of Disruptions dropped while subjects learned to more efficiently

10 Prebot et al.

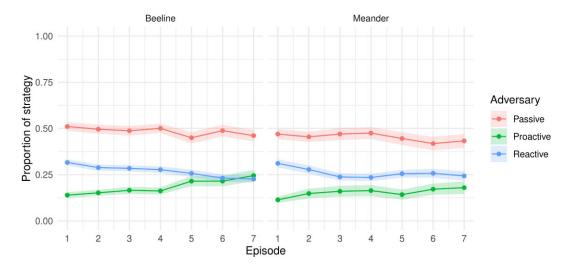


Figure 7. Average proportion of each strategy per episode.

Table 8. Results of the mixed ANOVA regarding the effect of adversary type and episodes on the proportion of defense strategies.

Strategy		NumDF	DenDF	F-value	P	P significance	$\eta^2$
Reactive							
	Adversary	1	117.00	0.18	.675		.00
	Episode	4.15	485.82	8.83	<.001	***	.03
	Adversary:Episode	4.15	485.82	2.30	.550		.01
Proactive	• •						
	Adversary	1	117.00	1.09	.299		.01
	Episode	3.03	354.99	9.23	<.001	***	.02
	Adversary:Episode	3.03	354.99	2.70	.045	*	.01
Passive	· -						
	Adversary	1	117.00	0.66	.417		.00
	Episode	3.73	436.85	3.51	.009	**	.01
	Adversary:Episode	3.73	436.85	1.11	.352		.00

<sup>\*</sup>P < .05, \*\*P < .01, and \*\*\*P < .001.

prevent the attacker from reaching the Operational Server and/or, (2) the Recovery Time improved, i.e. subjects became faster to recover the Operational Server from a disruption.

Results indicate a significant drop in the number of Disruptions recorded over time, while no amelioration is noticeable in terms of Recovery Time. This can be interpreted as the defenders learning to more efficiently block the progression of the attacker in the network, before it reaches the Operational Server.

Overall, participants confronted with a Beeline attacker learned to develop an efficient Proactive defense strategy to improve their performance, be it in terms of loss, number of disruptions, and recovery time. Our interpretation is that, even though both attack strategies are deterministic, Beeline is more direct and consistent, and routing through a smaller number of hosts than Meander. This makes the Beeline strategy easier for the defenders to form a mental representation of, and to predict the adversarial actions with increased defense experience. The predictability of the strategy of attack had a significant influence on how humans learn an effective defense strategy.

Although participants who faced the Beeline adversary seemed to significantly improve their performance over time, they only succeeded to achieve similar level of performance than participants who faced the Meander adversary. In some ways, the Beeline adversary leaves more room for improvement, which could also be a factor in the observed difference in learning pace. In past results involv-

ing experiments with cognitive models on the same task [33], defense agents showed accentuated learning curves when confronted to a Beeline attacker but similar final performance after a large number of episodes. It would be interesting to see how humans are able to improve their strategies and how their performance evolves with more episodes. Also, in future work, longer episodes (i.e. more than 25 steps) could allow us to use patterns identification methods and extended analysis of actions sequences, to refine the categorization of defense strategies and perhaps identify more complex heuristics.

In general, this study illustrates how the type of simulated adversary that human defenders face may influence the speed of learning and the development of an adequate defense strategy. A more aggressive but more predictive attacker was found to be easier to learn and exploit by human defender compared to a stealthy and less predictable adversary.

Cyber analysts have to work in a highly dynamic environment, with flawed and noisy information. Adversarial cyber-defense games and simulation tools like the IDG can help simulate such decision-making situations and better understand the cognitive demands faced by humans cyber defenders.

This experiment also aimed to provide human data to assess the accuracy of human-like IBL defense agents, as presented in refs. [33, 39, 40]. In this context, our work sheds light on the



Figure 8. Proportion of each strategy per subject and episode. Subjects are ordered by Loss. Least performing subject (maximum loss) in the top-left corner. The loss value is displayed above each graph.

12 Prebot et al.

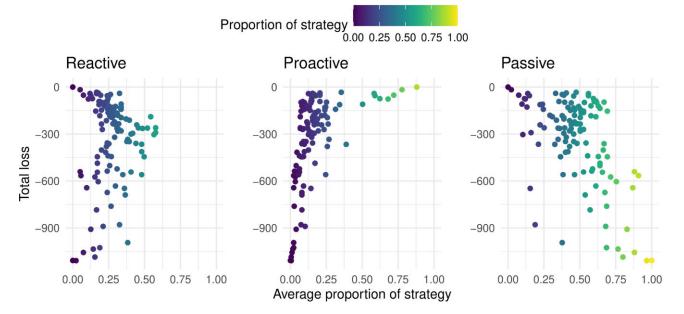


Figure 9. Scatter plot of subject's total Loss and proportion of strategy.

importance of providing less predictive attackers for the development and training of human defenders. These results support the findings of recent modeling experiments that have shown that dynamic attack strategies are a weakness for cognitive models and AI defense [33, 40].

Future work needs to look into the effect of such fully dynamic and adaptive attackers on the human development of defense strategies. We formulate the hypothesis that cognitive dynamic and adaptable attack agents that are able to learn, will present a bigger challenge to defenders, and thus, it provide a better training opportunity for defenders.

This is also a necessary evolution toward more realistic scenarios where expertise brings an advantage. The cybersecurity expertise in particular would be necessary in situations with complex environments and complex tools used in the workplace. In naturalistic settings, the diversification of strategies of attack and their dynamic adaptation to the opponent's actions is indeed more common, and becoming a prominent topic with AI-led cyberattacks.

Because participants with the skills and knowledge required to test highly technical tasks and sophisticated adversaries are hard to find and are often too busy to provide their time to test emulated adversaries, extensive care has been given to design a relevant cyber task that could be performed by a general population.

Future work will aim to improve the task design to be more representative of real-world environments, with an increased complexity of the scenario (e.g. larger networks, simulated regular user activity), by providing more diverse opponents strategies and by introducing teamwork. In particular, we will look into the development of a collaborative defense environment to further explore human—AI collaboration in cyber defense and address some of the challenges of the cyber battlefield of the future.

#### **Acknowledgments**

The authors thank the anonymous reviewers for their valuable suggestions. We thank Jeffrey Flagg, Dynamic Decision Making Laboratory, for research assistance in reviewing and running the study.

#### **Funding**

This research was supported by the Army Research Office and accomplished under Australia–US MURI Grant Number W911NF-20-S-000 and by the Army Research Laboratory under Cooperative Agreement Number W911NF-13-2-0045 (ARL Cyber Security CRA).

#### **Author contributions**

Baptiste Prebot (Conceptualization, Formal analysis, Methodology, Software, Writing – original draft, Writing – review & editing), Yinuo Du (Conceptualization, Software, Writing – original draft, Writing – review & editing), and Cleotilde Gonzalez (Conceptualization, Methodology, Writing – original draft, Writing – review & editing)

Conflict of interest statement: None declared.

#### References

- Li Y, Liu Q. A comprehensive review study of cyber-attacks and cyber security; emerging trends and recent developments. *Ener Rep* 2021;7: 8176–86.
- Thanh CT, Zelinka I. A survey on artificial intelligence in malware as next-generation threats. Mendel 2019:25:27–34.
- Colbert EJ, Kott A, Knachel LP. The game-theoretic model and experimental investigation of cyber wargaming. J Def Model Sim 2020;17: 21–38
- Ferguson-Walter K, Shade T, Rogers A. et al. The Tularosa study: an experimental design and implementation to quantify the effectiveness of cyber deception. Technical report, Albuquerque, NM: Sandia National Lab. (SNL-NM), 2018.
- Applebaum A, Miller D, Strom B. et al. Intelligent, automated red team emulation. In: Proceedings of the 32nd Annual Conference on Computer Security Applications. pp. 363–73, New York, NY, USA: Association for Computing Machinery, 2016.
- Kavak H, Padilla JJ, Vernon-Bido D. et al. Simulation for cybersecurity: state of the art and future directions. J Cybersecur 2021;7:tyab005.
- Varshney M, Pickett K, Bagrodia R. A live-virtual-constructive (LVC) framework for cyber operations test, evaluation and training. In: 2011-MILCOM 2011 Military Communications Conference. Baltimore, MD, USA: IEEE, 2011, 1387–92.

- Gutzwiller RS, Hunt SM, Lange DS. A task analysis toward characterizing cyber-cognitive situation awareness (CCSA) in cyber defense analysts.
   In: 2016 IEEE International Multi-Disciplinary Conference on Cognitive Methods in Situation Awareness and Decision Support (CogSIMA). San Diego, CA, USA: IEEE, 2016, 14–20.
- Veksler VD, Buchler N, LaFleur CG. et al. Cognitive models in cybersecurity: learning from expert analysts and predicting attacker behavior. Front Psychol 2020;11:1049.
- Veksler VD, Buchler N, Hoffman BE. et al. Simulations in cyber-security: a review of cognitive modeling of network attackers, defenders, and users. Front Psychol 2018;9:691.
- Cranford EA, Gonzalez C, Aggarwal P. et al. Towards a cognitive theory of cyber deception. Cogn Sci 2021;45:e13013.
- Johnson CK, Gutzwiller RS, Gervais J. et al. Decision-making biases and cyber attackers. In: 2021 36th IEEE/ACM International Conference on Automated Software Engineering Workshops (ASEW). Melbourne, Australia: IEEE, 2021, 140–4.
- Gonzalez C, Ben-Asher N, Oltramari A. et al. Cognition and technology.
   In: Cyber Defense and Situational Awareness. Cham: Springer, 2014, 93–117.
- Jones DN, Padilla E, Curtis SR. et al. Network discovery and scanning strategies and the Dark Triad. Comput Hum Behav 2021;122:106799.
- Curtis SR, Rajivan P, Jones DN. et al. Phishing attempts among the dark triad: patterns of attack and vulnerability. Comput Hum Behav 2018;87:174–82.
- Gutzwiller RS, Fugate S, Sawyer BD. et al. The human factors of cyber network defense. In: Proceedings of the Human Factors and Ergonomics Society Annual Meeting. Vol. 59. Los Angeles, CA: SAGE Publications, 2015, 322–6.
- Buchler N, Rajivan P, Marusich LR. et al. Sociometrics and observational assessment of teaming and leadership in a cyber security defense competition. Comput Secur 2018;73:114–36.
- 18. Strom BE, Applebaum A, Miller DP. *et al.* Mitre attack: design and philosophy. Technical report. The MITRE Corporation, 2018, Online.
- Gonzalez C, Vanyukov P, Martin MK. The use of microworlds to study dynamic decision making. Comput Hum Behav 2005;21:273–86.
- Aggarwal P, Gonzalez C, Dutt V. HackIt: a real-time simulation tool for studying real-world cyberattacks in the laboratory. In: *Handbook of Computer Networks and Cyber Security*. Cham: Springer, 2020, 949–59.
- Singh K, Aggarwal P, Rajivan P. et al. Training to detect phishing emails: effects of the frequency of experienced phishing emails. In: Proceedings of the Human Factors and Ergonomics Society Annual Meeting. Vol. 63. pp. 453–7, Los Angeles, CA:SAGE Publications, 2019.
- Ben-Asher N, Gonzalez C. Effects of cyber security knowledge on attack detection. Comput Hum Behav 2015;48:51–61.
- Moisan F, Gonzalez C. Security under uncertainty: adaptive attackers are more challenging to human defenders than random attackers. Front Psychol 2017;8:982.
- Hutchins EM, Cloppert MJ, Amin RM. et al. Intelligence-driven computer network defense informed by analysis of adversary campaigns and intrusion kill chains. Vol. 1, Bethesda, MD, USA: Lockheed Martin Corporation, 2011, 80.

- Zhang L, Thing VL. Three decades of deception techniques in active cyber defense-retrospect and outlook. Comput Secur 2021;106:102288.
- Tambe M. Security and game theory: algorithms, deployed systems, lessons learned. Cambridge: Cambridge University Press, 2011.
- Abbasi Y, Kar D, Sintov ND. et al. Know your adversary: insights for a better adversarial behavioral model.In: Proceedings of the 8th Annual Conference of the Cognitive Science Society, Austin, TX: Cognitive Science Society, 2016.
- Aggarwal P, Maqbool Z, Grover A. et al. Cyber security: a game-theoretic analysis of defender and attacker strategies in defacing-website games. In: 2015 International Conference on Cyber Situational Awareness, Data Analytics and Assessment (CyberSA). London, UK: IEEE, 2015, 1–8.
- Nochenson A, Heimann C. Simulation and game-theoretic analysis of an attacker-defender game. In: *International Conference on Decision and Game Theory for Security*. Berlin, Heidelberg: Springer, 2012, 138–51.
- Do CT, Tran NH, Hong C. et al. Game theory for cyber security and privacy. ACM Comput Surv (CSUR) 2017;50:1–37.
- Attiah A, Chatterjee M, Zou CC. A game theoretic approach to model cyber attack and defense strategies. In: 2018 IEEE International Conference on Communications (ICC). Kansas City, MO, USA: IEEE, 2018, 1–7.
- Wang Y, Wang Y, Liu J. et al. A survey of game theoretic methods for cyber security. In: 2016 IEEE First International Conference on Data Science in Cyberspace (DSC). Changsha, China: IEEE, 2016, 631–6.
- Du Y, Prébot B, Xi X. et al. Towards autonomous cyber defense: predictions from a cognitive model. Proc Hum Factor Ergon Soc 2022;66: 1121–5.
- Gonzalez C, Lerch FJ, Lebiere C. Instance-based learning in dynamic decision making. Cogn Sci 2003;27:591–635.
- Grisham J, Samtani S, Patton M. et al. Identifying mobile malware and key threat actors in online hacker forums for proactive cyber threat intelligence. In: 2017 IEEE International Conference on Intelligence and Security Informatics (ISI). Beijing, China: IEEE, 2017, 13–8.
- Bhuyan SS, Kabir UY, Escareno JM. et al. Transforming healthcare cybersecurity from reactive to proactive: current status and future recommendations. J Med Syst 2020;44:1–9.
- Samtani S, Abate M, Benjamin V. et al. Cybersecurity as an industry: a cyber threat intelligence perspective. In: Holt T, Bossler A (eds), The Palgrave Handbook of International Cybercrime and Cyberdeviance, Cham: Palgrave Macmillan. 2020, 135–54.
- Zarreh A, Saygin C, Wan H. et al. A game theory based cybersecurity assessment model for advanced manufacturing systems. Procedia Manuf 2018;26:1255–64.
- Prébot B, Du Y, Xi X. et al. Cognitive models of dynamic decision in autonomous intelligent cyber defense. In: International Conference on Autonomous Intelligent Cyber-defense Agents, Bordeaux, France, 2022
- Du Y, Prébot B, Gonzalez C. A cyber-war between bots: human-like attackers are more challenging for defenders than deterministic attackers. Proceedings of the 56th Hawaii International Conference on System Sciences (HICSS 2023), In: Bui T (ed), Honolulu, HI, USA: HICSS Conference Office, 2023.

FISEVIER

Contents lists available at ScienceDirect

## Computers in Human Behavior: Artificial Humans

journal homepage: www.journals.elsevier.com/computers-in-human-behavior-artificial-humans



#### Research paper

# Experimental evaluation of cognitive agents for collaboration in human-autonomy cyber defense teams

Yinuo Du a , Baptiste Prébot b, Tyler Malloy b, Fei Fang a, Cleotilde Gonzalez b

- <sup>a</sup> Department of Software and Societal Systems, Carnegie Mellon University, 4665 5th Ave, Pittsburgh, 15213, PA, USA
- <sup>b</sup> Department of Social and Decision Sciences, Carnegie Mellon University, 4815 Frew Street, Pittsburgh, 15213, PA, USA

#### ARTICLE INFO

#### Keywords: Human-autonomy teaming Cognitive agent Cybersecurity

#### ABSTRACT

Autonomous agents are becoming increasingly prevalent and capable of collaborating with humans on interdependent tasks as teammates. There is increasing recognition that human-like agents might be natural human collaborators. However, there has been limited work on designing agents according to the principles of human cognition or in empirically testing their teamwork effectiveness. In this study, we introduce the Team Defense Game (TDG), a novel experimental platform for investigating human-autonomy teaming in cyber defense scenarios. We design an agent that relies on episodic memory to determine its actions (Cognitive agent) and compare its effectiveness with two types of autonomous agents: one that relies on heuristic reasoning (Heuristic agent) and one that behaves randomly (Random agent). These agents are compared in a human-autonomy team (HAT) performing a cyber-protection task in the TDG. We systematically evaluate how autonomous teammates' abilities and competence impact the team's interaction and outcomes. The results revealed that teams with Cognitive agents are the most effective partners, followed by teams with Heuristic and Random agents. Evaluation of collaborative team process metrics suggests that the cognitive agent is more adaptive to individual play styles of human teammates, but it is also inconsistent and less predictable than the Heuristic agent. Competent agents (Cognitive and Heuristic agents) require less human effort but might cause over-reliance. A post-experiment questionnaire showed that competent agents are rated more trustworthy and cooperative than Random agents. We also found that human participants' subjective ratings correlate with their team performance, and humans tend to take the credit or responsibility for the team. Our work advances HAT research by providing empirical evidence of how the design of different autonomous agents (cognitive, heuristic, and random) affect team performance and dynamics in cybersecurity contexts. We propose that autonomous agents for HATs should possess both competence and human-like cognition while also ensuring predictable behavior or clear explanations to maintain human trust. Additionally, they should proactively seek human input to enhance teamwork effectiveness.

#### 1. Introduction

With advances in computational power, network robustness, cognitive modeling, and machine learning capabilities, a new form of team is on the rise: the human-autonomy team (HAT; McNeese, et al., 2018). HATs are composed of at least one member of a team that meets the definition of an 'autonomous agent', another member is human, and the team members depend on each other to achieve a collective goal. Wynne and Lyons (2018) define an 'ideal' autonomous agent teammate as: 'a highly altruistic, benevolent, interdependent, emotive, communicative and synchronized agent teammate, rather than simply an instrumental tool'. Although existing work has examined the

effects of agent performance (Bansal et al., 2019) and warmth (Harris-Watson, Larson, Lauharatanahirun, DeChurch, & Contractor, 2023) on the effectiveness of HAT, there is increasing recognition that human-likeness offers unique ways to improve human attitude toward agents and facilitates human-agent cooperation (Glikson & Woolley, 2020; Pelau, Dabija, & Ene, 2021; Zhang, Chong, Kotovsky, & Cagan, 2023). Human participants have shown a higher level of trust in the aid of computer agents with a human-like appearance (de Visser et al., 2012; Von der Pütten, Krämer, Gratch, & Kang, 2010), verbal communication (Kulms & Kopp, 2019), and a display of emotion (Boone & Buck, 2003; Kay, Keller, & Lehmann, 2020). However, there is very limited information, empirical investigations and data regarding how

E-mail addresses: yinuod@andrew.cmu.edu (Y. Du), Baptiste.Prebot@ensc.fr (B. Prébot), tylermal@andrew.cmu.edu (T. Malloy), feifang@cmu.edu (F. Fang), coty@cmu.edu (C. Gonzalez).

https://doi.org/10.1016/j.chbah.2025.100148

Received 30 September 2024; Received in revised form 23 March 2025; Accepted 3 April 2025 Available online 19 April 2025

<sup>\*</sup> Corresponding author.

humans and autonomous agents that are human-like at the *cognitive level* collaborate (Gonzalez, 2024). Musick, O'Neill, Schelble, McNeese, and Henke (2021) explored how team composition affects team processes and emergent team cognition. However, their teams consisted of only humans, and used a Wizard of Oz approach to make participants believe their teammates were autonomous agents. In contrast, our study implements genuine human-autonomy teams with actual autonomous agents working alongside human participants, providing empirical evidence about the true human-autonomy interaction rather than merely perceived interaction.

Cognitive architectures and theories of human decision making have made significant progress in emulating human-like behavior in dynamic environments (Ritter, Tehranchi, & Oury, 2019). Unlike typical computational algorithms that aim to make optimal decisions, cognitive architectures adhere to human constraints such as forgetting, limited attention, and bounded rationality (Gonzalez, 2024). Cognitive models such as instance-based learning theory (IBLT; Gonzalez, Lerch, & Lebiere, 2003) have been implemented in various domains and demonstrated similarity to human decision making processes, including repeated binary choice tasks (Gonzalez & Dutt, 2011; Lejarraga, Dutt, & Gonzalez, 2012), sequential decision-making (Bugbee & Gonzalez, 2022), and practical applications such as identifying phishing emails (Cranford, Lebiere, Rajivan, Aggarwal, & Gonzalez, 2019) and making decisions about cyber attacks (Aggarwal et al., 2022). More recently, cognitive agents have demonstrated a human-like theory of mind (Nguyen & Gonzalez, 2022), or the natural ability to predict the intentions and false beliefs of other agents (Geib & Goldman, 2009; Kautz, Allen et al., 1986). ToM has also been shown to be essential for HAT teamwork (Bendell, Williams, Fiore, & Jentsch, 2024)

The questions we pursue in this research are as follows: Do humanlike cognitive agents have any advantage over optimally performing non-cognitive agents in HAT collaborations? How do humans perceive the cooperativeness and trustworthiness of cognitive and non-cognitive agents in HATs? We explore the potential of cognitive theory to build human-like agents for HATs and compare such cognitive agents in a collaborative HAT applied to cyber defense. We developed a cognitive model that represents the human decision process and incorporated this agent into a HAT experiment where humans and autonomous agents interact and collaborate as a team to ensure the security of an computer network. To enable this investigation, we designed the Team Defense Game (TDG), a novel experimental platform that facilitates controlled studies of human-autonomy teaming in cyber defense scenarios. In an online experiment, we compare how a human defends against adversaries with the help of an autonomous defender that acts randomly, uses smart heuristics, or learns from experience interactively in the task (i.e., the cognitive agent). To evaluate team effectiveness (Hackman, 1978; O'Neill, Flathmann, McNeese, & Salas, 2023), we measure the HATs' performance in terms of their ability in a cyberdefense scenario that requires agents to prevent attacks and resolve network issues. In addition to performance in this task, we also measure human perception of agents in terms of trustworthiness and cooperativeness in a post-experiment questionnaire (Kocielnik, Amershi, & Bennett, 2019; Ragot, Martin, & Cojean, 2020). To better understand how different types of autonomous agents lead to different team outcomes, we measure three collaborative metrics during the teamwork process based on the dynamics of the cyberdefense task. At the general level, we hypothesize that the cognitive agent will be the best teammate, capable of cooperative interaction with humans during teamwork, lead to the best team performance, and be perceived as the most trustworthy and cooperative partner by humans. The following review of related work supports this general expectation.

#### 2. Related work

#### 2.1. Human-Autonomy Teaming (HAT) in cybersecurity

Autonomous systems have been used to control cyber operations (Stevens, 2020) and network security challenges (Bécue, Praça, &

Gama, 2021). Deep learning techniques have been used for the detection of anomalies and malware (Tayyab, Khan, Durad, Khan, & Lee, 2022). Bayesian networks have been applied for the identification of attack paths and the correlation of incident intrusion (Albasheer et al., 2022). Game-theoretic methods have been used to model the interaction between the defender and the adversary as security games and offer optimal allocation strategies of defense resources (Fang, 2021). However, it is clear that across cyber defense, many automated components are limited forms of adaptable automation; that is, they have low adaptability, or the adaptability is overly time-consuming. A typical example of the latter is an intrusion detection system in which a user configures alerts but has to manually adjust and maintain those settings and manage potential false alarms. Recently, the application of reinforcement learning has allowed adaptive cyber defense that is flexible to the dynamics of network/system security status (Du, Song, Milani, Gonzales, & Fang, 2022). However, much of this work still uses autonomous systems as decision-support tools. In this type of work, the autonomous system has no agency and is used to give recommendations to humans rather than to work with humans in a HAT collaboration.

Today, cyber analysts are a scarce resource and are often overloaded (Nobles, 2022). Security Operations Centers (SOCs) combat the growing problem of alert fatigue, where the sheer volume of alerts overwhelms SOC analysts and raises the risk of overlooking critical threats (Chhetri et al., 2024), creating ideal conditions for misallocation of attention (Parasuraman, Molloy, & Singh, 1993). To address this challenge and meet the demands posed by sophisticated adversaries and the need for agile responses, autonomous systems must evolve beyond mere recommender systems and operate with higher levels of agency (Linkov et al., 2023). The cyberdefense technology community is beginning to recognize the necessity of building autonomous agents that can act on their own (Kott, 2023).

However, it is essential to explore autonomous agents that can account for the decision-maker's values or specific mission needs. For example, following a cyber attack, an AI-generated decision engine may recommend disabling an application on the compromised computer system. This action may neutralize the threat posed by the compromised system, but could simultaneously endanger a mission, negatively impact a user's ability to perform critical tasks, or allow the adversary to extend the duration or scope of the cyber attack (Linkov et al., 2023). Human experts should stay in the loop to provide intuition, critical thinking, and contextual information by approving or denying recommendations from AI decision engines that may have negative impacts. Sarker, Janicke, Mohammad, Watters, and Nepal (2023).

In summary, autonomous agents should be able to operate with a degree of self-autonomy and self-directed behavior (agency) while at the same time working interdependently with humans to achieve a shared objective. For the partnership to be successful, the potential benefits of HAT must be weighed against foreseeable negative human-autonomy interactions. Unintended results of incorporating autonomous agents that must be addressed include creating more (not less) work for humans, failing to decrease required manpower, deskilling operators, reducing awareness, contributing to accidents, and loss of life (Gutzwiller, Clegg, & Blitch, 2013; Lyn Paul, Blaha, Fallon, Gonzalez, & Gutzwiller, 2019; Strauch, 2017). The critical factors for successful teamwork in HAT applications must be identified.

The growing literature on HATs in domains such as urban search and rescue (Wohleber, Stowers, Barnes, & Chen, 2023) and hospital management (Chiou, Lee, & Su, 2019) has identified some of these critical factors for successful teamwork, but little is known about HATs in cybersecurity. Teams in cybersecurity operations, especially those in 24/7 security operations centers, have specific dynamics (Paul, 2014). Different cyberdefense scenarios also pose unique challenges for humans and autonomous agents working in teams. For example, in incident response and recovery, autonomous agents might focus on information triage. They would leave the task of further analysis and strategic decision-making to humans. In adaptive defense, autonomous

agents can be more efficient in dynamically adjusting security mechanisms based on real-time threat intelligence. Humans would supervise and fine-tune agent decisions only when necessary (Linkov et al., 2023). Due to these varied and unique applications, a synthetic cyber task environment is needed to empirically evaluate HATs with different team compositions in various cyber scenarios. This environment would also provide humans with early exposure to demonstrations of simulated HATs before implementing them in the cyber workforce.

# 2.2. Cognitive models: Computational representation of human decision processes

The ability of autonomous systems to emulate human decision-making can benefit human-autonomy teaming (Jiao, Zhou, Gebraeel, & Duffy, 2020; McNeese, Demir, et al., 2018; Prebot, 2020). As suggested by Gutzwiller, Espinosa, Kenny, and Lange (2018) and Zhang, McNeese, Freeman, and Musick (2021), creating autonomous agents that can work efficiently with humans should involve modeling team interaction and human cognition. It has the potential to ease the coordination of actions, improve trust in autonomous agents, and increase the performance of the team. In fact, models of human cognition have already been used in tutoring systems, playing the role of a "simulated student". Computer tutors using cognitive models of students to build teaching instructions and provide directed feedback were shown to improve student performance by the same amount as conventional methods one third of the time (Anderson, Corbett, Koedinger, & Pelletier, 1995; Ferster, 2022; Matsuda et al., 2013).

Existing methods in human-autonomy teaming typically do not involve modeling the cognitive mechanisms that underlie dynamic decision making (Ren, Chen, & Qiu, 2023). As a result, the actions of these autonomous teammates can be difficult for end users to understand, even if they are theoretically more optimal than the decisions made by more cognitively inspired agents (Li et al., 2023). Overall, the research on decision support systems in cyber defense has put a strong preference for optimal decisions, rather than understandable and human-like decision making (Vegesna, 2023). In this work we investigate the ability of cognitively inspired autonomous agents to integrate with humans into a team, and compare it to more a optimally designed Heuristic model, which is described more fully in the following sections. Another motivation supporting the improved teaming afforded by more cognitively inspired agents, compared to deterministically optimal ones, is that they can better adapt to the natural variation in human behavior and theoretically result in improved teaming performance.

Instance-Based Learning Theory (IBLT) emerged from the need to explain dynamic human decision-making processes, where a sequence of interdependent decisions is made sequentially (Gonzalez et al., 2003). IBLT provides a general algorithm and mathematical formulation of memory retrieval related to the well-known cognitive architecture ACT-R (Anderson & Lebiere, 2014). The theory proposes a representation of decisions and outcomes in the form of instances. In the past decade, cognitive models based on IBLT have been applied to represent the dynamic decision-making process in various domains that require real-time interactivity between models and humans (Nguyen, Phan, & Gonzalez, 2023). With this increased use of IBLT, the application of models to tasks that involve multiple players is also becoming more common. The initial theoretical developments of IBLT in this direction involved two-person game theoretical models (Gonzalez, Ben-Asher, Martin, & Dutt, 2015). More recently, other interesting applications of IBLT models have been proposed, including the ability to predict other agents' goals, beliefs, and intentions through the Theory of Mind reasoning (Nguyen & Gonzalez, 2020). In summary, existing work on learning-based cognitive modeling, specifically IBLT, has shown that these models operate in ways similar to humans. This similarity can enhance trust and understandability in HATs, as it allows more relatable and predictable interactions between human operators and autonomous agents. In addition, researchers can use the same mechanisms to build human-like models of the theory of mind by observing others' behavior. Thus, IBLT provides critical building blocks for modeling shared cognition processes—memory, attention, and reasoning—central to Human-Autonomy Teaming.

In the context of cyber security, IBL models have been developed to represent individual human cyber defense decisions (Dutt, Ahn, & Gonzalez, 2011), human attacker decisions that can inform cyber defense strategies (Cranford, Gonzalez, et al., 2020; Cranford, Gonzalez, Aggarwal, et al., 2020; Gonzalez, Aggarwal, Lebiere, & Cranford, 2020), and end-user phishing classification decisions that can help improve cyber defense (Cranford et al., 2019; Xu, Singh, & Rajivan, 2022). However, there is no existing work that incorporates IBLT for the defense of human-autonomy teams.

#### 2.2.1. Instance-based learning theory

Although both the process and the mechanisms of IBLT have been published in multiple papers, we reproduce the mathematical formulations of the theory here for completeness. The central element of IBLT is the "instance". It represents a unit of memory resulting from evaluating potential choice alternatives. Each decision is stored in an instance, structured with three elements that are built over time: a situation state s which is composed of a set of features f; a decision or action a taken corresponding to an alternative in state s; and an expected utility or experienced outcome x of the action taken in a state. Concretely, for an IBL agent, an option k = (s, a) is defined by action a in state s. At time t, assume that  $n_{kt}$  different instances  $(k_i, x_{ik_it})$  for  $i = 1, ..., n_{kt}$ , associated with k. Each instance i in memory has an Activation value, which represents the ease of retrieving this information from memory (Anderson & Lebiere, 1998). Here, we consider a simplified version of the Activation equation, which only captures recency, frequency, and noise in memory:

$$\Lambda_{ik_{i}t} = \ln \left( \sum_{t' \in T_{ik_{i}t}} (t - t')^{-d} \right) + \sigma \ln \frac{1 - \xi_{ik_{i}t}}{\xi_{ik_{i}t}}, \tag{1}$$

where d and  $\sigma$  are the decay and noise parameters, respectively, and  $T_{ik_it} \subset \{0,\dots,t-1\}$  is the set of previous timestamps in which instance i was observed. The rightmost term represents noise to capture individual variation in activation, and  $\xi_{ik_it}$  is a random number drawn from a uniform distribution U(0,1) at each step and for each instance and option.

Activation of an instance i is used to determine the probability of retrieving an instance from memory. The probability of an instance where a soft-max function defines i:

$$P_{ik_{i}t} = \frac{e^{A_{ik_{i}t}/\tau}}{\sum_{j=1}^{n_{k_{i}}} e^{A_{jk_{j}t}/\tau}},$$
(2)

where  $\tau$  is the Boltzmann constant (i.e., the "temperature") in the Boltzmann distribution. For simplicity,  $\tau$  is often defined as a function of the same  $\sigma$  used in the activation equation  $\tau = \sigma \sqrt{2}$ .

The expected utility of option k is calculated based on *Blending* as specified in discrete choice tasks (Gonzalez & Dutt, 2011):

$$V_{kt} = \sum_{i=1}^{n_{kt}} P_{ik_i t} x_{ik_i t}. \tag{3}$$

The choice rule is to select the option corresponding to the maximum blended value. When the agent receives delayed results, the agent updates the expected utilities using a credit assignment mechanism (Nguyen, McDonald, & Gonzalez, 2021).

The Instance-Based Learning process is formalized in Algorithm 1, which outlines the sequential decision-making procedure of an IBL agent. The algorithm begins with initialization parameters that include a default utility, an empty memory dictionary to store instances, counters to track time and steps, and a flag indicating whether feedback is delayed. For each decision cycle, the agent observes the current state and enters an execution loop where it explores the available options.

Within this exploration, the agent calculates the activation values for each instance using Eq. (1), computes the retrieval probabilities with Eq. (2), and determines the blended values through Eq. (3). The agent then selects the action corresponding to the maximum blended value, executes it, observes the resulting state, and receives feedback. This feedback is stored in memory as a new instance, and the outcomes are updated through a credit assignment mechanism when the feedback is delayed. The process continues until a terminal state is reached or the step limit is exceeded, embodying the core principles of experience-based learning through memory activation and utility maximization.

**Input:** default utility  $u_0$ , a memory dictionary  $\mathcal{M} = \{\}$ , global counter t = 1, step limit L, a flag delayed to indicate whether feedback is delayed.

```
Initialize a counter (i.e., step) l = 0 and observe state s_l
    while s_l is not terminal and l < L do
        Execution Loop
            Exploration Loop k \in K do
                Compute activation values A_i(t) of instances
                 (k_i, T(i)) by Eq: (1)
                Compute retrieval probabilities P_i(t) by Eq. (2)
                Compute blended values V_k(t) corresponding to k by
            end
            Choose an action a corresponding to option
             k_l \in \arg\max\nolimits_{k \in K} V_k(t)
        end
        Take action a, move to state s_{l+1}, observe s_{l+1}, and receive
        Store t into instance corresponding to selecting k_t and
         achieving outcome u_{l+1} in \mathcal{M}
        If delayed is true, update outcomes using a credit assignment
         mechanism
        l \leftarrow l + 1 and t \leftarrow t + 1
    end
until task stopping condition
```

Algorithm 1: Instance-Based Learning Process

#### 3. Human-Autonomy Team Defense

#### 3.1. Paradigm as HAT

A definition of HAT provided by the systematic review of O'Neill, McNeese, Barron, and Schelble (2022) states that a team includes autonomous agents as individual members who are recognized and seen as performing a unique role on the team with humans. Their work identified both interdependence ("acting with other member activities and outcomes") and agency ("involving independence of actions among autonomous agent members") as two major criteria for autonomous agents to be seen as teammates rather than tools. The first criterion speaks to the nature of the interaction, and the second speaks to the perceived capacity of the agent.

Interdependence. Interdependence can arise from the interrelatedness of the task for each team member, the team structure that requires interaction among team members, or the common goal and shared outcomes (Johnson & Johnson, 1989; Ramamoorthy & Flood, 2004). In TDG, human participants are monetarily rewarded based on the outcomes achieved by the human participant and the autonomous agent together. The collaboration structure also creates an interdependence that necessitates an exchange of information (here, the intention of actions) between the autonomous agent and the human, similar to what happens in human-human cyber-teams.

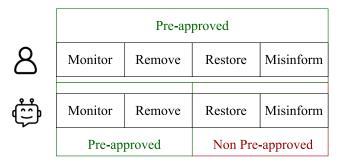


Fig. 1. Differences in autonomy of action between human and agent teammates. The autonomous teammate has to ask the human for approval of the *Restore* and *Misinform* actions.

Agency. One method of measuring the 'agency' of an autonomous agent is the Parasuraman Levels of Automation (LOA) paradigm (Parasuraman, Sheridan, & Wickens, 2000). O'Neill et al. (2022) assert that an autonomous agent in a HAT should at least have partial autonomy, i.e., they can perform actions for themselves instead of suggesting alternatives for humans. In the TDG, the level of autonomous agents falls between level 5, where the agent proposes a course of action but will not enact its decision without human approval ("non-preapproved" actions), and level 7, where the agent chooses and enacts its course of action, while notifying the human teammate.

#### 3.2. Team Defense Game (TDG)

We designed the Team Defense Game (TDG), an online cyber defense game developed to study how humans make decisions in collaboration with an automated agent to defend a network from cyberattacks. TDG is an extension of the interactive defense game (see Prebot, Du, & Gonzalez, 2023) adapted from the CAGE challenge (Standen, Lucas, Bowman, Richer, Kim et al., 2021), a cyber defense competition created to foster autonomous cyber defense research. In TDG, human participants play the role of cyber analysts. They are tasked with protecting the computer network of a fictitious manufacturing company against external malicious activity. Finding critical security incidents among a large number of false alerts generated from separate security products is cognitively demanding and stressful, often leading to frustration and performance degradation (Ban, Samuel, Takahashi, & Inoue, 2021; Dykstra & Paul, 2018; Nobles, 2022). To combat fatigue, human participants are paired with an artificial teammate, an autonomous cyber defense agent, who can make decisions and partially act independently to collaborate with the human defender in a team. Human and autonomous agents must collaborate efficiently to monitor the network, detect suspicious activity, and take appropriate actions to protect the network. The TDG provides human participants with a user interface to observe and analyze network events while interacting with an autonomous teammate protecting the same network and to supervise the actions of the teammate.

In human-human cyber teams (e.g., Cyber Protection Team, Incident Response Team, Security Operation Centers), analysts are provided with a set of pre-approved actions that they can execute without consulting their superior (Boyarchuk, Khudyntsev, Lebid, & Trofymchuk, 2021). If they consider that the situation requires actions that they cannot perform without approval from a superior, they need to submit their intention to their chain of command for validation before acting. It is a common practice to limit or expand the level of autonomy of human analysts through technical controls in the form of group memberships, password protection, firewalls, and even physical access. The idea that an autonomous teammate would have its level of autonomy controlled in a similar way makes it very similar to the technical controls that teams already use for human team members (Hauptman, Schelble,

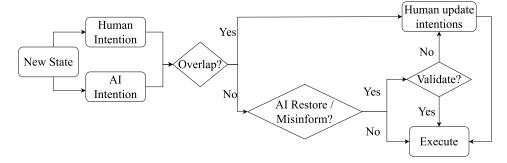


Fig. 2. Interaction flow.

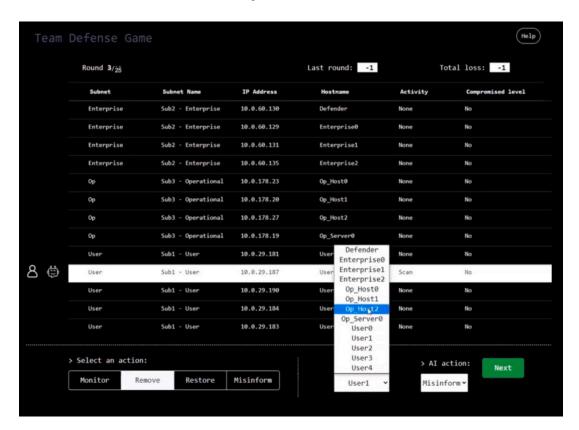


Fig. 3. User Interface example of the TDG.

McNeese, & Madathil, 2023). In TDG, we consider the autonomous agent a cyber analyst with a restricted set of pre-approved actions. As with its human counterparts in real-world settings, the autonomous agent can select an action that is not pre-approved but must submit for approval or modification to its (human) teammate. As shown in Fig. 1, the HAT consists of one human and one autonomous agent, each having a set of pre-approved actions. The agent takes the role of a low level analyst, and the human is the superior in the hierarchy. The agent's pre-approved actions are a subset of the human's pre-approved actions.

As shown in Fig. 2, in each step of the game, the decision to be made by a human analyst is a target (i.e., what computer or server to protect immediately) and an action to take on that target, and the same applies to the autonomous agent. Both make their decision without knowing the intention of their teammate. The human is presented with the agent's intention after he submits his intention. If the agent's intention involves one of the "pre-approved" actions (i.e., Monitor or Remove—see Fig. 1), it is "approved" by default, and the action is performed without the human's involvement. If the agent chooses to perform a "non-pre-approved" action (i.e., Restore or a Misinform), the human is prompted to validate or modify this intention by changing either the

selected action or the target. If the human and autonomous agent select the same target, the human is prompted to resolve this overlap. To do so, they must modify either one of the intentions (human's or agent's). The human decides when to proceed to the next step of the game which will execute the two actions selected by the human and the agent. The agent's decision is executed first. The order in which actions are executed does not affect the effects of these actions.

As an example shown in Fig. 3, the human decides to take the Remove action on host User1; the agent decides to take the Misinform action on the same host User1. In this situation, the human participant is prompted to resolve the conflict. The human participant is allowed to change the agent's action, target, or his own action or target. The team will not be able to move to the next step until the overlap of the human's target and the agent's target is resolved.

#### 3.3. Cyber scenario

As a simulated testbed, TDG allows customization of cyber scenarios. Each scenario consists of a simulated network, an adversary, and a defense team. TDG supports networks with arbitrary topologies

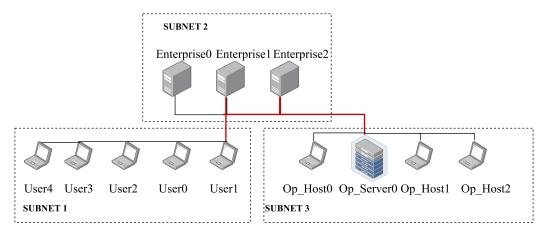


Fig. 4. Network configuration consisting of three subnets.

Table 1
Defender's loss caused by attacker's access and the cost of defense actions.

·		
Event or Action	Subnet	Cost (Points)
Attacker has administrator access on a Host	Subnet 1	-0.1
Attacker has administrator access on a Host	Subnet 2	-1
Attacker runs IMPACT attack on Operational Server	Subnet 3	-10
Defender Restore an Host	_	-1
Defender Misinform (deploy a decoy)	-	-0.5/step (max 5 steps)
Defender disturb a regular User	-	-0.5

and adversaries with various attack capabilities and strategies. In this section, we introduce the cyber scenario used in the human-subject experiment.

In this scenario, the network is a simplified version of a common corporate network. As shown in Fig. 4, the simulated computer network consists of 13 hosts (9 computer hosts and 4 servers) distributed over three subnets. Subnet 1 consists of user hosts that are not critical. Subnet 2 consists of enterprise servers designed to support user activities on Subnet 1. Subnet 3 contains the critical operational server (Server 0), which maintains a service that is key to the operations of the system owners and some other operational hosts. The goal of the adversary is to navigate through the network to the Operational Server (Op\_Server0) to steal valuable information and disrupt the network, which incurs a large cost to the defenders (see Table 1). The adversary algorithm follows a deterministic strategy that assumes prior knowledge of the network layout and is efficient, as it takes the shortest route to the operational server (see the red path in Fig. 4). Following this attack path, the adversary enters the network from a staff computer (i.e., User1) on the first subnet and makes its way to the critical Op\_Server0 by gaining administrator access to every host on their way. Each host on which an attacker has administrator-level access is costly for the defenders (represented in a loss of points).

The goal of the defense team is to minimize the number of points lost. To understand the costs associated with events and actions, the TDG interface provides human participants with a table representation of the computer network. Each element of the network is represented as a row in the table, associated with some static parameters (name, IP address, subnet to which it belongs) and dynamic parameters that represent the state of infection of the host (*Compromise level*) and the last *Activity* detected by the system, such as scans or exploitation attempts performed by the attacker on this host. To add realism and complexity to the task, regular network activity is generated by simulated regular users who perform random scans. Specifically, a regular user scanned a randomly-chosen asset in approximately 5 steps out of the 25 steps in each episode. Therefore, defenders must understand the observable

activity and compromise levels to discriminate suspicious activity and anticipate future actions of attackers. Based on these observable elements, Human participants and Agents in the defense team are given 4 possible actions to choose from in each step: (1) *Monitor* the network (i.e., do nothing), (2) *Remove* user-level adversary access to hosts, (3) *Restore* a system back to a standard configuration which will remove exploited privilege levels, (4) *Misinform* to deploy a decoy which can engage with the attacker to disrupt its operations and delay its progress. *Restoring* a system is guaranteed to remove adversary activity, but it is assumed that restoring a system disrupts the activities of legal users on that system. *Misinforming* with the honey service also requires careful quarantine and consumes computing resources. Thus, using these two commands is costly for the defense team.

The scenario consists of 7 episodes, each composed of 25 steps. To evaluate learning across episodes, at the end of each episode, the network and the attacker's progress are reset to the same initial state. At each step, the defense team takes two defense actions (one by the human and one by the autonomous agent), followed by an action by a regular user of the network. To add realism and complexity to the task, regular network activity is generated by simulated regular users who perform random scans. Specifically, a regular user scanned a randomly chosen asset in approximately 5 steps out of the 25 steps in each episode. Therefore, defenders must understand the observable activity and the levels of compromise to discriminate suspicious activity and anticipate future actions of attackers. Finally, the adversary algorithm observes the network and enacts an attack action. The cost to the defenders at each step is calculated based on the actions taken by the HAT as shown in Table 1, including the cost of losing hosts of various significance and the expense of Restore and Misinform actions. The cost of the previous round and the cumulative loss are displayed to the human participants. The goal of the defense team is to minimize cumulative costs.

#### 4. Autonomous Defender Agents

In this section, we introduce the three types of autonomous agents used in the TDG as partners with a human player, including: a Cognitive agent, a Heuristic agent, and a Random agent.

#### 4.1. Cognitive agent

The Cognitive agent determines the actions to take according to the IBLT algorithm introduced in Section 2.2.1. This agent makes decisions under the guidance of cognitive principles and has been shown to make defense decisions similar to humans (Gonzalez, 2024). The instances represent each decision made and are structured with the following three elements:

State: instance states are constructed to resemble the information that a human defender would have access to and use to determine the action they take. Specifically, there are two attributes for each host or server, representing the observed activity and the known compromised status of that host. The order of (Activity, Compromised Status) pairs for each host is fixed to encode the identity of each host, i.e., the Host name, IP address, and Subnet. The Step Index slot is included to indicate the step counter within each episode.

Action Space: The decision is for the Cognitive agent to choose a host to protect and the tool to protect it. Each action consists of a host and a command in the format of *cmd host*. The action space consists of each of the four actions that target each host in the network for a total of 40 possible actions.

*Utility:* The utility given to the Cognitive agent is the loss of the team from the last step in each step. As shown in Table 1, team loss is the sum of loss caused by attacker's access to the network and the cost of defense actions.

The Cognitive agent accumulates memory of defense decisions and their utility while defending a network against the adversary. The amount of noise added during instance activation computation is set to the default value 0.25. The rate at which activation for previously experienced instances in memory decay with the passage of time is set to the default value 0.5. The memory retention limit is set to 250000 instances. Note that the Cognitive agent is essentially a zero-parameter model, as the decay and noise parameters were not manipulated, fitted, or adapted in any way to its partner's behavior. During training, the cognitive agent gathers experience in an individual version of the TDG, where its actions are automatically approved. This training period lasts 500 episodes, after which the model achieves an average cumulative episode loss of  $68.95 \pm 7.23$  (see Table 1). When the agent makes a decision, during training and the main experiment, it does so by predicting the expected utility of each action that is available to it according to (3), and selecting the maximum utility action. During the training period, the agent successfully learns to defend against the attacker by selecting actions that maximize utility. Additionally, the memory capability of the IBL model allows it to continuously learn and change its action selection based on the utility that the team observes. It is important to note that the Cognitive agent does not directly observe or model its human teammate's specific actions. Rather, it observes: (1) the joint utility achieved by the team, which integrates both players' contributions, and (2) the resulting network states after both players act. Through its instance-based memory, it builds a repository of situation-action-outcome experiences, with more recent and frequent experiences having higher activation according to Eq. (1). This is how the agent indirectly incorporates human influence into its decision-making process without explicit teammate modeling.

```
Input: Current network state NetworkState
Output: Selected defense action
options \leftarrow \emptyset
exploited Hosts \leftarrow GetHostsWithStatus(NetworkState,
 "Exploited")
privEscedHosts \leftarrow GetHostsWithStatus(NetworkState,
 "PrivEsced")
attackPathHosts ← GetHostsOnAttackPath(NetworkState)
options.Add("Monitor")
if exploited Hosts \neq \emptyset then
   mostImportantExploited \leftarrow
    {\tt GetMostImportantHost}(\textit{exploitedHosts})
   options.Add("Remove " + mostImportantExploited)
end
if privEscedHosts \neq \emptyset then
   mostImportantPrivEsced \leftarrow
    GetMostImportantHost(privEscedHosts)
   options.Add("Restore " + mostImportantPrivEsced)
end
if attackPathHosts \neq \emptyset then
   relevant Host +
    GetMostImportantHost(attackPathHosts)
   options.Add("Misinform " + relevantHost)
selected\ Action \leftarrow RandomChoice(options)
return selected Action
Function GetMostImportantHost(hosts):
    Sort hosts by subnet importance (Subnet 3 > Subnet 2 >
    Subnet 1)
   return first host in sorted list
          Algorithm 2: Heuristic Agent Decision Making
```

#### 4.2. Heuristic agent

As a baseline to compare the Cognitive agent, we designed the Heuristic agent that was formalized as a set of rules, which are applied according to the state of the network. The algorithm of the Heuristic agents is shown in Algorithm 2.

At each step, the Heuristic algorithm identifies compromised hosts and randomly selects contextually appropriate defense actions based on the current state of the network. These heuristics assume that an agent has full knowledge of the network structure, the losses associated with each action taken by adversaries, and knows which actions will have the best chance of preventing the progress of the attacker. Furthermore, the performance of Heuristic agents share connections to the expected behavior of reinforcement learning (RL) after long periods of training (Sutton, Barto et al., 1998). After finding the optimal or near-optimal decision strategy during training, the RL models show stable behavior.

Both agents, Cognitive and Heuristic, are dynamic, but the Heuristic agent relies on pre-defined rules, which do not change over the course of an episode of the task. Also, the Heuristic agent is as competent as the cognitive agent when performing the task in isolation, achieving an average cumulative episode loss of (66.20  $\pm$  6.43). However, the Heuristic and Cognitive agents differ fundamentally in their underlying algorithms. The Heuristic agent uses fixed rule-based decision-making, selecting randomly from predetermined correct options based on the network state, while the Cognitive agent employs the IBLT algorithm to accumulate experiences and adapt its decision-making over time according to their cognitive mechanisms.

To illustrate the behavior of the Cognitive and Heuristic agents, Appendix (Behavior Analysis of the Autonomous Agents) shows the distribution of target selection patterns across simulation steps. The appendix suggests that the dynamics of the Cognitive and Heuristic agents are similar but also differ in concrete ways in which the agents

address the attacker's progression over the User, Enterprise, and Operational subnets. This suggests that the Heuristic agent is a good baseline comparison to the Cognitive agent.

#### 4.3. Random Agent

The key similarities between the Heuristic and Cognitive agents are their overall competence in the task and their ability to make decisions that directly depend on the network state of the cyber defense task. To get a more basic baseline, a Random agent was developed to evaluate teaming performance in the TDG, and to compare participants' perception of trustworthiness and cooperativeness.

The Random agent selects actions by choosing a resource as a target at random and then choosing an action to perform at random. All resources and actions have an equal probability of being selected by the random agent. As a result, the performance of the Random agent acting in isolation is considerably lower than that of the IBL or Heuristic models, achieving an average cumulative episodic loss of  $121.39 \pm 49.44$ . Due to this poor individual performance, the Random agent is not expected to achieve a high performance when paired with human participants, though results from this can clarify the relative improvements that are expected by using the Heuristic and cognitive agents.

Just as the Heuristic agent reflects the expected behavior of methods like RL, so does the Random agent. The Heuristic agent behaves similarly to an RL agent at the end of its training lifecycle, while the Random agent resembles an RL agent in the early stages of training. Early-stage RL agents typically exhibit roughly random action selection, which is often used in RL research to compare the performance of trained versus untrained agents (Sutton et al., 1998).

#### 5. Experiment method

#### 5.1. Participants

Participants were recruited through Amazon Mechanical Turk to participate in a cybersecurity study. The study was advertised to last about 60 min. The experiment took M =  $56.20 \pm 15.14$  Minutes on average. Participants received a base compensation of \$6, and up to \$12 in bonus payment based on their final score<sup>1</sup> for a total possible payment of \$18. The average bonus payment was 11.41  $\pm$  2.23. 156 participants (63 female, 90 male, 3 N/A) aged 22-65 years (M =  $39.49 \pm 9.35$ ) completed the study. 66 (27 female, 36 male, 3 N/A) were paired with the Random agent, 48 (16 female, 32 male) with the Heuristic agent, and 42 (20 female, 22 male) with the Cognitive agent. The different number of participants for each condition is caused by random assignment. We remove data from participants who did not fully complete the experiment and the post-experiment questionnaire. Although no formal attention checks were implemented, strict data cleaning was carried out. Participants with more than one missing value and participants whose completion time deviated by more than ±3 standard deviations from the mean were excluded from the analysis. We also excluded participants who showed signs of inadequate engagement, such as taking the same action throughout the game.

#### 5.2. Procedure

First, participants had to complete a demographic questionnaire and provide informed consent. Then, they received instructions for the task and were asked to complete a short quiz to verify their basic understanding of the instructions. The participants had to correctly answer all the questions before moving on to the next step of the experiment. They received feedback on the accuracy of their responses and were allowed to modify their responses if they were incorrect. There was no limit to the number of attempts by participants to answer the questions correctly. However, for control purposes, we recorded the score of their first attempt and the number of times they tried to answer the questions. The participants then watched a video introduction to the TDG explaining the interface, the game controls, and the dynamics of an episode. The participants were then led to a practice session consisting of 1 short episode (that is, a game) of 10 steps. The practice episode was intended to familiarize participants with the interface and game controls and present them with situations in which they must deal with supervision, overlap, and misinformation.

Following the practice session, each participant was randomly assigned to work in conjunction with one of the three types of autonomous agents. The participants performed the main task, which consisted of 7 episodes with 25 steps each. No time restrictions were imposed. The experimental conditions were constant throughout the episodes, which means that each participant played 7 episodes with the same autonomous teammate. The initial state of the network was the same for all participants and for each episode. Subsequently, the participants completed a post-experiment survey consisting of three parts: (1) a collaboration survey, (2) a trust survey, and (3) a background survey about their experience in computer science and cyber defense. Finally, the participants received their final score and were dismissed.<sup>2</sup>

#### 5.3. Dependent measures

A summary of the metrics, their units and their description is shown in Table 2.

#### 5.3.1. Team performance

We measured *team performance* with objective metrics: (1) *Loss*: the average cumulative episodic loss; (2) *Recovery time*: the average number of steps per episode that the adversary successfully impacts the operational server (that is, how long the team takes, on average, to stop an attack occurring on the operational server).

#### 5.3.2. Collaborative process metrics

To evaluate how human and autonomous team members synchronize each other's activities in the interdependent team task, we measure their collaborative process in three situations: (a) *Overlap*, (b) *Supervision*, (c) *Backup*, which have the greatest potential to impact team performance.

Overlap. Overlap refers to situations in which the human player and the autonomous agent choose the same target. A high number of overlaps could suggest a lack of coordination. We measure (1) Frequency of Overlap: the number of overlap cases per episode (2) Adjustment Rate: the rate of two possible types of adjustment human participants could use to resolve the overlap: adjust their own action or adjust the agent's action. Self-adjustment demonstrates the willingness of human participants to adapt to the autonomous agent for cohesive collaboration.

<sup>&</sup>lt;sup>1</sup> As the score used in this experiment is negative (loss), the bonus payment was calculated by using the difference to the maximum possible loss and attributing 0.005\$ per point: bonus=(total loss+1120)\*0.005. 1120 is the maximum total loss for seven episodes

<sup>&</sup>lt;sup>2</sup> The experimental instructions, quizzes, and surveys, along with the data and analysis scripts, can be accessed on Open Science Framework: https://osf.io/xk624/.

Table 2

Metrics summary.		
Metric	Unit/Scale	Description
Team Performance Metrics		
Loss	Points	Average cumulative episodic loss (negative values indicate points lost). Lower values (closer to zero) indicate better performance.
Recovery Time	Steps	Average number of steps per episode required to stop an attacker's impact on the operational server. Lower values indicate faster recovery and better performance.
Human-Agent Interaction Metrics	S	
Frequency of Overlap	Proportion	Rate of instances where human and agent selected the same target per episode (0–1 scale). Lower values indicate better coordination.
Adjustment Rate	Proportion	Rate at which humans adjusted their own actions (vs. the agent's) when resolving target conflicts (0–1 scale). Higher values indicate human adaptability.
Frequency of Supervision	Proportion	Rate of instances per episode where the agent required human validation for non-pre-approved actions $(0-1 \text{ scale})$ .
Agreement Rate	Proportion	Rate at which humans approved the agent's actions without modification during supervision (0–1 scale). Higher values indicate greater trust.
Frequency of Multiple Breaches	Proportion	Rate of instances per episode when multiple hosts were compromised simultaneously (0–1 scale). Lower values indicate better defense.
Backup Rate	Proportion	Rate at which humans attempted to recover breached hosts when multiple breaches occurred (0–1 scale). Higher values indicate better backup behavior.
Human Effort and Efficiency Met	rics	
Human Effort	Actions/Episode	Average number of active actions (Remove, Restore, Misinform) taken by humans per episode. Lower values may indicate less human workload.
Human Efficiency	Points/Action	Loss reduction divided by human effort. Higher values indicate humans achieved better results with fewer actions.

Supervision. Supervision refers to the cases where the agent intends to take an action that requires the validation of its human teammate (i.e., Restore and Misinform). We measure (1) Frequency of supervision: the number of times per episode the Agent needs a validation of their action from the human; (2) Agreement Rate: the rate of times the human participant allows the agent to execute its intended action without modification, out of the total number of times the Agent needs a validation of their action from the human. A high Agreement Rate could indicate that the agent's decision making aligns well with the human participant's judgment and is trusted.

Backup. The challenging task of defending multiple hosts legitimately calls for backing-up behavior. We measure (1) Frequency of Multiple Breaches: the rate of times that more than one host in the network are compromised; (2) Backup Rate: the rate of times that the human participant attempts to recover a breached host, out of the total number of times there are multiple breaches in the network. A high Frequency of Multiple Breaches shows a high demand for both members of the HATs to contribute and recover the compromised host. A high Backup Rate suggests that human participants can rise to the challenge and support their agent teammate when necessary.

#### 5.3.3. Human Effort and Efficiency

Engagement (Sidner, Lee, Kidd, Lesh, & Rich, 2005) is another key process that underlies how effectively autonomous agents can interact with human partners (Holroyd, Rich, Sidner, & Ponsler, 2011; Sidner, Lee, & Lesh, 2003). To evaluate whether human participants are actively and effectively engaged in teamwork, we measure (1) Human Effort: The frequency of humans taking active actions (Remove, Restore, Misinform) rather than passively monitoring per episode; and (2) Human Efficiency: Loss reduction divided by the total human effort during the episode. The loss reduction is calculated by the difference between the observed loss and the maximum loss over all participants in the experiment, to result in a positive value for our Efficiency metric.

#### 5.3.4. Human perception of the autonomous agent

In the post-experiment questionnaire shown in Table 3, we measured *Cooperativeness* and *Trustworthiness*. The cooperativeness and trustworthiness survey questions were inspired by previous studies on automation trust, including discussions in Glikson and Woolley (2020), Schelble et al. (2022). We measured the perceived cooperativeness of the autonomous teammate through a home-made survey composed of 6 items, each rated on a 5-step Likert scale. For the trustworthiness survey we kept all 6 items of Merritt's trust scale (Merritt, Heimbaugh, LaChapell, & Lee, 2013) and adapted the questions to make reference to the "teammate" rather than the "automation". This survey is also based on a 5-step Likert scale. Given the novelty of the current study, where automation is used as a teammate, we could not rely on well established metrics. We used these questions simply to have a subjective metric of cooperativeness and trust, in addition to the objective metrics.

#### 6. Results

In this section, we present the experimental results comparing the three agent types (Random, Heuristic, and Cognitive) in human-autonomy teaming for cyber defense. We analyze team performance metrics, collaborative process patterns, human effort indicators, and participant perceptions of agent trustworthiness and cooperativeness. Table 4 summarizes the descriptive statistics in each of the three agent conditions, which will discuss and test for their significance.

#### 6.1. Team performance

Fig. 5 presents the average loss of the team and the average recovery time of the team per episode for all HATs under the three conditions. The minimum possible loss per episode is 0 and the maximum is 160. As shown in Table 4, the HAT loss is largest when humans are paired

 Table 3

 Cooperativeness and Trustworthiness Questionnaire. Participants had to rate each item on a 5-step Likert

scale.	
Cooperation	My teammate and I coordinated our actions well together My teammate and I coordinated our actions better over the course of the episodes. My teammate and I contributed equally to the team performance I had to carry the weight to make the team better My teammate perceived accurately what task I was trying to accomplish.
	I was able to understand and predict what task my teammate was trying to accomplish.
	I believe my teammate is a competent performer.
	I trust my teammate.
Trust	I have confidence in the choices taken by my teammate.
Trust	I can depend on my teammate.
	I can rely on my teammate to behave in consistent ways.
	I can rely on my teammate to do their best whenever I validate its decision.

 Table 4

 Descriptive Statistics: Mean and Standard Deviation of HATs with three types of agen

Metric	Random	Heuristic	Cognitive
Team Performance			
Loss	M = -79.69	M = -59.69	M = -52.85
	SD = 49.11	SD = 28.54	SD = 27.42
Recovery time	M = 3.92	M = 2.15	M = 1.45
	SD = 3.89	SD = 2.09	SD = 1.89
Human-Agent Interaction			
Frequency ofOverlap	M = 0.11	M = 0.08	M = 0.078
	SD = 0.08	SD = 0.08	SD = 0.09
Adjustment rate(of human actions)	M = 0.14	M = 0.29	M = 0.21
	SD = 0.24	SD = 0.34	SD = 0.33
Frequency ofSupervision	M = 0.49	M = 0.33	M = 0.37
	SD = 0.10	SD = 0.11	SD = 0.13
Agreement Rate	M = 0.35	M = 0.74	M = 0.73
	SD = 0.012	SD = 0.017	SD = 0.02
Frequency ofMultiple Breaches	M = 0.41	M = 0.14	M = 0.14
	SD = 0.02	SD = 0.021	SD = 0.013
Backup rate	M = 0.20	M = 0.17	M = 0.11
	SD = 0.02	SD = 0.02	SD = 0.03
Human Effort and Efficiency			
Human Effort	M = 10.67	M = 9.28	M = 8.38
	SD = 7.13	SD = 7.18	SD = 7.52
Human Efficiency	M = 16.80	M = 29.82	M = 39.40
	SD = 27.02	SD = 38.32	SD = 46.53

with Random agents, comparatively lower with Heuristic agents, and minimal when humans are paired with Cognitive agents. The HAT recovery time follows a similar pattern, with a longer time to recover a breached host when humans are paired with Random agents, a shorter duration with Heuristic agents, and the shortest duration when paired with Cognitive agents. These observations strongly suggest that Cognitive partners in HATs are human's most effective team collaborators. It is worth noting that Fig. 5 shows a sudden drop in loss from episode 4 to episode 5. Upon inspecting the human data, we found that two participants were able to stop the adversary earlier in episode 5 compared to episode 4. This early intervention limited the adversary's opportunities to launch Impact actions, reducing the average loss by approximately 10 points—roughly equivalent to preventing one Impact action. This finding aligns with the observed decrease in recovery time from episode 4 to episode 5.

A two-way mixed measures ANOVA was performed to compare the team loss of HATs with three types of agents over 7 episodes. The results indicated a significant main effect for the type of agent,  $F(2,152)=11.037, p<.05, \eta_G^2=0.086$ . There was no significant interaction between the type of agent and the episode,  $F(11.26,855.40)=0.938, p=0.504, \eta_G^2=0.005$ . There was also no significant effect of the episode  $F(5.63,855.40)=0.938, p=0.821, \eta_G^2=0.001$ . Post hoc tests using Tukey's HSD indicated that HATs with Cognitive agent achieve a

significantly lower loss than HATs with Heuristic agent (p = 0.029) and lower than Random agent (p < .001).

Similarly, the two-way mixed measures ANOVA for the recovery time of HAT indicated a significant main effect for the agent type,  $F(2,152)=18.297, p<.05, \eta_G^2=0.117$  but no effect of the interaction between agent type and episode,  $F(11.38,864.78)=0.95, p=0.744, \eta_G^2=0.003$ ; and no effect of episode,  $F(5.69,864.78)=0.95, p=0.313, \eta_G^2=0.005$ . Post hoc testing using Tukey's HSD indicated that HATs with a Cognitive agent achieve a significantly shorter recovery time than HATs with a Heuristic agent (p=0.007) and shorter with a Random agent (p<0.001).

#### 6.2. Conflict resolution

As illustrated in Table 4, overlap occurs more often in the *Random* condition. A two-way mixed measures ANOVA was conducted to compare the overlap frequency of HATs with three types of agents in 7 episodes. The results confirm a significant effect of the agent type,  $F(2,152)=4.912, p=0.009, \eta_G^2=0.046$ , and Post hoc tests using Tukey's HSD indicated that HATs with Random agent indeed observe significantly more overlap situations than HATs with Cognitive agent (p<.001) and HATs with Heuristic agents (p<.001).

Humans appear to mostly modify the action of agents when there are overlaps, even when they are paired with Heuristic and Cognitive agents. A two-way mixed measures ANOVA was performed to compare the rates of Human self-adjustment with three types of agents over 7 episodes. The results indicated a significant effect for the type of agent, F(2,65) = 4.607, p = 0.013,  $\eta_G^2 = 0.084$ , and a significant effect for the episode,  $F(6,390) = 3.670, p = 0.001, \eta_G^2 = 0.014$ . There was no significant interaction between the type of agent and the episode,  $F(12,390) = 0.662, p = 0.788, \eta_G^2 = 0.005$ . As shown in Table 4, post hoc testing using Tukey's HSD indicated that in overlap situations, when paired with a Heuristic agent, humans adjust their own decision significantly more than when paired with Cognitive (p = 0.015) and Random (p < .01)). When paired with Cognitive agent, humans also tend to adjust their own decision more than when paired with Random (p = 0.14). Further analysis revealed that about 80% of the time, they tend to resolve the conflicts by changing the agent's command to Monitor.

#### 6.3. Supervision

As shown in Table 4, random agents require the most supervision. A two-way mixed ANOVA was performed to compare the frequency of HAT supervisory cases with three types of agents in 7 episodes. The results indicated a significant main effect for the type of agent,  $F(2,152)=76.811, p<.001, \eta_G^2=0.277$ . There was no significant interaction between the type of agent and the episode,  $F(11.49,872.98)=0.957, p=0.133, \eta_G^2=0.011$ . There was also no significant effect of the episode  $F(5.74,8872.98)=0.957, p=0.424, \eta_G^2=0.004$ . Post hoc testing using Tukey's HSD indicated that HATs with Random agents

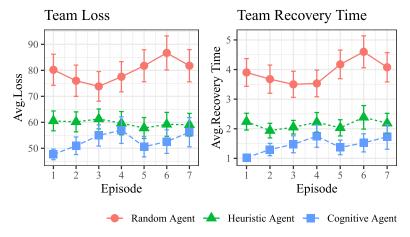


Fig. 5. Evolution of team performance across the 7 episodes of the experiment composed of Team Loss per episode (left) and Team Recovery Time (right).

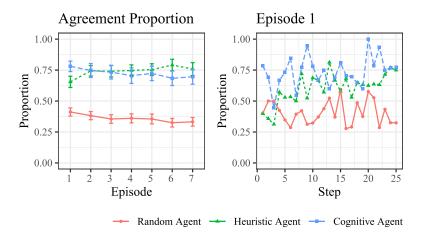


Fig. 6. Evolution of the rate of Human's agreement with Agent's intention during supervisory situations, throughout the first episode (right) and the entire experiment (left).

require significantly more supervision than HATs with Heuristic agents (p < .001) and Cognitive agents (p < .001).

The left panel of Fig. 6 presents the agreement rate throughout 7 episodes. Humans agreed significantly more with the Cognitive and Heuristic agents starting from the first episode. The agreement rate increase in the Heuristic condition is probably due to the consistency and, thus, predictability of the Heuristic agents. A two-way mixed measures ANOVA was conducted to compare the rate of human agreement of HATs with three types of agents across 7 episodes. There was a significant interaction between agent type and episode, F(8.73, 663.49) = $0.728, p < .05, \eta_G^2 = 0.254$ . Post hoc testing using Tukey's HSD indicated that Cognitive agents and Heuristic agents get significantly more agreements than random agents. Further inspection of the agreement rate within episode 1 is presented in the right panel of Fig. 6. By the end of the first episode, the agreement rate increased to approximately 75% in the Cognitive and Heuristic condition and dropped to approximately 35% in the random condition, which indicates that humans are able to observe the competencies of teammates rapidly (Abele, Ellemers, Fiske, Koch, & Yzerbyt, 2021) during the initial exploration episode.

Closer examination of the agreement rate regarding the Restore vs. Misinform actions showed that human participants are able to differentially place their trust in the autonomous agent intentions that vary in reliability. We observed that, independently of the type of teammate, humans tend to agree more with the Restore action (respectively 81%, 80%, and 49% in Heuristic, Cognitive, and Random HATs), which has a deterministic benefit to the team, compared to the Misinform action (71%, 70% and 21% of agreement) that deploy traps which might not catch the adversary successfully every time. At the same time, human participants have the tendency to merge their trust across

actions despite their differences in the actions' reliability. Overtime, the consistent and accurate Restore actions from the Heuristic agent earned itself more trust and led to a higher false agreement rate to its Misinform actions in later episodes.

#### 6.4. Backup behavior

We consider a Backup to be a state of the environment in which there are at least two hosts that are compromised simultaneously. As shown in Table 4, situations requiring backup behavior are more common in HATs with Random agents than with Cognitive or Heuristic agents. A 1-way ANOVA was performed to compare the frequency of multiple breaches between HATs with the three types of agents. The results confirmed a significant effect of the agent type, F(2, 147) = 46.85, p < .001,  $\eta_G^2 = 0.174$ . A post-hoc test using Tukey's HSD indicated that HATs with Random agents get significantly more multiple breach situations than HATs with Cognitive (p < .001) and Heuristic (p < .001).

Also shown in Table 4, Random agents encounter more backup states from their human teammates, with a higher rate of times (20%) when there is more than one breached host. Humans are not good at identifying the need to backup the agents, especially when the agents (i.e., Heuristic and Cognitive agents) appear capable. A two-way mixed measures ANOVA was conducted to compare the rate of backup behaviors of Humans in HATs with three types of agents across 7 episodes. The results indicated a significant effect for the agent type,  $F(2,680)=18.621, p<0.001, \eta_G^2=0.093$ . There was also a significant effect of the episode,  $F(1,680)=5.018, p=0.025, \eta_G^2=0.065$ . The backup rate decreases overtime. There was no significant interaction between the type of agent and the episode, F(2,680)=1.758, p=0.018

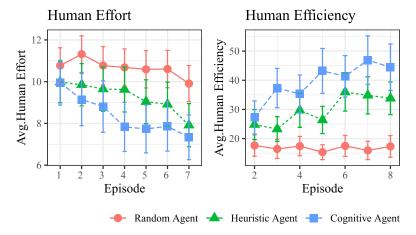


Fig. 7. Evolution of Human Effort (left) and Human efficiency (right) across the experiment.

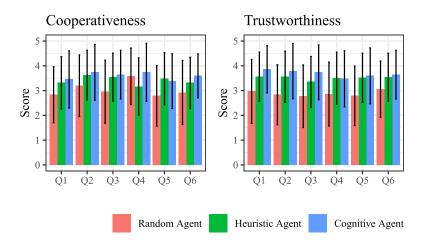


Fig. 8. Left: Cooperativeness Likert Scale, Right: Trustworthiness Likert Scale (1 - Strongly disagree, 2 - Disagree, 3 - Neither agree or disagree, 4 - Agree, 5 - Strongly Agree) - see Table 3.

0.173,  $\eta_G^2=0.002$ . Post hoc testing using Tukey's HSD indicated that Humans paired with a Cognitive agent provide significantly less backup to their teammates than Humans paired with a Heuristic agent (p<.001) and a Random agent (p<.001). Humans contribute the most when paired with Random agents(20%); however, not enough to compensate for the deficiency of the agents.

#### 6.5. Effort and Efficiency

Fig. 7 left panel shows a decrease in the frequency of humans taking active actions (Remove, Restore, Misinform) rather than passively monitoring. A two-way mixed measures ANOVA was conducted to compare human effort in HATs with three types of agents across 7 episodes. In fact, there is a significant effect of episode  $F(3.46,526.09)=5.463,p<0.05,\eta_G^2=0.036$ , where human effort decreases over episodes. There was no significant interaction between the type of agent and the episode,  $F(6.92,526.09)=0.561,p=0.79,\eta_G^2=0.003$  and there was no significant effect of the agent types  $F(2,152)=1.801,p=0.169,\eta_G^2=0.001$ .

Fig. 7 shows a higher efficiency in participants paired with Cognitive agents compared to both Heuristic or Random agents. A two-way mixed measures ANOVA was conducted to compare the human efficiency of the participants in HATs with three types of agents across the 7 episodes. The results indicated a significant main effect for the type of agent,  $F(2,152)=7.949, p<0.001, \eta_G^2=0.145$ , episode  $F(3.46,630.77)=5.322, p<0.001, \eta_G^2=0.037$  and interaction effect

 $F(8.30,630.77)=2.435, p=0.012, \eta_G^2=0.092$ . Tukey's HSD indicated that HATs with Cognitive agents achieve significantly higher human efficiency than HATs with Heuristic agents (p=0.0004) and Random agents (p<0.001).

Although participants' efforts did not show a significant difference between conditions, their efficiency did. This demonstrates the important relationship between autonomous agent strategy and the efficiency of humans in HATs. For HATs with Cognitive teammates, human participants achieved a high level of efficiency. This translates into a relatively low effort from the human participant compared to the improvement in performance observed. One possible explanation for this observed difference is that Cognitive agents are better suited to learning the individual play styles of participants, which can vary from more to less effortful, while the Heuristic and Random agents could not learn these differences. This is evidenced by the fact that in Cognitive and Heuristic HATs, participants begin with nearly the same efficiency on the first episode, which quickly diverges and remains higher for participants in Cognitive HATs.

#### 6.6. Perception of agent cooperativeness and trustworthiness

As shown in Fig. 8, Cognitive and Heuristic agents are rated significantly more cooperative and trustworthy than the Random agent in general. A one-way ANOVA was conducted to evaluate human's perception of the cooperativeness and trustworthiness of the three

types of agents. The results indicated a significant effect of the type of agent on perceived cooperativeness, F (2, 153) = 7.23, p<.05, and trustworthiness F(2, 153) = 10.125, p<.05. 51.9% of participants in the Cognitive condition *Agree* or *Strongly Agree* that the Cognitive agent is cooperative and trustworthy, 49.8% in the Heuristic condition, and 32.9% in the Random condition.

The open-ended feedback in the questionnaire revealed that participants' perception of agents might be affected by their various expectations of the autonomous agent before the experiment. As shown in the following quote, a participant that partnered with a Heuristic agent expected some level of communication with the agent to facilitate decision-making and planning. The absence of such a feature caused frustration. This participant answered *Neither agree or disagree* to most of the survey questions.

"It was challenging because you had to work with an AI that cannot communicate. Made it difficult to come up with a plan or strategize".

Another participant that teamed with a Random agent, on the other hand, has a high level of confidence in the agent's decision-making process. This participant selected *Agree* or *Strongly agree* to all of the survey questions.

"I trusted the partner more than I trusted myself. I was glad it was AI, I assumed they were using an algorithm to make the decisions".

### 7. Discussion

Our results support three general implications for the design of autonomous agents that collaborate with humans in teams: (1) dynamic agents that emulate human-like cognitive processes are beneficial for HAT effectiveness, the benefit can be enhanced through providing explanations when the agents adjust their behavior; (2) competent agents are more trusted and can lead to better HAT performance but might cause over-reliance; (3) human trust in autonomous agents is dynamically shaped by their interactions during teamwork, and should be measured before, during, and after the collaboration.

Human-like cognition. Our first finding concerns the benefits of replicating human cognitive processes in autonomous cognitive teammates. From the comparison between HATs with Cognitive and Heuristic agents (the Cognitive agent achieves the same level of competence as the Heuristic agent after learning), we found that humans who partnered with cognitive agents became more efficient in reducing team loss over time. HATs with Cognitive agents also delivered better team performance. In the post-experiment questionnaire, the Cognitive agent scores slightly higher than the Heuristic agents on trustworthiness. However, there was no significant difference in human trust in autonomous agents. This suggests that human perception of teammate behavior may not be an effective metric for HAT performance. In supervision situations, human participants agreed to approximately 75% agent decisions from the first episode in both conditions. The agreement rates slightly increased over time in the Heuristic condition but slightly dropped in the Cognitive condition. In the presence of conflicts, humans tend to resolve the overlap by adjusting their own actions rather than the agent's actions. One explanation is that humans are less likely to adjust the behavior of the Heuristic agent more because of its consistent behavior. People can easily recognize the 'rules' it follows, making its actions predictable and understandable. In contrast, the Cognitive agent, despite its adaptability to the environment and the human teammate, exhibits less consistency and predictability. In summary, we found that autonomous agents endowed with human-like cognition can improve team effectiveness. However, agents must also employ certain behaviors, such as explanations, to maintain human trust.

Competence. Our second finding concerns the competence of auto nomous teammates. From the comparison between the competent agents (i.e., Cognitive and Heuristic) versus the incompetent Random agent, we found that Humans demonstrate a significantly higher level of trust toward the competent agent, both behaviorally and subjectively. HATs with competent agents were significantly more effective; however, it is worth noticing that humans are not good at identifying the deficiencies and failures of competent agents. Most of the participants did not contribute in situations that require backup behavior. Consistent with previous empirical studies on system-wide trust (Walliser, de Visser, & Shaw, 2023), humans also present a tendency to apply trust broadly rather than specifically to each specific function of the agent. The consistent and accurate Restore actions from the Heuristic agent earned itself more trust and led to a higher false agreement rate to its Misinform actions in later episodes. In sum, the competence of autonomous agents is essential to gain human trust and achieve desirable HAT performance. However, for effective collaboration between humans and autonomous agents, it is crucial that agents actively signal when they require human assistance and transparently communicate the level of uncertainty (Demir et al., 2019) in their decision-making processes (Tomsett et al., 2020).

Trust measurement. Our last insight concerns the measurement of human trust in autonomous agents. First, we observed a significant effect of the episode on the rate of agreement, which echoes the findings of previous research that human trust can be affected by their interaction experience with agents (Kulms & Kopp, 2019). When establishing trust, affective and competency-based dimensions interact to dynamically shape human behavior toward the agent. In teamwork settings, the improved or deteriorated trust caused by previous episodes can affect the interaction in subsequent episodes and spiral into a success or vicious cycle that drags down both team confidence and performance over time. Therefore, it is crucial to regularly assess the different dimensions of trust to adapt the agent's behavior accordingly.

Second, we found that significant differences in agent cognitive ability and competence are more evident in behavioral measures than in self-reported scores. The post-experiment questionnaire showed that some participants had unrealistically high expectations of the agent, which later led to blind trust or great disappointment in their interaction with the agents. We also observed a discrepancy between the actual contribution and the perceived contribution. Many participants Agree or Strongly Agree that "I had to carry the weight to make the team better" while, actually, their agent made a greater contribution to the team. This could be due to the moral assumption that "even if the autonomous agent gains more agency, humans remain responsible" (Cummings, 2014).

Lastly, reinforcing the first two points of the discussion, our findings suggest that a high level of competency, when paired with strong predictability, can lead to over-reliance. We observed that humans tended to align more closely with the strategies of the heuristic agent, occasionally resulting in misplaced trust and reduced performance. In contrast, participants who worked with the cognitive agent were more critical of its decisions, which may indicate a more engaged and adaptive interaction. This suggests that cognitive agents could play a valuable role in dynamic trust calibration, helping to maintain team efficiency while preventing over-reliance. Measurement of behavioral adaptation, such as changes in human oversight, agreement with agent decisions, and response to unexpected agent actions, could then provide richer insights into trust dynamics than static self-reported trust ratings.

In summary, to better leverage human feedback on agents' trust-worthiness and cooperativeness, we recommend refining subjective measures to distinctly capture different dimensions of trust, including propensity to trust (Hoff & Bashir, 2015; Schoorman, Mayer, & Davis, 2007), affective trust, and competency-based trust. Furthermore, trust levels should be unobtrusively monitored during teamwork using behavioral markers, while post-experiment questionnaires should clearly differentiate between taskwork and teamwork contributions.

#### 7.1. Implications on HAT for cybersecurity

In the context of cybersecurity, the system-wide trust problem is more prominent. Autonomous cybersecurity systems tend to integrate more than one function, such as analyzing alerts, validating security controls, and resolving incidents based on standardized procedures and playbooks. Each task can be accomplished with different levels of trustworthiness and might evolve to a different level of risk. Thus, it is important to further investigate system-wide trust mediation mechanisms and calibrate the reliance on the different functionalities of the same autonomous defense agents. This calibration helps prevent over-reliance on potentially fallible automation while leveraging its strengths, ultimately enhancing the safety and effectiveness of HATs. The assignment of credits and responsibilities is also important in settings such as cybersecurity, where each decision can be of high risk. Cyber workforce includes staff from very different specializations trained in cybersecurity. Special caution is needed when this type of staff is involved in human-autonomous teaming. As they may not have a "fair" judgment on the autonomous agents.

## 7.2. Limitations and future work

Human-like cognition. We designed an autonomous agent with human cognition based on instance-based learning theory (IBLT), which has enjoyed many successes in replicating individual human-like behavior in the cyber domain and other contexts (Gonzalez, 2024). However, the definition of human likeness in a team context needs further exploration. Future research may focus on conducting "Turing experiments" in a team setting to better understand the desirable teamwork behavior and human-likeness. It is possible that "human-like" is not equivalent to "ideal autonomous teammate", and we might want the autonomous agents to avoid the limitations of humans in teamwork. For example, humans may prefer to reach consensus rather than propose objections and alternative solutions, which has a negative impact on teamwork effectiveness. Secondly, in our current design of the Cognitive agent, the agent leverages the accumulation of instancebased knowledge, recognition-based retrieval, adaptive strategies, and feedback update mechanisms in IBLT. It is unclear how each individual underlying cognitive mechanism contributes to the human-likeness of agents. In future research, we will investigate the instance-based learning mechanisms in more detail to better understand the link between cognitive mechanisms and desirable human-like teamwork behavior. Finally, the current Cognitive agent treats the human teammate as part of the environment, while humans recognize each other's strengths, weaknesses, and working styles to coordinate their actions and adapt their behavior to other team members. In future research, we plan to develop a Cognitive agent that includes a direct observation of human actions, so that we can give the Cognitive agent the ability to make predictions ahead of time regarding human intentions as in the work on Theory of Mid (Nguyen & Gonzalez, 2022). In addition to investigating specific cognitive mechanisms, future work could explore comparisons with other learning approaches, such as reinforcement learning, which has been applied in cybersecurity contexts (Du et al., 2022). While our current study focused on IBL-based cognitive agents, understanding the relative advantages of different learning frameworks for human-autonomy teaming in cybersecurity would provide valuable insights for the field. Specifically, research could examine how different algorithmic approaches affect team dynamics, human trust, and overall performance in various cybersecurity scenarios and team configurations.

Human-automation communication. Our current interactive online team game only supports the action phase of teamwork, where the agent and human participants collectively protect a network through four defensive actions. Teams, in reality, also engage in another important phase of teamwork: the transition phase. It refers to periods of time

when teams focus primarily on evaluation and/or planning activities to guide their achievement of a team goal or objective. Teams must take a look at how well they performed during the previous episode and prepare for the next. They compare current performance levels with goals and derive performance gaps. Closing these gaps, in combination with current and anticipated future assignments, guides the development of future performance goals and strategies to achieve them. This phase is especially important in adversarial cyber scenarios, where team members need to develop a threat model and collectively develop a strategy against the attacker. In future work, we will enhance the platform with human–agent communication facilities to support a richer HAT teamwork experience.

HAT for cybersecurity. In the experiment reported in this work, we recruited participants with cybersecurity knowledge from Amazon MTurk. We preselected these participants based on their technical background and screened them according to their cybersecurity knowledge, but they are not necessarily experts in cyber defense. Another limitation comes from the simple cyber defense scenario, which involves a single attacker on a small network. Previous research shows that factors such as workload, stress, and risk can affect teamwork. In cybersecurity, high workload, stress caused by the presence of powerful attackers, and risks involved in defense decisions are very common and cannot be ignored. In future work, we will explore game scenarios with various levels of complexity and workload to better understand the unique challenges of HAT in cybersecurity. Our platform is in continuous development as we refine its capabilities and usability. Motivated by the need to develop open source testbeds and platforms for HAT research, we plan to make our development general and provide documentation for others to consume. The improvements will make it easier for researchers to conduct controlled experiments with minimal coding requirements. We are working to create more accessible interfaces, comprehensive documentation and ready-to-use experimental templates that will lower the entry barrier for researchers of diverse backgrounds. This enhanced platform will allow us to conduct experiments that resemble team configurations in real-world settings, such as including more than two members of the human or autonomous agent team. Such configurations would more accurately reflect operational security teams in which multiple analysts and automated tools work together to defend the network infrastructure.

#### 8. Conclusion

In this work, we found that Cognitive agents are indeed better teammates than Heuristic and Random agents. HATs with Cognitive agents are more effective in protecting the network from breaches, recovering exploited hosts faster, and thus losing fewer points in the team task. Human participants partnered with Cognitive agents also showed higher efficiency. On the other hand, the adaptivity of Cognitive agents renders their behavior relatively unpredictable, which can undermine the trust between humans and agents.

This research makes several important contributions to the field of human-autonomy teaming. First, unlike previous studies that have relied on Wizard of Oz approaches in which humans simulate agent behavior (e.g., Musick et al., 2021), our work uses actual autonomous agents that operate independently according to their programmed cognitive or heuristic mechanisms. This methodological advancement allows us to investigate real human-autonomy interactions rather than perceived ones, providing more ecologically valid insights into the dynamics of HATs.

Second, we introduced the Team Defense Game (TDG), a novel experimental platform specifically designed to study human-autonomy teaming in cyber-defense contexts. The TDG provides researchers with a controlled environment to examine various aspects of HAT collaboration, including conflict resolution, supervision dynamics, and backup behaviors. This platform offers significant advantages over existing

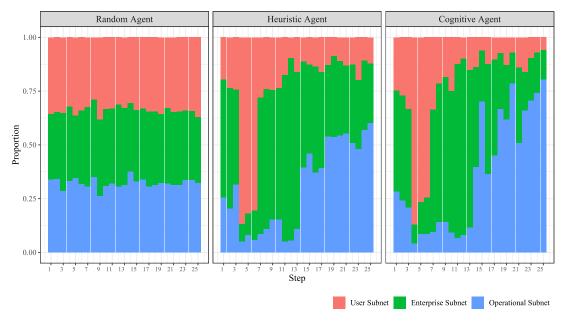


Fig. 9. Distribution of Agent Target Selection Across Steps by Agent Type. The plot shows the proportion of different host types (User, Enterprise, and Op) targeted by Random, Heuristic, and Cognitive agents over 25 steps of the simulation, averaged across all episodes.

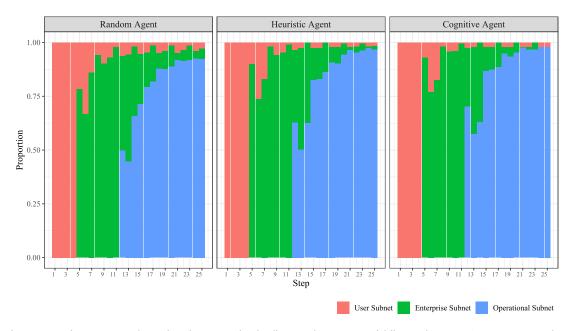


Fig. 10. Temporal Progression of Agent Target Selection by Subnet Type. The plot illustrates the proportion of different subnet types (User, Enterprise, and Operation) targeted by the attacker against Random, Heuristic, and Cognitive agents across 25 simulation steps, averaged across all episodes, demonstrating their strategic priorities throughout the cyber defense scenario.

testbeds by incorporating realistic cyber defense scenarios while maintaining experimental control, making it a valuable resource for future HAT research.

Third, our empirical findings demonstrate that cognitive architectures based on human learning mechanisms offer substantial benefits for HAT effectiveness compared to both random and optimally designed heuristic agents. The cognitive agent's ability to adapt not only to the task environment but also to its human teammate's behavior patterns creates a more complementary partnership, leading to improved team performance despite the potential unpredictability introduced by its adaptive nature.

Using the case of cyber protection teams, this study demonstrated the possibility of effective human-autonomy teamwork and provided evidence of the value of cognitive modeling approaches in agent design.

Our results highlight that the implementation of human-like cognitive processes in autonomous agents represents a promising direction for the development of collaborative AI systems. The balance between adaptability and predictability remains a key challenge, but our findings suggest that the benefits of cognitive adaptation outweigh the costs in terms of overall team effectiveness. As autonomous systems become increasingly integrated into human teams in various domains, these insights can inform the design of agents that function not only as tools but as genuine teammates capable of complementing human capabilities and adapting to human work styles.

## CRediT authorship contribution statement

**Yinuo Du:** Writing – original draft, Visualization, Software, Methodology, Data curation, Conceptualization. **Baptiste Prébot:** Writing

review & editing, Methodology, Data curation, Conceptualization.
 Tyler Malloy: Writing – review & editing, Writing – original draft,
 Methodology, Formal analysis, Conceptualization. Fei Fang: Writing – review & editing, Supervision. Cleotilde Gonzalez: Writing – review & editing, Supervision, Project administration, Methodology, Funding acquisition, Conceptualization.

### Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

### Acknowledgments

This research was sponsored by the Army Research Office and accomplished under Australia–US MURIGrant Number W911NF-20-S-000 and by the Army Research Laboratory under CooperativeAgreement Number W911NF-13-2-0045 (ARL Cyber Security CRA).

## Appendix. Behavior analysis of the autonomous agents

Fig. 9 illustrates the target selection patterns of three different agent types: Random, Heuristic, and Cognitive on the 25 steps of an episode, averaged across all episodes in the experiment.

As evidenced by the consistent distribution across all steps, the Random agent maintains approximately equal proportions of actions targeting the User, Enterprise, and Operational subnets throughout an episode.

In contrast, both the Heuristic and Cognitive agents demonstrate similar dynamic behaviors regarding the action proportions on the three subnets. The Heuristic agent initially executes more actions on the Enterprise subnet, gradually shifting their attention toward the Operational subnet as the episode progresses.

Interestingly, while the Cognitive agent initially demonstrates a similar frequency of actions as those of the Heuristic agent, the Cognitive agent appears to prioritize Operational hosts slightly earlier (around Step 15) and more significantly in later steps than the Heuristic agent. Thus, both the Heuristic and Cognitive agents demonstrate dynamic allocation of their resources to protect the most vulnerable and valuable network assets as the attack progresses. The Cognitive agent demonstrates superior performance given its earlier prioritization of the actions on the Operational subnet.

Importantly, the progression of the adversary through the network topology (shown in Fig. 10) is essentially the same, regardless of the type of agents (Random, Heuristic, or Cognitive) used as a human teammate. The adversary begins in the User subnet, proceeding through the Enterprise subnet, and ultimately targets the critical Operational subnet. These two figures illustrate the competency of the Heuristic and the Cognitive agents to dynamically follow the attacker path of action to block it.

#### References

- Abele, A. E., Ellemers, N., Fiske, S. T., Koch, A., & Yzerbyt, V. (2021). Navigating the social world: Toward an integrated framework for evaluating self, individuals, and groups. *Psychological Review*, 128(2), 290.
- Aggarwal, P., Thakoor, O., Jabbari, S., Cranford, E. A., Lebiere, C., Tambe, M., & Gonzalez, C. (2022). Designing effective masking strategies for cyberdefense through human experimentation and cognitive models. *Computers & Security*, 117, Article 102671.
- Albasheer, H., Md Siraj, M., Mubarakali, A., Elsier Tayfour, O., Salih, S., Hamdan, M., Khan, S., Zainal, A., & Kamarudeen, S. (2022). Cyber-attack prediction based on network intrusion detection systems for alert correlation techniques: a survey. Sensors, 22(4), 1494.
- Anderson, J. R., Corbett, A. T., Koedinger, K. R., & Pelletier, R. (1995). Cognitive tutors: Lessons learned. The Journal of the Learning Sciences, 4(2), 167–207.

- Anderson, J. R., & Lebiere, C. J. (1998). In J. R. Anderson, & C. J. Lebiere (Eds.), The atomic components of thought (p. 504). New York: Psychology Press, http://dx.doi.org/10.4324/9781315805696.
- Anderson, J. R., & Lebiere, C. J. (2014). The atomic components of thought. Psychology Press.
- Ban, T., Samuel, N., Takahashi, T., & Inoue, D. (2021). Combat security alert fatigue with ai-assisted techniques. In Proceedings of the 14th cyber security experimentation and test workshop (pp. 9-16).
- Bansal, G., Nushi, B., Kamar, E., Lasecki, W. S., Weld, D. S., & Horvitz, E. (2019).
  Beyond accuracy: The role of mental models in human-AI team performance. 7,
  In Proceedings of the AAAI conference on human computation and crowdsourcing (pp. 2–11).
- Baruwal Chhetri, M., Tariq, S., Singh, R., Jalalvand, F., Paris, C., & Nepal, S. (2024).
  Towards Human-Al Teaming to Mitigate Alert Fatigue in Security Operations
  Centres. ACM Transactions on Internet Technology.
- Bécue, A., Praça, I., & Gama, J. (2021). Artificial intelligence, cyber-threats and industry 4.0: Challenges and opportunities. Artificial Intelligence Review, 54(5), 3849–3886.
- Bendell, R., Williams, J., Fiore, S. M., & Jentsch, F. (2024). Individual and team profiling to support theory of mind in artificial social intelligence. *Scientific Reports*, 14(1), 12635.
- Boone, R. T., & Buck, R. (2003). Emotional expressivity and trustworthiness: The role of nonverbal behavior in the evolution of cooperation. *Journal of Nonverbal Behavior*, 27, 163–182.
- Boyarchuk, R., Khudyntsev, M., Lebid, O., & Trofymchuk, O. (2021). Organizational and technical model of national cybersecurity and cyber protection. In *CPITS* (pp. 37–46).
- Bugbee, E. H., & Gonzalez, C. (2022). Making predictions without data: How an instance-based learning model predicts sequential decisions in the balloon analog risk task. In Proceedings of the annual meeting of the cognitive science society: Vol. 44, (44).
- Chiou, E. K., Lee, J. D., & Su, T. (2019). Negotiated and reciprocal exchange structures in human-agent cooperation. Computers in Human Behavior, 90, 288–297.
- Cranford, E., Gonzalez, C., Aggarwal, P., Cooney, S., Tambe, M., & Lebiere, C. (2020). Adaptive cyber deception: Cognitively informed signaling for cyber defense. In Proceedings of the 53rd hawaii international conference on system sciences.
- Cranford, E. A., Gonzalez, C., Aggarwal, P., Tambe, M., & Lebiere, C. (2020). What attackers know and what they have to lose: Framing effects on cyber-attacker decision making. 64, In Proceedings of the human factors and ergonomics society annual meeting (1), (pp. 456–460). SAGE Publications Sage CA: Los Angeles, CA.
- Cranford, E. A., Lebiere, C., Rajivan, P., Aggarwal, P., & Gonzalez, C. (2019). Modeling cognitive dynamics in (End)-user response to phishing emails.
- Cummings, M. M. (2014). Man versus machine or man+ machine? *IEEE Intelligent Systems*, 29(5), 62–69.
- de Visser, E. J., Krueger, F., McKnight, P., Scheid, S., Smith, M., Chalk, S., & Parasuraman, R. (2012). The world is not enough: Trust in cognitive agents. 56, In Proceedings of the human factors and ergonomics society annual meeting (1), (pp. 263–267). Los Angeles, CA: Sage Publications Sage CA.
- Demir, M., McNeese, N. J., Johnson, C., Gorman, J. C., Grimm, D., & Cooke, N. J. (2019). Effective team interaction for adaptive training and situation awareness in human-autonomy teaming. In 2019 IEEE conference on cognitive and computational aspects of situation management (pp. 122–126). IEEE.
- Du, Y., Song, Z., Milani, S., Gonzales, C., & Fang, F. (2022). Learning to play an adaptive cyber deception game. 6, In Proc. of the 21st international conference on autonomous agents and multiagent systems. Auckland, New Zealand.
- Dutt, V., Ahn, Y.-S., & Gonzalez, C. (2011). Cyber situation awareness: Modeling the security analyst in a cyber-attack scenario through instance-based learning. In Y. Li (Ed.), Data and applications security and privacy XXV (pp. 280–292). Berlin, Heidelberg: Springer.
- Dykstra, J., & Paul, C. L. (2018). Cyber operations stress survey ({{{{COSS)}}}}}: Studying fatigue, frustration, and cognitive workload in cybersecurity operations. In 11th USeNIX workshop on cyber security experimentation and test.
- Fang, F. (2021). Game theoretic models for cyber deception. In *Proceedings of the 8th ACM workshop on moving target defense* (pp. 23–24).
- Ferster, B. (2022). Intelligent tutoring systems. Routledge.
- Geib, C. W., & Goldman, R. P. (2009). A probabilistic plan recognition algorithm based on plan tree grammars. Artificial Intelligence, 173(11), 1101–1132.
- Glikson, E., & Woolley, A. W. (2020). Human trust in artificial intelligence: Review of empirical research. Academy of Management Annals, 14(2), 627–660.
- Gonzalez, C. (2024). Building human-like artificial agents: A general cognitive algorithm for emulating human decision-making in dynamic environments. *Perspectives on Psychological Science*, 19(5), 860–873.
- Gonzalez, C., Aggarwal, P., Lebiere, C., & Cranford, E. (2020). Design of dynamic and personalized deception: A research framework and new insights. In Proceedings of the 53rd Hawaii international conference on system sciences.
- Gonzalez, C., Ben-Asher, N., Martin, J. M., & Dutt, V. (2015). A cognitive model of dynamic cooperation with varied interdependency information. *Cognitive Science*, 39(3), 457–495.
- Gonzalez, C., & Dutt, V. (2011). Instance-based learning: integrating sampling and repeated decisions from experience. *Psychological Review*, 118(4), 523.
- Gonzalez, C., Lerch, J. F., & Lebiere, C. (2003). Instance-based learning in dynamic decision making. Cognitive Science, 27(4), 591–635.

- Gutzwiller, R. S., Clegg, B. A., & Blitch, J. G. (2013). Part-task training in the context of automation: Current and future directions. *The American Journal of Psychology*, 126(4), 417–432.
- Gutzwiller, R. S., Espinosa, S. H., Kenny, C., & Lange, D. S. (2018). A design pattern for working agreements in human-autonomy teaming. 8, In Advances in human factors in simulation and modeling: proceedings of the AHFE 2017 international conference on human factors in simulation and modeling, July 17–21, 2017, the Westin Bonaventure Hotel, Los Angeles, California, USA (pp. 12–24). Springer.
- Hackman, J. R. (1978). The design of work in the 1980s. *Organizational Dynamics*, 7(1), 3–17.
- Harris-Watson, A. M., Larson, L. E., Lauharatanahirun, N., DeChurch, L. A., & Contractor, N. S. (2023). Social perception in Human-AI teams: Warmth and competence predict receptivity to AI teammates. *Computers in Human Behavior*, 145, Article 107765.
- Hauptman, A. I., Schelble, B. G., McNeese, N. J., & Madathil, K. C. (2023). Adapt and overcome: Perceptions of adaptive autonomous agents for human-AI teaming. Computers in Human Behavior, 138, Article 107451.
- Hoff, K. A., & Bashir, M. (2015). Trust in automation: Integrating empirical evidence on factors that influence trust. Human Factors, 57(3), 407-434.
- Holroyd, A., Rich, C., Sidner, C. L., & Ponsler, B. (2011). Generating connection events for human-robot collaboration. In *2011 RO-MAN* (pp. 241–246). IEEE.
- Jiao, J. R., Zhou, F., Gebraeel, N. Z., & Duffy, V. (2020). Towards augmenting cyber-physical-human collaborative cognition for human-automation interaction in complex manufacturing and operational environments. *International Journal* of Production Research, 58(16), 5089–5111. http://dx.doi.org/10.1080/00207543. 2020.1722324.
- Johnson, D. W., & Johnson, R. T. (1989). Cooperation and competition: Theory and research. Interaction Book Company.
- Kautz, H. A., Allen, J. F., et al. (1986). Generalized plan recognition. 86, In AAAI (3237), (p. 5). Philadelphia, PA.
- Kay, T., Keller, L., & Lehmann, L. (2020). The evolution of altruism and the serial rediscovery of the role of relatedness. *Proceedings of the National Academy of Sciences*, 117(46), 28894–28898.
- Kocielnik, R., Amershi, S., & Bennett, P. N. (2019). Will you accept an imperfect ai? exploring designs for adjusting end-user expectations of ai systems. In *Proceedings* of the 2019 CHI conference on human factors in computing systems (pp. 1–14).
- Kott, A. (2023). 87, Autonomous intelligent cyber defense agent (AICA): A comprehensive guide. Springer Nature.
- Kulms, P., & Kopp, S. (2019). More human-likeness, more trust? The effect of anthropomorphism on self-reported and behavioral trust in continued and interdependent human-agent cooperation. In *Proceedings of mensch und computer 2019* (pp. 31–42).
- Lejarraga, T., Dutt, V., & Gonzalez, C. (2012). Instance-based learning: A general model of repeated binary choice. *Journal of Behavioral Decision Making*, 25(2), 143–153.
- Li, S., Zheng, P., Liu, S., Wang, Z., Wang, X. V., Zheng, L., & Wang, L. (2023). Proactive human-robot collaboration: Mutual-cognitive, predictable, and self-organising perspectives. Robotics and Computer-Integrated Manufacturing, 81, Article 102510.
- Linkov, I., Stoddard, K., Strelzoff, A., Galaitsi, S., Keisler, J., Trump, B. D., Kott, A., Bielik, P., & Tsankov, P. (2023). Toward mission-critical AI: Interpretable, actionable, and resilient AI. In 2023 15th international conference on cyber conflict: meeting reality (pp. 181–197). IEEE.
- Lyn Paul, C., Blaha, L. M., Fallon, C. K., Gonzalez, C., & Gutzwiller, R. S. (2019). Opportunities and challenges for human-machine teaming in cybersecurity operations. 63, In Proceedings of the human factors and ergonomics society annual meeting (1), (pp. 442–446). Los Angeles, CA: SAGE Publications Sage CA.
- Matsuda, N., Yarzebinski, E., Keiser, V., Raizada, R., Cohen, W. W., Stylianides, G. J., & Koedinger, K. R. (2013). Cognitive anatomy of tutor learning: Lessons learned with SimStudent. *Journal of Educational Psychology*, 105(4), 1152.
- McNeese, N. J., Demir, M., Cooke, N. J., & Myers, C. (2018). Teaming with a synthetic teammate: Insights into human-autonomy teaming. *Human Factors*, 60(2), 262–273.
- McNeese, N. J., Demir, M., Cooke, N. J., & Myers, C. (2018). Teaming with a synthetic teammate: Insights into human-autonomy teaming. *Human Factors*, 60(2), 262–273. http://dx.doi.org/10.1177/0018720817743223, PMID: 29185818.
- Merritt, S. M., Heimbaugh, H., LaChapell, J., & Lee, D. (2013). I trust it, but I don't know why: Effects of implicit attitudes toward automation on trust in an automated system. *Human Factors*, 55(3), 520–534.
- Musick, G., O'Neill, T. A., Schelble, B. G., McNeese, N. J., & Henke, J. B. (2021).
  What happens when humans believe their teammate is an AI? An investigation into humans teaming with autonomy. *Computers in Human Behavior*, 122, Article 106852.
- Nguyen, T. N., & Gonzalez, C. (2020). Cognitive machine theory of mind. In *CogSci*. Nguyen, T. N., & Gonzalez, C. (2022). Theory of mind from observation in cognitive
- models and humans. Topics in Cognitive Science, 14(4), 665–686.

  Nguyen, T. N., McDonald, C., & Gonzalez, C. (2021). Credit assignment: Challenges and
- opportunities in developing human-like AI agents: Technical Report, Carnegie Mellon University.

  Nguyen, T. N., Phan, D. N., & Gonzalez, C. (2023). Learning in cooperative multiagent
- systems using cognitive and machine models. ACM Transactions on Autonomous and Adaptive Systems, 18(4), 1–22.
- Nobles, C. (2022). Stress, burnout, and security fatigue in cybersecurity: A human factors problem. HOLISTICA-Journal of Business and Public Administration, 13(1), 49–72.

- O'Neill, T. A., Flathmann, C., McNeese, N. J., & Salas, E. (2023). Human-autonomy teaming: Need for a guiding team-based framework? Computers in Human Behavior, Article 107762.
- O'Neill, T., McNeese, N., Barron, A., & Schelble, B. (2022). Human–autonomy teaming: A review and analysis of the empirical literature. *Human Factors*, 64(5), 904–938.
- Parasuraman, R., Molloy, R., & Singh, I. L. (1993). Performance consequences of automation-induced'complacency'. The International Journal of Aviation Psychology, 3(1), 1–23.
- Parasuraman, R., Sheridan, T. B., & Wickens, C. D. (2000). A model for types and levels of human interaction with automation. *IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems and Humans*, 30(3), 286–297.
- Paul, C. L. (2014). Human-centered study of a network operations center: experience report and lessons learned. In *Proceedings of the 2014 ACM workshop on security information workers* (pp. 39–42).
- Pelau, C., Dabija, D.-C., & Ene, I. (2021). What makes an AI device human-like? The role of interaction quality, empathy and perceived psychological anthropomorphic characteristics in the acceptance of artificial intelligence in the service industry. Computers in Human Behavior, 122, Article 106855.
- Prebot, B. (2020). Représentation partagée et travail collaboratif en contexte C2: monitoring d'opérateurs en situation simulée de command and control (Ph.D. thesis), Bordeaux.
- Prebot, B., Du, Y., & Gonzalez, C. (2023). Learning about simulated adversaries from human defenders using interactive cyber-defense games. *Journal of Cybersecurity*, 9(1), tyad022.
- Ragot, M., Martin, N., & Cojean, S. (2020). Ai-generated vs. human artworks. A perception bias towards artificial intelligence? In Extended abstracts of the 2020 CHI conference on human factors in computing systems (pp. 1–10).
- Ramamoorthy, N., & Flood, P. C. (2004). Individualism/collectivism, perceived task interdependence and teamwork attitudes among Irish blue-collar employees: a test of the main and moderating effects? *Human Relations*, *57*(3), 347–366.
- Ren, M., Chen, N., & Qiu, H. (2023). Human-machine collaborative decision-making: An evolutionary roadmap based on cognitive intelligence. *International Journal of Social Robotics*, 15(7), 1101–1114.
- Ritter, F. E., Tehranchi, F., & Oury, J. D. (2019). ACT-R: A cognitive architecture for modeling cognition. Wiley Interdisciplinary Reviews: Cognitive Science, 10(3), Article e1488.
- Sarker, I. H., Janicke, H., Mohammad, N., Watters, P., & Nepal, S. (2023). AI potentiality and awareness: A position paper from the perspective of human-AI teaming in cybersecurity. arXiv preprint arXiv:2310.12162.
- Schelble, B. G., Flathmann, C., McNeese, N. J., O'Neill, T., Pak, R., & Namara, M. (2022). Investigating the effects of perceived teammate artificiality on human performance and cognition. *International Journal of Human–Computer Interaction*, 1, 16
- Schoorman, F. D., Mayer, R. C., & Davis, J. H. (2007). An integrative model of organizational trust: Past, present, and future. Academy of Management Review, 32(2), 344–354.
- Sidner, C. L., Lee, C., Kidd, C. D., Lesh, N., & Rich, C. (2005). Explorations in engagement for humans and robots. *Artificial Intelligence*, 166(1–2), 140–164.
- Sidner, C. L., Lee, C., & Lesh, N. (2003). Engagement when looking: behaviors for robots when collaborating with people. In Diabruck: proceedings of the 7th workshop on the semantic and pragmatics of dialogue (pp. 123–130). Citeseer.
- Standen, M., Lucas, M., Bowman, D., Richer, T. J., Kim, J., & Marriott, D. (2021). CybORG: A Gym for the development of autonomous cyber agents. In IJCAI-21 1st international workshop on adaptive cyber defense. arXiv.
- Stevens, T. (2020). Knowledge in the grey zone: Al and cybersecurity. *Digital War*, 1, 164–170.
- Strauch, B. (2017). Ironies of automation: Still unresolved after all these years. *IEEE Transactions on Human-Machine Systems*, 48(5), 419–433.
- Sutton, R. S., Barto, A. G., et al. (1998). 1, Reinforcement learning: An introduction. (1), Cambridge: MIT Press.
- Tayyab, U.-e.-H., Khan, F. B., Durad, M. H., Khan, A., & Lee, Y. S. (2022). A survey of the recent trends in deep learning based malware detection. *Journal of Cybersecurity* and Privacy, 2(4), 800–829.
- Tomsett, R., Preece, A., Braines, D., Cerutti, F., Chakraborty, S., Srivastava, M., Pearson, G., & Kaplan, L. (2020). Rapid trust calibration through interpretable and uncertainty-aware AI. *Patterns*, 1(4).
- Vegesna, V. V. (2023). Comprehensive analysis of AI-enhanced defense systems in cyberspace. *International Numeric Journal of Machine Learning and Robots*, 7(7).
- Von der Pütten, A. M., Krämer, N. C., Gratch, J., & Kang, S.-H. (2010). "It doesn't matter what you are!" Explaining social effects of agents and avatars. Computers in Human Behavior, 26(6), 1641–1650.
- Walliser, J. C., de Visser, E. J., & Shaw, T. H. (2023). Exploring system wide trust prevalence and mitigation strategies with multiple autonomous agents. *Computers in Human Behavior*, 143, Article 107671.

- Wohleber, R. W., Stowers, K., Barnes, M., & Chen, J. Y. (2023). Agent transparency in mixed-initiative multi-UxV control: How should intelligent agent collaborators speak their minds? Computers in Human Behavior, 148, Article 107866. http://dx. doi.org/10.1016/j.chb.2023.107866, URL https://www.sciencedirect.com/science/ article/pii/S0747563223002170.
- Wynne, K. T., & Lyons, J. B. (2018). An integrative model of autonomous agent teammate-likeness. *Theoretical Issues in Ergonomics Science*, 19(3), 353–374.
- Xu, T., Singh, K., & Rajivan, P. (2022). Modeling phishing decisions using instance based learning and natural language processing.
- Zhang, G., Chong, L., Kotovsky, K., & Cagan, J. (2023). Trust in an AI versus a human teammate: The effects of teammate identity and performance on human-AI cooperation. *Computers in Human Behavior*, 139, Article 107536.
- Zhang, R., McNeese, N. J., Freeman, G., & Musick, G. (2021). "An ideal human" expectations of AI teammates in human-AI teaming. *Proceedings of the ACM on Human-Computer Interaction*, 4(CSCW3), 1–25.

# Emergent Cooperative Decision-making in Triadic Prisoner's Dilemmas: Effects of Incentives and Information

## ARTICLE INFO

Keywords:

Prisoner's Dilemma; Triads; Cooperation; Reciprocity; Strategy; Simulation

### ABSTRACT

While pairwise cooperation has been extensively studied through the Prisoner's Dilemma (PD), our understanding of how cooperation emerges in small groups remains limited. We extend the classic dyadic PD framework to a triadic framework, examining two sets of PD games per individual and how individual strategies and relationships aggregate to group cooperation. Through two experiments (N=519), we investigate: (1) how structural incentives shape cooperation by varying the K-index (0.4 / 0.8), a theoretical value that predicts greater cooperation for higher K values, and (2) how different degrees of information about mutual interdependence affect group behavior. We find that, under minimal information conditions, a higher K-index promotes sustained cooperation in the triadic setting, in alignment with the theoretical definition of the K-index. However, while experiential information (observing others' actions/outcomes) enhances cooperation, descriptive information (complete payoff matrices) paradoxically reduces cooperation. Analysis of triadic interactions reveals that selective cooperation by a third player in the group can stabilize cooperative dyadic relationships and destabilize defective dyadic relationships. These findings provide insights for designing cooperative systems, particularly in contexts where organizations must balance information sharing benefits against strategic risks.

# 1. Introduction

Cooperation is a cornerstone of human society, enabling collective endeavors from business partnerships to global collaborations. Despite its importance, cooperation has been frequently studied in dyads, and understanding how cooperation emerges and persists in groups of more than two individuals remains a significant challenge. Although the Prisoner's Dilemma (PD) has served as a foundational model for studying cooperation, its traditional focus on two-person interactions leaves critical gaps in our understanding of group dynamics. Expanding this framework to capture real-world complexity is essential, but doing so requires carefully balancing analytical tractability with ecological validity.

Previous research has extensively studied cooperation in large groups to capture realistic social dynamics. However, as the size of the group increases, interactions between individuals become difficult to disentangle, and aggregate results often obscure individual contributions Barrett and Dannenberg (2017). Studies of the N-person Prisoner's Dilemma have revealed how group size affects cooperation rates Capraro, Jordan and Rand (2013) and strategy evolution Grujić, Gracia-Lázaro, Milinski, Semmann, Traulsen, Cuesta, Sánchez and Moreno (2014), but aggregating multiple relationships obscures individual decision-making. To address this limitation, researchers have turned to the study of pairwise interactions when individuals work in groups, which provides clarity and allows a detailed study of strategies such as tit-for-tat Taylor and Nowak (2007).

Building on these approaches, our research focuses on a Triadic PD, where each player in a group of three faces a PD game with each of the other two individuals in a group. Previous research Juvina et al. (2011) demonstrated the feasibility of studying PD in a group, but their focus on collective outcomes left open questions about pairwise relationships and individual dynamics. We developed a framework that preserves the analytical clarity of pairwise interactions while capturing essential group dynamics. Our approach examines how each participant simultaneously manages relationships with two other players within a triad context, allowing the analysis of emergent properties such as third-player effects and relationship imbalances. By examining three-person groups, we extend previous work to reveal how individual strategies and dyadic relationships contribute to group cooperation. We examine two key factors that influence cooperation in the Triadic PD: structural incentives and information availability.

Early work established the K-index Rapoport, Chammah and Orwant (1965) as a theoretical predictor of cooperation by quantifying the interdependence between players. Subsequent research validated its predictive power in dyadic settings - Moisan, ten Brincke, Murphy and Gonzalez (2018a) demonstrated a correlation with cooperation rates

ORCID(s):

in two-person games, while Hilbe, Wu, Traulsen and Nowak (2014) showed how the K-index shapes the evolution of the strategy. However, empirical validation in small groups remains scarce. Although some studies examined cooperative dynamics in multiplayer contexts, the role of structural incentives captured by the K-index remains unexplored. By experimentally varying the K-index in triadic PD games, this study provides novel evidence for how these incentives shape dyadic and group-level cooperation.

Information availability represents another crucial factor. In dyadic settings, research has shown how the completeness of information shapes cooperation strategies, from forgiveness under imperfect information to strict reciprocity with perfect information Romano, Balliet, Yamagishi and Liu (2017). Gonzalez and et al. (2015) showed how different levels of information affect learning and adaptation in two-person games, while Nax, Burton-Chellew, West and Young (2023) extended this analysis to larger groups but focused on aggregate outcomes. These effects become particularly complex in small groups, where the design of strategic information can be crucial to maintaining cooperation. However, we lack a systematic understanding of how different levels of information about mutual dependence affect cooperation as we move from dyads to larger groups. Building on these information effects in dyadic settings, our study examines how participants navigate more complex information environments in small groups. By systematically varying the visibility of interactions between group members, we reveal novel dynamics in how individuals process and integrate information about multiple relationships simultaneously. This approach illuminates previously hidden interaction effects between information structures and strategic choices, providing insight into how information availability shapes cooperative behavior when actors must manage multiple interdependent relationships.

These theoretical contributions have significant practical implications in multiple contexts. In this research, we study the particularly relevant context of cybersecurity, where organizations must balance the collective benefits of information sharing regarding cyber defense with competitive risks Tosh, Shetty, Sengupta and Bagchi (2015). Previous research shows that information sharing improves collective security Garrido, Sanner and Löhr (2016); however, many organizations do not share information about their vulnerabilities and experienced attacks because revealing such details could expose them to reputational damage, legal liabilities, or competitors' exploitation. This competitive tension makes it difficult for organizations to trust each other fully, even when collaboration would enhance overall security. Existing models addressing these dynamics often oversimplify bilateral relationships by assuming uniform trust or cooperation and struggle to account for the complexities of larger networks with multiple interdependencies.

The paper proceeds as follows. Section 2 reviews related work on cooperation in groups, the effects of incentives and information levels, and outlines our proposed Triadic PD. Section 3 presents the design of the Triadic PD in a multidefender game (MDG) in the cybersecurity domain. Section 4 presents the experimental design and methods to study the effects of incentives and information availability in the MDG. Section 5 presents our experimental results at the individual, dyad, and triad levels. Section 6 discusses the results and their implications for fostering cooperation in practical settings, particularly within the cybersecurity domain.

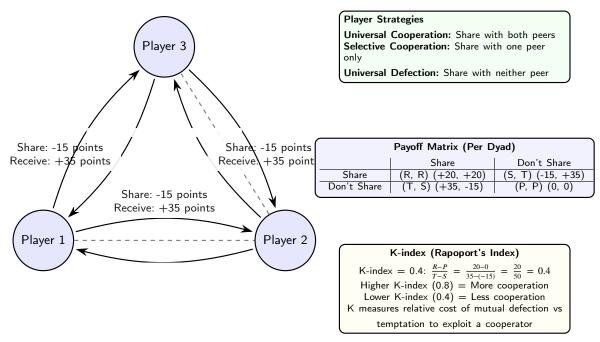
## 2. A Triadic PD Framework

Cooperation within groups has been widely studied in multiple disciplines, including economics, social psychology, and evolutionary biology. In game theory, group cooperation is often studied through public goods games and collective action problems, which introduce complexities not present in dyadic interactions. Public goods games, for example, examine how individuals contribute to a shared resource pool while dealing with incentives to ride free Fehr and Gächter (2000); Santos, Santos and Pacheco (2008). Recent research highlights factors such as group size, communication, and mechanisms for punishing free riders as critical to maintaining cooperation Perc, Jordan, Rand, Wang, Boccaletti and Szolnoki (2017); Szolnoki and Perc (2019). These studies show that group cooperation is more complex than dyadic interactions due to greater mutual dependencies and coordination challenges.

The dyadic PD has also been extended to multiplayer versions, sometimes referred to as N-person PD. In these settings, individuals must decide whether to cooperate with the entire group rather than just one other player. Research suggests that cooperation in such multiplayer PD scenarios is influenced by factors such as reciprocity, reputation, and social norms Nowak (2010); Capraro et al. (2013). However, multi-player dilemmas introduce unique challenges, including coordination issues and an increased impact of individual decisions on group outcomes Grujić et al. (2014); Van Lange, Balliet and Joireman (2020).

We extend the classic dyadic PD to triadic groups (as shown in Fig. 1), where each player in a three-person group faces a PD game with each of the other two individuals in a group. This should allow us to examine the individual, dyadic and group dynamics of cooperation and their effects of incentives and information availability. By

incorporating insights about pairwise accountability and strategic information use, this framework addresses gaps in existing game-theoretic models while providing actionable insights regarding collective, group behaviors. The triadic structure captures essential features of information sharing networks while maintaining analytical clarity.



**Figure 1:** Structure of the Triadic Prisoner's Dilemma. Each player engages in separate PD interactions with the other two players. When sharing information, a player incurs a cost of 15 points while providing a benefit of 35 points to the recipient. The K-index of 0.4 (low condition) indicates the cooperation incentive structure, calculated as the ratio between mutual cooperation benefit and exploitation temptation. In our experiment, we compare this with a higher K-index of 0.8, which theoretically promotes greater cooperation. Players must navigate these incentives while managing strategies across multiple relationships simultaneously.

In contrast to N-person PD, our Triadic PD retains the fundamental 2-by-2 PD interactions within a triad context. This allows us to focus on pairwise decision-making while examining how these interactions aggregate to influence the group. The advantage of this design is that it enables us to capture both dyadic and triadic dynamics, providing insights into how individual relationships affect broader group behavior. Unlike typical N-person PD scenarios, where cooperation is assessed at the collective level, our approach provides a detailed examination of the interplay between dyads within the group, revealing the conditions under which cooperation is stabilized or disrupted from individual behavior.

Our Triadic PD framework is particularly valuable for understanding dynamic phenomena that emerge in evolutionary game theory beyond static equilibrium states. While dyadic models can demonstrate basic cooperation patterns, they often miss complex dynamics such as oscillations between strategies that occur in multiplayer contexts. For example, the rock-paper-scissor dynamics observed in side-blotched lizards Sinervo and Lively (1996) show how frequency-dependent selection can drive cyclic changes in strategy prevalence, similar to potential oscillations in cooperation within small groups. Other relevant phenomena our framework can illuminate include contagious outbreaks of cooperation or defection Helbing and Yu (2009) and the emergence of polarization where subgroups adopt increasingly extreme strategies Yang (2023). By examining how incentive structures and information availability affect these dynamics in triads, we bridge the gap between overly simplified dyadic models and intractably complex large-group models, offering insights into how cooperation stabilizes or breaks down in real-world multi-agent systems.

## 2.1. Incentive Structure and Social Preferences

The interplay between incentive structures and social preferences significantly shapes cooperative behavior in group settings. Incentive structures, defined by the potential rewards or costs players face based on their decisions to cooperate or defect, directly influence individual motivations. Moisan et al. Moisan, ten Brincke, Murphy and Gonzalez (2018b) demonstrated that as players' cooperativeness increases, there is a sharp transition from defection to cooperation, with

the transition point depending on the game's payoff matrix. Their work showed that inequality aversion among players promotes cooperation by transforming perceived incentives.

A well-established measure of expected cooperation in PD games is Rapoport's K-index Rapoport (1967), defined as K = (R-P)/(T-S), where R represents the reward for mutual cooperation, P the punishment for mutual defection, T the temptation payoff for unilateral defection, and S the sucker payoff for unilateral cooperation (see Fig. 1). The K-index captures the expected cooperation by considering how much players benefit from defecting versus the cost of mutual defection. When K is high (i.e. when T is not much larger than S or P is numerically large), defection is less rewarding, and mutual defection is more costly, making cooperation more likely.

Prosociality Michael, McEllin and Felber (2020) (e.g., Social Value Orientation (SVO)) adds another layer to this dynamic by reflecting how individuals weigh their results against others. SVO can be represented through a utility function  $u(\pi_{\text{self}}, \pi_{\text{opponent}}) = u_{\text{self}} + \alpha \cdot u_{\text{opponent}}$ , where  $\alpha$  represents the weight given to the opponent's payoff. For any PD game, there exists a threshold  $\bar{\alpha}$  such that players with  $\alpha > \bar{\alpha}$  will prefer cooperation regardless of their beliefs about the behavior of others, while those with  $\alpha < \bar{\alpha}$  will consistently choose defection.

Furthermore, it is worth noting that recent developments in evolutionary game theory have expanded beyond the K-index to characterize social dilemmas more comprehensively. The universal dilemma strength framework distinguishes between Chicken-type dilemmas (Dg' = (T - R)/(R - P)) and Stag Hunt-type dilemmas (Dr' = (P - S)/(R - P)), offering additional dimensions to analyze the dynamics of cooperation Wang, Kokubo, Jusup and Tanimoto (2015); Ito and Tanimoto (2018). In our experimental design, we use a structure similar to the Donor & Recipient game, where these dilemma types maintain equal strength (Dg' = Dr'), allowing us to meaningfully interpret changes in cooperation through the K-index while acknowledging these broader theoretical developments.

# 2.2. Information Levels and Decision Making

The effectiveness of incentive structures in promoting cooperation is highly dependent on the information available to players about their mutual interdependence Vuolevi and Van Lange (2012). The Hierarchy of Social Information (HSI) framework proposed by Gonzalez and Martin Gonzalez and Martin (2011) conceptualizes three main levels of interpersonal information. At the Minimal Information level, players know that they are interdependent, but lack details about how their actions affect others. The Experiential Information level allows players to observe others' actions and outcomes over time, enabling learning through experience about their interdependencies. The Descriptive Information level provides complete information about the payoff structure upfront, in addition to experiential feedback.

This framework suggests that providing more detailed information about interaction structures can foster cooperation more effectively than limited or no social information. Gonzalez et al. Gonzalez, Ben-Asher, Martin and Dutt (2015) found that continued visibility of the payoff matrix helps clarify the trade-off between short-term and long-term rewards, while experiential feedback strengthens the understanding of reciprocal relationships.

The combination of incentive structures and information levels creates a complex decision environment where players must balance individual and collective interests. These factors are relevant in many practical problems. In particular, in cybersecurity contexts, organizations must decide whether to share threat information with others. Research has shown that rewarding and punishing certain actions can significantly affect information exchange behavior Tosh et al. (2015) and in cyberdefense these incentives together with the type of information exchange have shown that clarity of feedback on interdependencies influences cooperation rates Garrido et al. (2016).

# 3. Triadic PD in a Multi-Defender Game (MDG) for Cybersecurity

To test our Triadic PD framework and the effects of incentives and information, we developed the Multi-Defender Game (MDG). This game simulates a cybersecurity scenario in which three defenders must make decisions about sharing threat information. Each participant plays through 50 rounds of decision making, managing their resources while facing possible cyber attacks. Fig. 1 shows the Triadic interaction structure in our game and Fig. 2 illustrates three steps during each round of decision making in the game.

First, Fig. 2a shows the initial interface where players receive their status information and make sharing decisions. Second, when choosing to share information, players can select specific defenders or share with all members of the group, as shown in Fig. 2b. This granular control over information sharing allows players to implement selective sharing strategies.

The sharing interaction between each pair of defenders creates a PD interaction between each of the two players. Sharing information costs the sender 15 points, but provides the receiver with 35 points. Therefore, mutual sharing

results in a net gain of 35 - 15 = 20 points for each defender. When one defender shares while the other does not, the sharer loses 15 points, while the receiver gains 35 points. If neither shares, both receive 0 points.

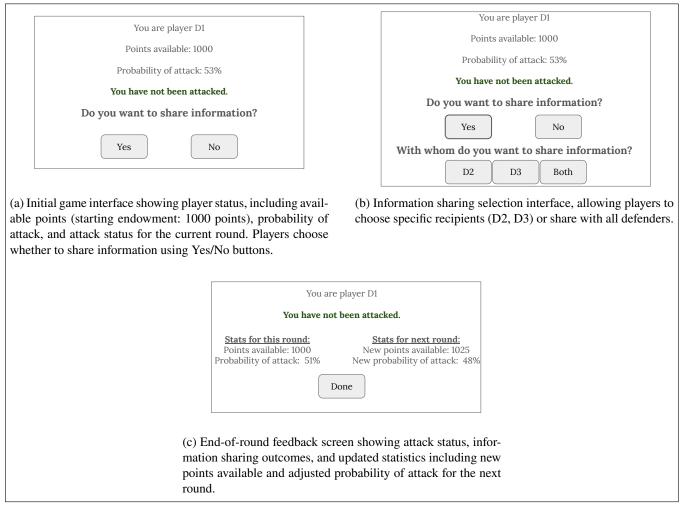


Figure 2: Interfaces used during the game: (a) Initial game interface, (b) Information sharing interface, and (c) End-of-round feedback screen.

Information sharing has both immediate and long-term effects in the game. Beyond direct point exchanges, receiving information helps strengthen a defender's security posture, reducing their probability of being breached in subsequent rounds. The probability of a breach in round t+1 is reduced according to the formula:

$$\Pr_{t+1} = \Pr_{t} -(0.95 \cdot Z_i^t / 2000) \tag{1}$$

where  $Z_i^t$  represents the accumulated reward of defender i in trial t.

The accumulated reward of each player is updated each round based on three factors: their sharing decisions, attack costs (-30 points when attacked), and information-sharing rewards:

$$R_i^t = R_i^{t-1} + Z_i^t + (-30) \cdot C_a^t \tag{2}$$

where  $C_a^t$  represents the attack status in round t.

## 4. Methods

This experiment examines how incentive structures and information availability shape cooperative behavior in Triadic PD. The MDG frames Triadic PD in the context of cybersecurity, as it is an important area to illustrate this tension between individual and collective interests. Our multilevel analysis approach allows us to examine how cooperation emerges and evolves at the individual, dyadic, and group levels.

To investigate how incentive structures and information levels influence cooperative behavior, we conducted two studies using the MDG. Study 1 examined the effect of incentive structures by varying the K-index (0.4 vs 0.8) under minimal information conditions. Study 2 investigated the impact of information availability by comparing three levels of feedback (Minimal, Experiential, Descriptive) while maintaining a constant K-index of 0.4.

The three different levels of information conditions are shown in Fig. 3. In the Minimal Condition, the participants received only basic feedback on whether other players shared information with them. The Experiential condition provided detailed information through a feedback table showing each defender's actions, attack status, and the resulting point exchanges. This allowed participants to learn about the consequences of their decisions and others' behaviors through direct experience. The Descriptive condition supplemented the experiential information with the complete pay-off matrix, allowing participants to understand the full range of possible outcomes before making their decisions.

### A. Minimal Information Level

Game Status Update:

- Defender 1 shared information with me
- Defender 2 didn't share any information

# B. Experiential Information Level

My Actions	Defender 1	Defender 2
<ul> <li>Did not share with Defender 1</li> <li>Shared with Defender 2</li> </ul>	Shared information     Was attacked	<ul><li>Did not share</li><li>No points exchanged</li></ul>
0.14.54 2 6.6.146. 2	• My gain: +35	• My gain: 0
	• Their cost: -15	• Their gain: 0

## C. Descriptive Information Level

Additional to experiential information, players see the payoff matrix: (left: K=0.4; right: K=0.8)

	Share	Don't Share		Share	Don't Share
Share	(+20, +20)	(-15, +35)	Share	(+30, +30)	(-15, +35)
Don't Share	(+35, -15)	(0, 0)	Don't Share	(+35, -15)	(-10, -10)

Figure 3: Three levels of information provided to players. (A) Minimal Information provides only basic sharing status. (B) Experiential Information shows detailed outcomes of interactions with each defender. (C) Descriptive Information adds the complete payoff matrix to help players understand potential outcomes.

This design resulted in four experimental conditions:

- Game I: K-index = 0.4, Minimal Information
- Game II: K-index = 0.8, Minimal Information
- Game III: K-index = 0.4, Experiential Information
- Game IV: K-index = 0.4, Descriptive Information

Our hypotheses examine how the relative costs and rewards of cooperation influence group behavior and how different levels of information about interdependence affect group cooperation.

**Hypothesis 1.** Groups in the high K index condition (K = 0.8) will demonstrate higher cooperation rates than those in the low K index condition (K = 0.4), as the increased reward-to-cost ratio makes cooperation more attractive.

**Hypothesis 2.** A higher K-index will lead to stronger "lock-in" effects between dyads, where the pairs maintain consistent cooperation or defection. This is because the reduced temptation to defect makes established cooperative relationships more stable, while coordination challenges and the stabilization of defection as a safe strategy reinforce existing defective relationships.

Building on the Hierarchy of Social Information framework, we examine how different levels of information about interdependence affect group cooperation.

**Hypothesis 3.** Cooperation rates will increase with information level:

- Descriptive (complete payoff information) will show the highest cooperation
- Experiential (observing others' actions/outcomes) will show moderate cooperation
- Minimal (basic awareness) will show the lowest cooperation

This progression reflects how greater awareness of mutual interdependence promotes cooperative behavior.

**Hypothesis 4.** In triadic interactions, increased information will lead to:

- Greater disparity in how individuals treat their two peers, as participants can make more informed choices about selective cooperation
- Stronger mediation effects from third players, due to a better understanding of group dynamics
- a more balanced triad, reducing disparities in relationship strength as participants can better calibrate their cooperative behaviors.

# 4.1. Participants

A total of 519 participants (173 groups of 3 individuals) from Amazon Mechanical Turk participated in groups of 3. About 32% reported having high school education, 53% a bachelor's degree, 11% a master's degree, and 4% reported other forms of education. Of the 173 groups, 51 groups participated in Game I (K = 0.4, Info = minimum), 51 groups in Game II (K = 0.8, Info = minimum), 34 groups in Game III (K = 0.4, Info = experiential), and 36 groups in Game IV (K = 0.4, Info = descriptive). About 46% of the participants identified as female. Participants received a base payment of \$3 and could earn up to \$1.75 additional bonuses based on their performance. The average time to complete the experiment was 20 minutes (SD = 2.3).

This study was approved by the Carnegie Mellon University Institutional Review Board (IRB ID: IRB-STUDY2015\_00000418, "Social Cognitive Aspects of Decision Making in Cyber-Security"). All participants provided informed consent before participating in the experiment.

## 4.2. Procedure

Participants provided informed consent and basic demographic information before receiving task instructions. Their main objective was to maximize their points throughout the game. Before starting, participants completed a comprehension test to test their understanding of the dynamics of the game. They received feedback on incorrect answers and proceeded only after selecting the correct responses. They were informed that they would be part of a three-player group and would receive feedback on information sharing in each trial. The experiment consisted of 50 trials, although the participants were not informed of this number in advance. After completion of all trials, the participants completed a survey about their strategies.

## 4.3. Dependent Variables

To study the effects of incentives and information on emergent cooperative behavior in 3-person groups, we perform multilevel analysis. We start from the individual level, moving on to pair-level analyses, and finally consider the triad-level analyses to examine the role of a third player, and reveal how players affect each other during repeated interactions.

## 4.3.1. Individual metrics

The choice between cooperation and defection is difficult, as its costs and consequences are not immediately clear. Rapid decision making often contradicts the principles of a well-considered policy. The behavior of the participants can be spontaneous and largely influenced by their basic attitudes. Therefore, it is important to understand whether participants are naturally inclined to cooperate or to defect. We measure the behavior of the participants that might lead to good individual performance (*Success*), their cooperation frequencies contingent on the preceding play, including the other player's response (*Decision-conditioned Probabilities*), the payoff (*Outcome-Conditioned Probabilities*), and the attack status (*Post-Attack Cooperation Probability*).

Past research has shown that for an individual to perform well in a durable iterated Prisoner's Dilemma, the rule of thumb is to avoid unnecessary conflict by initiating cooperative behavior. Additionally, it is important to present a predictable pattern and make it clear to the other player that both cooperation and defect will be reciprocated, encouraging mutual responsiveness. We measure the participant's (1) *First move*, which refers to their initial decision in the game: sharing with both group mates; sharing with one group mate; and sharing with none. We quantify the (2) *Predictability* of a participant's behavior (X) given the behavior of the other player in the last round (Y) with conditional entropy:

$$H(X \mid Y) = P(Y = C) \cdot \left( -\sum_{x \in \{C, D\}} P(X = x \mid Y = C) \cdot \log_2 P(X = x \mid Y = C) \right) + P(Y = D) \cdot \left( -\sum_{x \in \{C, D\}} P(X = x \mid Y = D) \cdot \log_2 P(X = x \mid Y = D) \right)$$
(3)

This metric captures how consistently a participant's actions can be anticipated based on the actions of the other player. Specifically, it measures the uncertainty in the response of a participant (cooperate or defect) given the previous move of the other player (cooperate or defect), providing a sense of how predictable and stable their strategy is.

Post-Attack Cooperation Probability. In MDG, participants are faced with the possibility of being attacked at each trial. To understand the effect of an attack on each individual's behavior, we measured the probability of cooperating after a play in which the participant was attacked (P(C|Attacked)).

Decision-conditioned probabilities. To gain insight into how individuals adjust their strategies to share or not information with other defenders based on their behavior, we measured the following conditional probabilities for each individual in a pair of defenders: (1) Cooperation Inertia: the probability that Player 1 responds cooperatively following their own cooperative response on the preceding play, regardless of Player 2's behavior. (2) Defection Inertia: the probability that Player 1 responds defectively following his own defecting response on the preceding play, regardless of Player 2's behavior.

Outcome-Conditioned probabilities. We also evaluate how a participant's behavior changes depending on the past reward: (1) Trustworthiness: This refers to the probability of cooperating following a play in which the participant received the outcome R. (2) Forgiveness: This denotes the probability of cooperating following a play in which the player received the outcome S. (3) Repentance: This describes the probability of cooperating after a play in which the player received the result T. (4) Trust: This is the probability of cooperating after a play in which the player received P out.

## 4.3.2. Dyad metrics.

To find out the strength of interaction, we evaluate the correlation between the frequency of cooperative choices of two dyad members (*Imitation*) and patterns presented in the sequence of plays (Player 1 Cooperate, Player 2 Cooperate - CC; Player 1 Cooperate, Player 2 Defect, Player 1 Defect, Player 2 Cooperate - DC; Player 1 Defect, Player 2 Defect, DD) generated by the dyads (*Lock-in*).

Imitation. If participants tend to imitate each other, we expect to see matched responses such as CC (both cooperate) and DD (both defects) in each test. To examine whether the imitation effect operates in the sequence of individual plays, we calculate the *proportion of matched responses*  $(\rho_x)$ . Here, we define the decisions of player 1 and player 2 as random variables, where A represents the decision of player 1 and B represents the decision of player 2 lagged by

x plays. The formula for  $\rho_x$  is:

$$\rho_{X} = \frac{I_{CC}I_{DD} - I_{CD}I_{DC}}{\sqrt{(I_{CC} + I_{CD})(I_{CC} + I_{DC})(I_{DD} + I_{CD})(I_{DD} + I_{DC})}} \tag{4}$$

Where  $I_{CC}$ ,  $I_{CD}$ ,  $I_{DC}$ , and  $I_{DD}$  represent the frequencies of each possible outcome pair in the sequence of plays:  $I_{CC}$  is the frequency of both players cooperating,  $I_{CD}$  is the frequency of player 1 cooperating while player 2 defects,  $I_{DC}$  is the frequency of player 1 defecting while player 2 cooperates, and  $I_{DD}$  is the frequency of both players defecting. This formula computes the correlation coefficient between the players' decisions, normalized to account for the different frequencies of cooperation and defection by each player. In this context,  $\rho_1$  represents the reaction to what the other player did on the immediately preceding play, while  $\rho_2$ ,  $\rho_3$ , and  $\rho_4$  represent the degree of interaction with the other player's responses from 2, 3, and 4 plays ago, respectively. The window of x plays is important because it allows us to understand how far back the influence of one player's decision extends in affecting the other player's behavior.

Lock-in Effect. To inspect the patterns of interaction between dyads, we group the sequence of plays into blocks of 25 consecutive trials (a block) as a unit of analysis and characterize each unit by (1)  $F_{CC}$ ,  $F_{DD}$ : the fraction of times the CC, DD state occurs within the 25 trials; and (2)  $L_{CC}$ ,  $L_{DD}$ : whether the dyad is predominantly in a particular state toward the end of the block.  $L_{CC} = 1$  if there are 10 or more CC from the last 13 plays in the block. The choice of 25 plays as the size of the block is to balance the need for sufficient data to observe patterns while keeping the analysis manageable. Focusing on the last 13 plays within the block helps in understanding recent behavior and whether the dyad has settled into a stable pattern.

We expect to see *the imitation effect* in all games. We expect to observe stronger *Lock-in Effect* in Game II where the temptation to defect is relatively weak. We expect to see a weaker *lock-in effect* in Game IV where participants are repeatedly presented with the payoff matrix, which potentially leads participants to break the pattern. The *Experiential* information in Game III might reinforce the outcomes of actions and strengthen the *Lock-in Effect*, or prompt participants to behave more strategically thus weaken the *Lock-in Effect*.

### 4.3.3. Group metrics

The behavior in triads can differ significantly from that in dyads. For example, a third player can act as a mediator and influence the decisions of the other two. Faced with two other players, participants might adopt a dual strategy, treating their interactions with each participant separately based on previous interactions. Participants may also consider fairness and equality in their decisions, especially when the actions and consequences of the two peers are displayed side by side. Alternatively, they might intentionally favor one peer due to personal preferences or previous interactions. To study emergent group behaviors in triads, we evaluated how the presence of a third participant affects cooperative behavior towards their two peers (*Disparity* and *Sequential Dependence*), the strength of dyads (*Mediation Effect*), and the friendship dynamics in the group ( *balance*). Each of these metrics is explained below.

We expect a third player to enforce cooperation between the other two players, especially when the incentive structure encourages cooperation, as in Game II. In Games III and IV, where detailed information about actions and payoffs is provided, we expect participants to adopt more differential attitudes toward their peers. This detailed information might lead participants to adjust their strategies based on the perceived level of cooperation of their peers.

Disparity. Disparity occurs when one player treats the other two players differently, leading to imbalances in cooperation, trust, or reciprocity. A participant may favor one player over another if they believe that one player is more likely to reciprocate cooperation than the other. We verify the existence of disparity using two metrics: (1)  $Gap_{AB,AC}$ , which calculates the difference between player A's propensity to cooperate with each of the other two players B and C, indicating a greater disparity with a larger gap; and (2) Sequential Dependence, which analyzes player decisions based on actions of each of the other two players measuring specific dyad metrics. We identify patterns in the decisions of the participants by calculating the proportion of four possible outcomes: the player cooperates with both peers, cooperates with one peer and defects with the other, defects with one peer and cooperates with the other, and defects with both peers, conditioned on the decisions of the peers in the preceding round. These metrics help us understand the extent of the disparity in participant behavior by showing how differently a player treats their two peers based on past interactions and the influence of the actions of each peer.

**Table 1**Regression Table

Predictor	Odds Ratio	CI	Р
(Intercept)	1.154	1.112 ~ 1.198	< 0.001
Point	1.000047	1.000034 ~ 1.000053	< 0.001
Probability of breach	0.953	0.914 ~ 0.993	< 0.05
Attack status	1.185	1.160 ~ 1.211	< 0.001
Last.received.peer1	1.807	1.778 ~ 1.837	< 0.001
Last.received.peer2	1.877	1.847 ~ 1.907	< 0.001

Mediation Effect. To evaluate how a third participant influences cooperation within a dyad over time, we calculate (1)  $\rho_{third-dyad}$ , which measures the impact of the third participant's type (prosocial, neutral, or self-interested) on the level of cooperation within the dyad. This impact is analyzed using a mixed-effects model that accounts for the repeated measures over time. We define the dyads and third party in groups as follows: dyad 1 consists of player 1 and player 2 with player 3 as the third party; dyad 2 consists of player 1 and player 3 with player 2 as the third party; dyad 3 consists of player 2 and player 3 with player 1 as the third party.

$$C_{dyad,t} = \beta_0 + \beta_1 \cdot \text{Category}_{third} + \beta_2 \cdot t + \beta_3 \cdot (\text{Category}_{third} \times t) + u_i + \epsilon_t$$
 (5)

Where  $C_{dyad,t}$  represents the cooperation level within a dyad at time t;  $\beta_0$  is the intercept (baseline cooperation level);  $\beta_1$  represents the fixed effect of the third player's category on cooperation;  $\beta_2$  captures the fixed effect of time;  $\beta_3$  represents the interaction effect between the third player's category and time;  $u_i$  is the random effect for group i that accounts for between-group variation; and  $\epsilon_t$  is the residual error term at time t.

Balance. When there is disparity or favoritism in a triad, it often leads to an imbalanced state. Such an imbalance can create tension and drive changes in behavior. If two players have strong ties, the third player could either form a cooperative relationship with at least one of them, creating a closed triad with strong dyads all around, or remain excluded, leaving the triad unbalanced with strong and weak dyads coexisting. To examine the interplay of dyad relationships in triads, we take advantage of the concept of "tie strength" from Granovetter (1973) and operationalize it as (1) Dynamic strength of dyads, defined as  $DS_t = \frac{C_t + R_t}{2}$ . Here,  $C_t$  is the cooperation rate up to time t, calculated as  $C_t = \frac{\text{Cumulative cooperative rounds up to t}}{t}$ , and  $R_t$  is the reciprocity index up to time t, calculated as  $R_t = \frac{\text{Number of reciprocated cooperative actions up to t}}{\text{Total opportunities for reciprocation up to t}}$ .

# 5. Results

### 5.1. Individual-Level Behavior

Regression Analysis. A logistic regression model with the sharing decision as the dependent variable and the six predictors as independent variables (Cumulative points *Point*, Probability of Breach *Pb*, Decisions from peer 1 (Last.received.peer1) and peer 2 (Last.received.peer2) in the preceding play, and MTurkId) helped determine the effects of the various factors on the likelihood of cooperation. The regression model includes the MTurk id as an error term to indicate the error margin in the model:

shared  $\sim Point + Pb + Atk + last.received.peer1 + last.received.peer2 + (1|MturkId).$ 

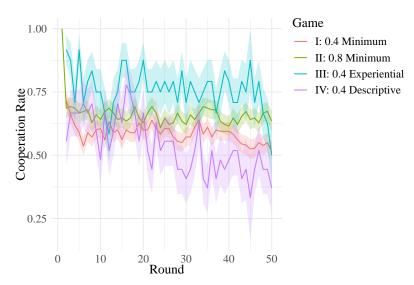
The regression model suggests that points, probability of breach, attack status, and peers' decisions in the preceding play are predictive features of cooperation, as shown in Table 1. The odds ratio suggests that, overall, receiving information from peers in the preceding game increases the chance of cooperation around  $80\% \sim 87\%$ . The participants being attacked cooperate approximately 18.5% more than the safe participants. This could be explained by realistic conflict theory Jackson (1993), when individuals or groups perceive a shared threat, they are more likely to put aside internal conflicts and unite against the common enemy. Interestingly, the probability of being attacked (Pb) has a mildly negative effect on cooperation, which reveals a discrepancy between decisions based on direct experience of risks versus abstract descriptions of probabilities Rakow and Newell (2010). Points have a positive effect on cooperation. For every extra 100 points, the odds that the outcome occurs increase by approximately 4.7%.

Game	Mean	SD
I (0.4 Minimum)	0.572	0.374
II (0.8 Minimum)	0.681	0.345
III (0.4 Experiential)	0.717	0.350
IV (0.4 Descriptive)	0.496	0.381

**Table 2**Mean and standard deviation of cooperation rate across all 50 rounds for each experimental condition (n=51, 51, 34, and 36 groups for Games I-IV, respectively).

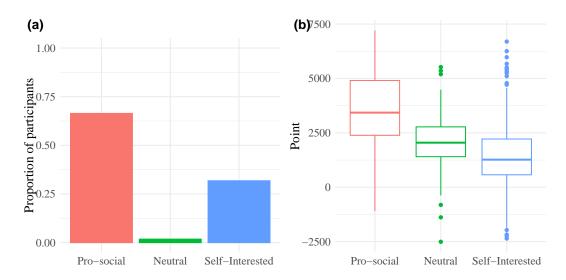
Overtime Cooperation Rate. As shown in Table 2, III > II > IV in terms of individual cooperation rate, suggesting that both a higher K-index (IV) and more information (II) can promote cooperation as expected. The effect of information is stronger. However, in contrary to our expectation, descriptive information (i.e., the display of the payoff matrix in Game IV) backfired and the cooperation frequencies are the lowest among the four games.

As shown in Fig. 4, the cooperation rate dropped sharply during the first ten rounds in all four games, similar to previous studies of PD. It then continued to decline notably in Game IV and remained relatively stable in the rest of the game conditions. A two-way ANOVA was performed to examine the effects of the number of K-indices and the information level and the round index on cooperation. There was a significant main effect of the number of K-index on cooperation [F(1, 19) = 67.656, p < 0.001] and a significant main effect of the information level [F (2, 30) = 87.179, p < 0.001]. Furthermore, the round also had a nearly significant effect on cooperation [F(49, 14) = 1.255, p < 0.109]. The interaction effects are not significant.



**Figure 4:** Average Individual Cooperation Rate Over Time. The x-axis represents the game rounds (50 total), while the y-axis shows the average cooperation rate for individuals under each game condition: Game I (K-index=0.4, Minimal information, n=51 groups), Game II (K-index=0.8, Minimal information, n=51 groups), Game III (K-index=0.4, Experiential information, n=34 groups), and Game IV (K-index=0.4, Descriptive information, n=36 groups).

The high rate of cooperation at the beginning of the game indicates that most of the participants entered the experiment with a cooperative mind. As shown in Fig. 5(a), a further inspection of their first move revealed that most of the participants (66%) are prosocial and share with both groupmates. The rest are Self-Interested (33%) and share with no one. Very few participants start with selective sharing. This finding is consistent with other studies Ackermann, Fleiß and Murphy (2016) that explicitly measure the SVO of participants. A bimodal pattern is commonly found in the SVO index distributions, where most participants are quite self-regarding or rather prosocial. The nice gesture paid off in this game. As shown in Fig. 5 (b), nicer participants perform significantly better than Self-Interested participants, under all conditions (Pro-social: m = 3495, Neutral: m = 2098, Self-Interested: m = 1505).



**Figure 5**: Participant Strategies and Performance. (a) The proportion of participants categorized by their initial sharing strategy: Pro-social (shared with both groupmates), Neutral (shared with one groupmate), or Self-Interested (shared with no one). (b) Cumulative points earned by participants at the end of the game, categorized by their initial strategy.

Conditional Propensities. To examine the cooperativeness of participants in more detail, we measure the components of the cooperative rate, the conditional probabilities. As shown in Fig. 6(a), the overall tendency to retaliate and to persist in defection is stronger than the tendency to respond cooperatively and to persist in cooperation. As shown in Table 3, the only exception is Game III (Info = Experiential, K = 0.4), where the tendency to respond cooperatively and persist in cooperation triumphs over the tendency to retaliate and persist in defection. One-way ANOVA showed that the level of information on mutual interdependence has a significant or almost significant effect on persistence [F(2, 347) = 3.953, p < 0.05], reciprocate [F(2,347) = 2.812, p = 0.061], and continue to tend [F(2,347) = 2.449, p = 0.087]. The K-index of the incentive structure has a significant effect on revenge [F(1, 347) = 4.086, p < 0.05], reciprocates [F(1, 347) = p < 0.05] and continues the trend [F(1, 347) = p < 0.01].

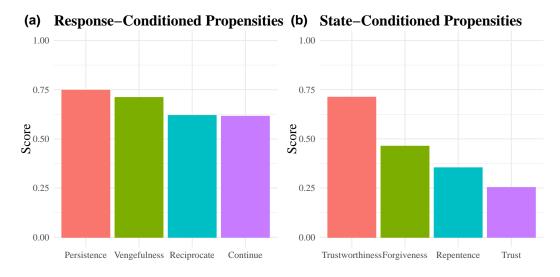


Figure 6: Behavioral Propensities Based on Past Decisions. (a) Cooperation rates conditioned on an individual's or their peer's prior actions (e.g., persistence and vengefulness tendencies). (b) Cooperation rates conditioned on the outcomes of prior interactions (e.g., trustworthiness and forgiveness tendencies).

In terms of state-conditioned propensities, as shown in Fig. 6(b),  $Trustworthiness\ x > Forgiveness\ y > Repentance\ z > Trust\ w$  in all four games. These four conditional probabilities correspond to a memory-1 strategy framework as

	Persistence	Vengefulness	Reciprocate	Continue
ı	0.743	0.716	0.685	0.655
Ш	0.730	0.676	0.658	0.667
Ш	0.647	0.681	0.721	0.708
IV	0.735	0.697	0.530	0.560

 Table 3

 Average cooperation rate conditioned on the actions in preceding play.

	Trustworthiness (x)	Forgiveness (y)	Repentence (z)	Trust (w)	Attacked	Safe
I	0.742	0.466	0.418	0.324	0.556	0.559
П	0.733	0.550	0.434	0.266	0.672	0.684
III	0.650	0.434	0.303	0.278	0.765	0.783
IV	0.650	0.434	0.303	0.278	0.455	0.479

 Table 4

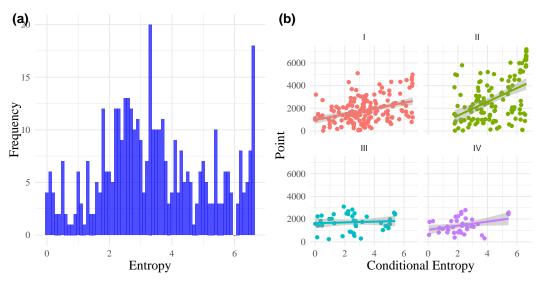
 Average cooperation rate conditioned on the reward and attack status in preceding play.

formalized by Hilbe, Martinez-Vaquero, Chatterjee and Nowak (2017), who demonstrate that strategies based only on the previous round's outcome can be highly effective in repeated social dilemmas. In their framework, a player's strategy is fully described by four parameters representing the probability of cooperation following each possible outcome state. Since x, y, 1-z, 1-w represent the tendencies to repeat the previous response in each of the four states, respectively, CC, CD, DC, DD, taking the perspective of Player 1, as T > R > P > S, we expect 1 - z > x > 1 - w > y. This theoretical expectation aligns with the adaptive dynamics model of LaPorte, Hilbe and Nowak (2023), who showed that memory-1 strategies evolve differently depending on the payoff structure and population composition. The violation of (a) 1-z>x indicates a greater propensity to cooperate than one would expect from the payoffs. This cooperative bias has been identified by Hilbe et al. (2017) as a characteristic of successful memory-1 strategies in environments that favor long-term reciprocity over immediate exploitation. The violation of (b) x > 1 - w indicates a greater propensity to defect than one would expect from the rewards. The violation of (c) 1-z>1-w is ambivalent since both are propensities to defect. As shown in Table 4, the violation of the cooperative bias of (a) is present in all except the Info descriptive condition with the payoff matrix presented. The violation of the defecting bias of (b) is present in conditions without information about the mutual interdependence between participants. (c) is present in all conditions, indicating that fear of receiving S rather than the hope of receiving T is the most important factor in the persistent response to defection. Interestingly, the different patterns observed across our experimental conditions align with the findings of LaPorte et al. (2023), who found that different information environments can alter the adaptive value of specific memory-1 strategies, particularly affecting how players weight immediate versus future payoffs. It is also worth noting that participants are biased towards cooperation in all games except Game IV (Info = descriptive, K = 0.4). Finally, unlike our expectation, being attacked or not in the preceding round has no significant effect on the cooperation rate, as shown in Table 4.

In light of these findings, the next question naturally arises: Are these conditional cooperative behaviors noticeable to the other party in the game as discernible patterns? To investigate this, we calculated the conditional entropy to measure the randomness of the behavior of the participants. As shown in Fig. 7(a), the participants behave randomly, that is, when the opponent chooses to cooperate, the participants might cooperate or defect. However, unlike our expectation, the randomness of the behavior of the participants does not harm their performance. In contrast, as demonstrated in Fig. 7(b) higher randomness leads to better performance.

## 5.2. Dyad-Level Behavior

Imitation Effect As shown in Table 5, correlation  $(\rho)$  between the random variables C1 (player 1 share or not share) and C2 (player 2 share or not share) showed that there is a significant correlation between the decision to share of the participants and the decisions of their opponents in previous trials. The strength of the interaction decays with the interval. The correlation  $(\rho)$  is the smallest in Game IV (k = 0.4, Info = descriptive). Evidently, presenting the payoff matrix obscures the effect of imitating the other's last response.



**Figure 7:** Relationship Between Behavioral Randomness and Performance. (a) Distribution of participants' conditional entropy values, representing the unpredictability of their responses based on their peers' prior actions. (b) Relationship between conditional entropy and cumulative points earned, showing how greater randomness in decision-making influences performance across game conditions.

	$\rho_1$	$ ho_2$	$\rho_3$	$ ho_4$	$\rho_5$
	0.274	0.221	0.177	0.162	0.163
Ш	0.287	0.231	0.209	0.168	0.155
Ш	0.299	0.217	0.143	0.120	0.067
IV	0.156	0.110	0.108	0.073	0.076

**Table 5**Correlation between decisions of two paired participants.

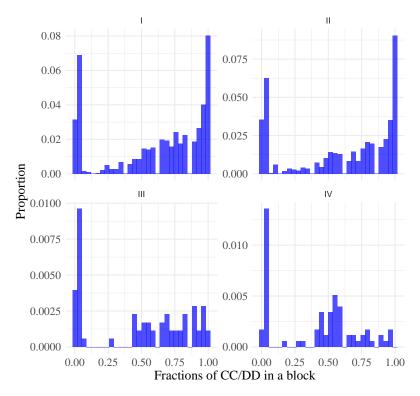
Lock-in Effect As shown in Fig. 8, a bimodal distribution is observed with respect to the fraction of CC and DD responses. It shows that the interaction between two participants tends to throw the performance toward one extreme, lasting mutual cooperation, or mutual defect. Pairs of participants in games with higher K-index payoffs are more prone to be locked in CC or DD traps. Pairs of participants in games with more information about their mutual interdependence are less likely to be locked in CC or DD traps.

## 5.3. Triad-Level Behavior

Unlike dyadic interactions, players in triads must manage multiple relationships simultaneously. We first examine how individuals distribute their cooperation between two peers. Then analyze how they develop differentiated strategies based on each peer's behavior. Building on these individual patterns, we investigate how the strategy of a third player influences cooperation between pairs. Finally, we examine how these pairwise interactions combine to shape the overall balance of relationships within the group.

Disparity We first examine whether participants tend to share information equally with both peers or favor one over the other. Under high incentives (k-index = 0.8), participants shared with both peers 43% of the time and with only one peer 26% of the time. Under low incentives (k-index = 0.4), the sharing was split more evenly: 29% with both peers and 27% with one peer.

We performed repeated measures ANOVA and found a statistically significant effect of condition (F(3, 515) = 2.148, p < 0.001) and round (F(49, 25235)= 5.148, p < 0.001), while there was no significant effect of the interaction between condition and round (F(147, 25235) = 2.849, p=1). A Tukey post hoc test suggested that participants increase sharing with one defender over rounds and decrease sharing with two defenders under all conditions. Fig. 9 shows that the participants started sharing with both players; however, the proportion of sharing with both reduced, while sharing with one increased overtime (specifically in the first 25 trials).

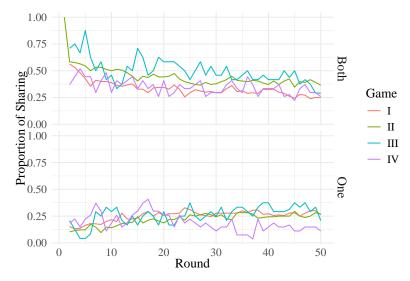


**Figure 8:** Horizontal: Fraction of CC/DD responses in a block of 25 plays; Vertical: Fraction of 25 blocks corresponding to each fraction of DD responses

Sequential Dependence Building on the observed shift towards selective sharing, we analyzed how the treatment of each peer evolved for the participants based on previous interactions. Fig. 10 shows the complete dynamics of participants' responses to their peers' previous actions. After mutual cooperation (CC), the participants maintained high levels of cooperation, with this tendency strengthening over time under all conditions. After mutual defection (DD), participants in Game III showed a greater willingness to initiate cooperation (DD  $\rightarrow$  CC, CD, or DC), although this decreased with time. When peers made different choices (CD or DC), responses varied by information condition: In Game III with experiential feedback, participants often attempted to restore mutual cooperation (CD/DC  $\rightarrow$  CC), while in Game IV with complete payoff information, they more frequently matched their peers' previous actions (CD  $\rightarrow$  CD, DC  $\rightarrow$  DC).

Figure 11 examines the behavioral differences between the treatment of the two peers by participants in several metrics: persistence, vengefulness, forgiveness, and trust. The "gap" metrics reveal fluctuations in how participants differentiated their cooperative or retaliatory responses toward peers, with some variations observed across games and rounds. For example, persistence and vengefulness metrics exhibit shifts that may reflect adjustments in participants' approaches to maintaining or retaliating against cooperation, while forgiveness and trust show varying tendencies to rebuild cooperation after negative interactions. These differences appear to be more pronounced in Games III and IV, where higher levels of information could have encouraged more nuanced strategies. Although the results do not point to a single clear pattern, they highlight the dynamic nature of triadic interactions and the potential for information availability to influence relational strategies.

Overall, the individual's behavior is not uniform across peers, with significant disparities in how they respond to cooperation and defection. As expected, the disparity is more significant in games with higher levels of information (Game III and IV), especially in terms of *Persistence, Vengefulness, Forgiveness*, and *Trust*. In Game III, the participants become more *Persistent* in continuing their defection toward one peer than to another. Overtime, they view a peer as more *trustworthy* and become more willing to take a risk with a peer after mutual defection in the preceding round (*Trust*).



**Figure 9:** Overtime sharing preferences of individual participants across the four experimental conditions. The y-axis represents the proportion of participants who chose to share information with both groupmates (top) or with only one groupmate (bottom) over 50 rounds. Game I (K-index=0.4, Minimal information, n=51 groups), Game II (K-index=0.8, Minimal information, n=51 groups), Game III (K-index=0.4, Experiential information, n=34 groups), and Game IV (K-index=0.4, Descriptive information, n=36 groups).

Mediation Effect Moving from individual behaviors to group dynamics, we examine how a third player's cooperation strategy affects the relationship between the other two players. We categorize third-player strategies as universal cooperation ("All"), universal defection ("None"), or selective cooperation ("One"). Fig. 12 shows how these strategies influence pair cooperation over time under different initial conditions.

Surprisingly, selective cooperation by the third player leads to the highest levels of pair cooperation when pairs start with mutual cooperation (CC) or mutual defection (DD). However, this selective strategy can destabilize pairs that begin with mixed strategies (CD/DC). In contrast, universal strategies (either "All" or "None") lead to more stable but generally lower levels of cooperation, particularly when pairs start with unilateral cooperation.

These patterns suggest that selective cooperation by the third player can effectively promote cooperation in stable pairs but may disrupt already unstable relationships. This finding highlights how third-party behavior can reinforce or destabilize the dynamics of existing relationships.

*Balance* To understand how pairwise relationships interact within triads, we examine the relative strength of all three dyadic relationships over time. Fig. 13 shows these dynamics under different initial conditions.

When all members begin with cooperation ("All, All, All"), relationships stabilize with minimal differences in strength between pairs. Mixed initial strategies lead to growing disparities between pairs, creating persistent imbalances. Even when one member cooperates while two defect ("All, None, None"), overall cooperation increases, but with noticeable gaps between pair relationships, the cooperative member tends to form stronger bonds with whichever peer reciprocates first.

Groups that start with universal defection can improve over time, although they plateau with persistent differences between pairs. These findings align with social balance theory: Although consistent cooperation promotes stable, balanced relationships, asymmetric strategies create lasting imbalances that resist equilibration.

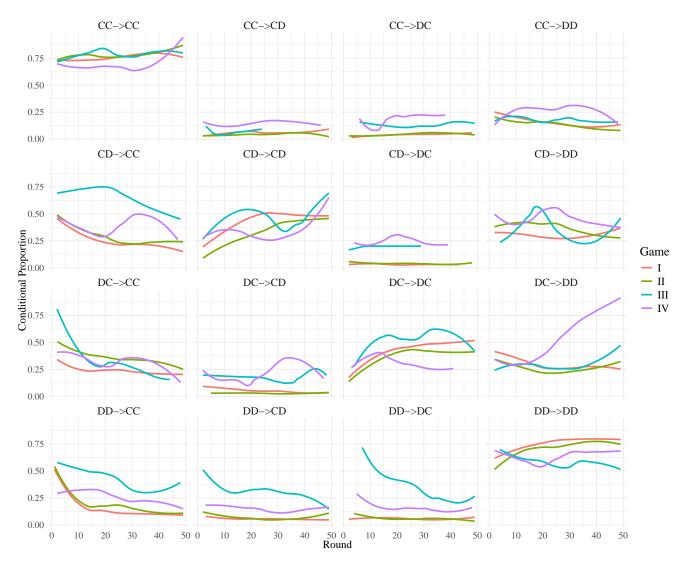
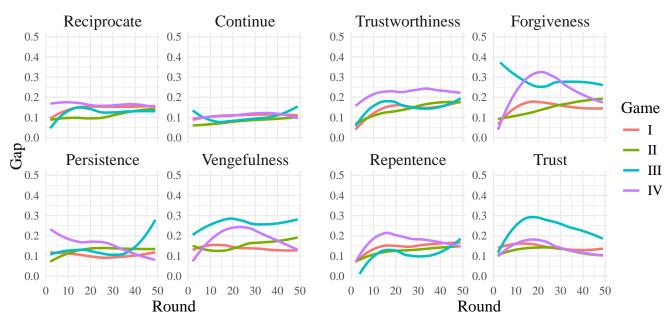


Figure 10: Sequential dependency between each individual's sharing decisions in a round and his peers' decisions in the preceding round. E.g., Panel  $CD \to DC$  represents when peer 1 cooperated and peer 2 defected in the preceding round (CD), the proportion of participants choosing to defect peer 1 and cooperate with peer 2 (DC) in the current round  $(P(DC \mid CD))$ . Game I (K-index=0.4, Minimal information, n=51 groups), Game II (K-index=0.8, Minimal information, n=51 groups), Game III (K-index=0.4, Experiential information, n=34 groups), and Game IV (K-index=0.4, Descriptive information, n=36 groups).



**Figure 11:** The 'gap' quantifies differences in how individuals treat their two peers, measured through four dimensions: persistence, vengefulness, forgiveness, and trust. A higher gap indicates greater differentiation in responses. The x-axis shows the game rounds, and the y-axis represents the gap value for each metric across different game conditions. Game I (K-index=0.4, Minimal information, n=51 groups), Game II (K-index=0.8, Minimal information, n=51 groups), Game III (K-index=0.4, Experiential information, n=34 groups), and Game IV (K-index=0.4, Descriptive information, n=36 groups).

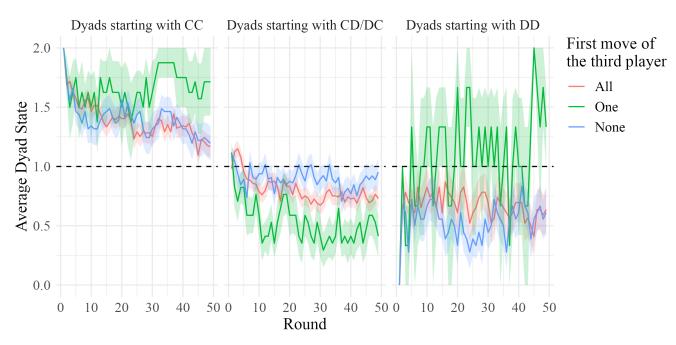


Figure 12: Influence of the Third Player on Dyadic Cooperation. The y-axis shows the average mutual sharing behavior (0: mutual defection (DD), 1: unilateral sharing (CD/DC), 2: mutual sharing (CC)). Lines represent changes in dyadic cooperation across game conditions and rounds. Game I (K-index=0.4, Minimal information, n=51 groups), Game II (K-index=0.8, Minimal information, n=51 groups), Game III (K-index=0.4, Experiential information, n=34 groups), and Game IV (K-index=0.4, Descriptive information, n=36 groups).

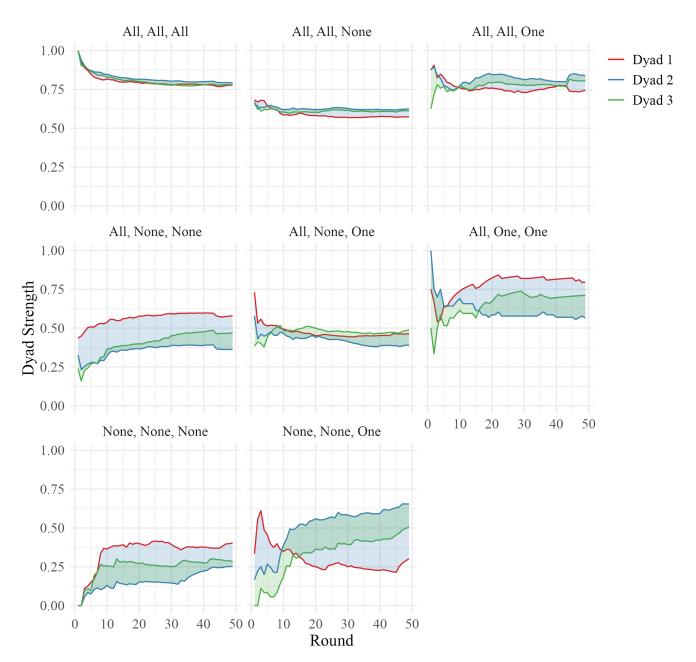


Figure 13: Evolution of Dyadic Relationship Strengths in Triads. Each panel represents different initial cooperation and defection patterns. Lines show the strength of individual dyadic relationships, and the shaded ribbon illustrates the disparity between the strongest and weakest relationships over time.

## 6. Discussion

Our research extends the current understanding of cooperation by examining how dyadic relationships aggregate and evolve within three-person groups. The findings both support and challenge existing theoretical frameworks in social dynamics and learning.

Study 1 demonstrated that structural incentives, represented by the K index, significantly influence information-sharing behavior. Groups operating with higher K-index values (0.8) showed consistently higher rates of cooperation compared to those with lower K-index values (0.4). This aligns with Rapoport's predictions and extends the findings of dyadic studies Moisan et al. (2018b) to triadic interactions. However, the effect was moderated by information availability, suggesting that structural incentives alone cannot fully explain cooperative behavior in groups.

Study 2 revealed a nuanced relationship between information availability and cooperation. Contrary to our initial expectations and previous findings in dyadic settings (Gonzalez et al., 2015), more information did not always lead to better outcomes. Although experiential information promoted cooperation as predicted by Instance-Based Learning Theory Gonzalez and Martin (2011), descriptive information decreased cooperation. This may occur because explicit payoff matrices focus participants' attention on short-term strategic calculations rather than long-term relationship building. When participants can clearly see the immediate benefits of defection, they may be more tempted to exploit cooperative partners, despite the long-term advantages of sustained cooperation. This aligns with research showing that making payoff structures explicit can trigger more competitive mindsets (Chen, Geng, Chen and Fu (2024)). Furthermore, the complexity of managing two relationships simultaneously can lead participants to default to simpler competitive strategies when presented with complete strategic information.

The evolution from universal to selective cooperation strategies observed in our study provides empirical support for theoretical models of cooperation development Nowak (2010). However, the mediating effects of the third player extend beyond current theoretical frameworks. Although social balance theory predicts stability in balanced relationships, our findings reveal how selective cooperation by a third party can either stabilize or destabilize existing dyadic relationships, depending on initial conditions.

These findings advance our theoretical understanding in several ways. First, they demonstrate that triadic structures fundamentally alter cooperation dynamics compared to dyads, particularly in how information is processed and used strategically. Second, they suggest that learning mechanisms in group settings may differ from those in dyadic interactions, with implications for cognitive modeling approaches. Third, they reveal how individual strategies aggregate to produce emergent group patterns that cannot be predicted from dyadic interactions alone.

Our findings suggest specific mechanisms for improving cybersecurity information-sharing systems. Higher cooperation in experiential learning indicates that platforms should provide clear feedback on successful threat mitigations resulting from shared intelligence. For example, organizations could receive detailed metrics showing how their shared indicators helped prevent attacks on partner organizations. The effectiveness of selective cooperation strategies suggests the implementation of tiered sharing frameworks where organizations can maintain different levels of information exchange with different partners based on reciprocity. Furthermore, the effects of the K-index indicate that policy-makers should consider tax incentives or liability protections to improve the cost-benefit ratio of sharing security intelligence.

## 6.1. Limitations and Future Work

Several limitations warrant consideration. The artificial nature of the laboratory setting may not fully capture the complexity of real-world information-sharing decisions. Our pool of MTurk participants may not represent how cybersecurity professionals make sharing decisions. The 50-round experimental design, while allowing for strategy evolution, may not reflect the indefinite time horizons of real organizational relationships. Additionally, our focus on three-person groups, while providing analytical clarity, may not be generalized to larger information-sharing networks with more complex interdependencies.

From a theoretical perspective, our study focused on a specific class of social dilemma where the chicken-type dilemma strength (Dg' = (T - R)/(R - P)) equals the stag hunt-type dilemma strength (Dr' = (P - S)/(R - P)), as is characteristic of Donor & Recipient games Wang et al. (2015). This methodological choice provided analytical clarity, but represents a limitation in generalizing our findings. In real-world scenarios, cooperative decisions may be influenced by unequal fears of exploitation (stag hunt-type) versus temptations to exploit (chicken-type). The surprising negative effect of descriptive information could be moderated by the relative balance between these dilemma types, potentially explaining the variance in information-sharing behaviors in different contexts Ito and Tanimoto (2018).

Future research directions should examine how these dynamics scale to larger networks and investigate how cognitive processes influence the simultaneous management of multiple cooperative relationships. Particularly important is understanding how Instance-Based Learning Theory can be extended to account for the concurrent management of multiple relationships, and how social balance theory can incorporate the dynamic effects of selective cooperation strategies.

Further studies should systematically vary the strengths of the dilemma independently to explore how they interact with the availability of information and the effects of third parties in triadic structures. For instance, selective cooperation by third players might have different mediating effects when fear of exploitation dominates versus when greed dominates. This approach could help address an important question raised by our findings: Why does the third player's selective cooperation strategy produce different effects depending on the dyad's initial state? The universal framework for the strength of dilemmas suggests that the effectiveness of conditional strategies can depend on whether a particular dilemma emphasizes fear of exploitation or the temptation to defect?

Future work should also test these findings with cybersecurity professionals in more realistic information-sharing scenarios to validate their applicability to real-world contexts. These extensions would provide deeper insights into the cognitive mechanisms underlying the strategic management of multiple interdependent relationships, with important implications for designing effective information-sharing systems.

# 7. Acknowledgements

This research was sponsored by the Army Research Office and accomplished under Australia-US MURI Grant Number W911NF-20-S-000 and by the Army Research Laboratory under Cooperative Agreement Number W911NF-13-2-0045 (ARL Cyber Security CRA).

# 8. Declaration of the use of AI

The authors did not use AI for any part of the work related to the manuscript submitted.

# 9. Data Availability

The data that support the findings of this study are available from the corresponding author upon reasonable request.

## References

Ackermann, K.A., Fleiß, J., Murphy, R.O., 2016. Reciprocity as an individual difference. Journal of Conflict Resolution 60, 340-367.

Barrett, S., Dannenberg, A., 2017. Tipping versus cooperating to supply a public good. Journal of the European Economic Association 15, 910–941. doi:10.1093/jeea/jvw022.

Capraro, V., Jordan, J.J., Rand, D.G., 2013. Group size effect on cooperation in one-shot social dilemmas. Scientific Reports 3, 1526.

Chen, Z., Geng, Y., Chen, X., Fu, F., 2024. Unbending strategies shepherd cooperation and suppress extortion in spatial populations. arXiv preprint arXiv:2405.19565.

Fehr, E., Gächter, S., 2000. Cooperation and punishment in public goods experiments. American Economic Review 90, 980-994.

Garrido, P., Sanner, M., Löhr, H., 2016. Shall we collaborate? a game-theoretic analysis of resistance to collaborative intrusion detection, in: 2016 IEEE Conference on Communications and Network Security (CNS), IEEE. pp. 280–288.

Gonzalez, C., et al., 2015. Effects of descriptive and experiential information on decision-making in dyadic games. Decision 2, 121–136. doi:10.1037/dec0000025.

Gonzalez, C., Ben-Asher, N., Martin, J.M., Dutt, V., 2015. A cognitive model of dynamic cooperation with varied interdependency information. Cognitive science 39, 457–495.

Gonzalez, C., Martin, J.M., 2011. Scaling up instance-based learning theory to account for social interactions. Negotiation and Conflict Management Research 4, 110–128.

Granovetter, M.S., 1973. The strength of weak ties. American journal of sociology 78, 1360–1380.

Grujić, J., Gracia-Lázaro, C., Milinski, M., Semmann, D., Traulsen, A., Cuesta, J.A., Sánchez, A., Moreno, Y., 2014. A comparative analysis of spatial prisoner's dilemma experiments: Conditional cooperation and payoff irrelevance. Scientific Reports 4, 4615.

Helbing, D., Yu, W., 2009. The outbreak of cooperation among success-driven individuals under noisy conditions. Proceedings of the National Academy of Sciences 106, 3680–3685.

Hilbe, C., Martinez-Vaquero, L.A., Chatterjee, K., Nowak, M.A., 2017. Memory-n strategies of direct reciprocity. Proceedings of the National Academy of Sciences 114, 4715–4720.

Hilbe, C., Wu, B., Traulsen, A., Nowak, M.A., 2014. Cooperation and control in multiplayer social dilemmas: Zero-determinant strategies in prisoner 2019s dilemma games. Proceedings of the National Academy of Sciences 111, 16425–16430.

- Ito, H., Tanimoto, J., 2018. Scaling the phase-planes of social dilemma strengths shows game-class changes in the five rules governing the evolution of cooperation. Royal Society open science 5, 181085.
- Jackson, J.W., 1993. Realistic group conflict theory: A review and evaluation of the theoretical and empirical literature. The Psychological Record 43, 395.
- Juvina, I., et al., 2011. The dynamics of trust and cooperation in repeated games: An experimental study of intergroup and intragroup interactions. Journal of Conflict Resolution 55, 939–965.
- LaPorte, P., Hilbe, C., Nowak, M.A., 2023. Adaptive dynamics of memory-one strategies in the repeated donation game. PLoS Computational Biology 19, e1010987.
- Michael, J., McEllin, L., Felber, A., 2020. Prosocial effects of coordination-what, how and why? Acta psychologica 207, 103083.
- Moisan, F., ten Brincke, R., Murphy, R.O., Gonzalez, C., 2018a. Not all prisoner 2019s dilemma games are equal: Incentives, social preferences, and cooperation. Decision 5, 306–318.
- Moisan, F., ten Brincke, R., Murphy, R.O., Gonzalez, C., 2018b. Not all prisoner's dilemma games are equal: Incentives, social preferences, and cooperation. Decision 5, 306.
- Nax, H.H., Burton-Chellew, M.N., West, S.A., Young, H.P., 2023. Information design and human cooperation in strategic interactions. Nature Communications 14, 1–9. doi:10.1038/s41467-023-36847-9.
- Nowak, M.A., 2010. The evolution of cooperation: A retrospective. Philosophical Transactions of the Royal Society B: Biological Sciences 365, 73–79.
- Perc, M., Jordan, J.J., Rand, D.G., Wang, Z., Boccaletti, S., Szolnoki, A., 2017. Statistical physics of human cooperation. Physics Reports 687, 1–51.
- Rakow, T., Newell, B.R., 2010. Degrees of uncertainty: An overview and framework for future research on experience-based choice. Journal of Behavioral Decision Making 23, 1–14.
- Rapoport, A., 1967. A note on the" index of cooperation" for prisoner's dilemma. Journal of Conflict Resolution 11, 100-103.
- Rapoport, A., Chammah, A.M., Orwant, C.J., 1965. Prisoner's dilemma: A study in conflict and cooperation. volume 165. University of Michigan press.
- Romano, A., Balliet, D., Yamagishi, T., Liu, J.H., 2017. Information content and cooperation in social dilemmas. Nature Human Behaviour 1, 1–7. doi:10.1038/s41562-017-0114.
- Santos, F.C., Santos, M.D., Pacheco, J.M., 2008. Social diversity promotes the emergence of cooperation in public goods games. Nature 454, 213–216
- Sinervo, B., Lively, C.M., 1996. The rock-paper-scissors game and the evolution of alternative male strategies. Nature 380, 240-243.
- Szolnoki, A., Perc, M., 2019. Modelling the impact of group size on cooperation in spatial public goods games. Physical Review E 100, 052312.
- Taylor, C., Nowak, M.A., 2007. Transforming the dilemma. Evolution 61, 2281–2292. doi:10.1111/j.1558-5646.2007.00196.x.
- Tosh, D.K., Shetty, S., Sengupta, S., Bagchi, S., 2015. An evolutionary game-theoretic framework for cyber-threat information sharing. IEEE Transactions on Information Forensics and Security 10, 2835–2850.
- Van Lange, P.A., Balliet, D., Joireman, J., 2020. Cooperative interactions in groups of different sizes and the role of leadership. Nature Human Behaviour 4, 916–923.
- Vuolevi, J.H., Van Lange, P.A., 2012. Boundaries of reciprocity: Incompleteness of information undermines cooperation. Acta Psychologica 141, 67–72.
- Wang, Z., Kokubo, S., Jusup, M., Tanimoto, J., 2015. Universal scaling for the dilemma strength in evolutionary games. Physics of life reviews 14, 1–30.
- Yang, Z., 2023. Role polarization and its effects in the spatial ultimatum game. Physical Review E 108, 024106.

# Toward a Cognitive Theory of Interdependent Decisions in Groups: Dynamic Prosociality, Categorization, and Contrast

Yinuo Du, Palvi Aggarwal, Kuldeep Singh, Fei Fang, and Cleotilde Gonzalez

Carnegie Mellon University

University of Texas at El Paso

# **Author Note**

Correspondence concerning this article should be addressed to Cleotilde Gonzalez, Social and Decision Science Department, Carnegie Mellon University, Pittsburgh, PA. Email: coty@cmu.edu

## Abstract

We analyze the dynamics of strategic interaction among a group of human agents through a novel cognitive model that integrates three key psychological mechanisms: dynamic prosociality, category learning, and contrast effects. The dynamic prosociality mechanism enables individuals to adjust how much they value others' choices and outcomes based on expectation-reality discrepancies. The category learning mechanism captures how people efficiently organize their social experiences into behavioral prototypes through hierarchical clustering. The contrast effect sharpens the distinctions between these behavioral categories by amplifying perceived differences between groups based on their relative positions along behavioral dimensions. Using data from online group experiments, we demonstrate that the model successfully reproduces human behavior patterns without parameter fitting. Through detailed analysis of dynamic prosociality, we gain insight into the psychological processes underlying how individuals evaluate and respond to others in group settings. These findings advance our understanding of human cognition in complex social environments and suggest ways to improve collective outcomes in real-world applications.

 $\label{lem:keywords:} \textit{Keywords:} \ \text{Instance-Based Learning Theory, Prosociality, Category learning,}$  Contrast effect

# Toward a Cognitive Theory of Interdependent Decisions in Groups: Dynamic Prosociality, Categorization, and Contrast

## Introduction

Human social systems are defined by strategic interdependence, where individual choices collectively shape broader outcomes. While much research has focused on simple dyadic interactions, real-world contexts demand the management of numerous simultaneous relationships, all under cognitive limitations. Whether in cybersecurity information sharing, organizational resource allocation, or international diplomacy, individuals must continuously monitor, assess, and respond to multiple partners—often with constrained cognitive bandwidth (CISA, 2023; Stevens et al., 2018; P. A. Van Lange et al., 2013).

Current approaches to modeling behavior in strategic team interactions typically fall into three broad categories. Evolutionary approaches (Li et al., 2023) effectively capture population-level outcomes, but often abstract away individual cognitive processes. Game-theoretic frameworks provide precise mathematical formulations, but generally rely on strong rationality assumptions that do not always align with human behavior. In contrast, social dilemma research (D. Balliet & Van Lange, 2013; D. P. Balliet et al., 2017; P. A. M. Van Lange et al., 2014) offers psychologically grounded insights into human cooperation, often focusing on motivational and contextual factors rather than formal rationality. Cognitive modeling has emerged as a promising fourth route, explicitly addressing how humans navigate social learning under cognitive constraints (Gallotti & Grujić, 2018; Martin et al., 2014; Shum et al., 2019). Although this approach has shown success in dyadic settings (Gonzalez et al., 2015), extending these models to multi-agent contexts presents unique challenges.

A central challenge is that humans often struggle to track interactions with multiple partners simultaneously, a limitation shaped by fundamental cognitive constraints in working memory and attention. As it was established long time ago, the human working memory capacity is limited (Cowan, 2001; Miller, 1956; Simon, 1974), making it

challenging to simultaneously track the detailed behavioral histories of numerous interaction partners (Stiller & Dunbar, 2007). Evidence from social network studies confirms this constraint, with some researchers demonstrating that cognitive limitations restrict the number of concurrent relationships humans can maintain (Dunbar, 1998). Macrae and Bodenhausen (2000) showed that when cognitive resources are stretched thin, humans resort to categorical processing of social information rather than individual processing. For instance, when managing multiple workplace relationships simultaneously, people may classify colleagues as 'reliable team players' or 'self-interested actors' rather than maintaining detailed records of each person's specific actions and motivations. This fundamental constraint on human information processing is especially problematic in multi-agent contexts, where the number of unique dyadic relationships increases quadratically with the number of agents (Dziura et al., 2023). When interacting with only one partner, people can dedicate sufficient cognitive resources to track detailed sequential behavioral patterns, but this capacity can quickly become overwhelmed as the social environment becomes more complex.

Our research addresses these challenges by proposing a novel cognitive framework for interdependent decisions in teams. Our new model integrates three key psychological mechanisms to navigate complex social environments: (1) dynamic prosociality, a refined version of the surprise-driven weight update mechanism of (Gonzalez et al., 2015) that addresses limitations in the original formulation, allowing more discriminative responses to different partners; (2) category learning that allows individuals to efficiently manage multiple relationships by grouping similar partners rather than tracking each one individually (Rosch, 1978); and (3) contrast effects that sharpen distinctions between different sequential behavior patterns, facilitating more effective discrimination between cooperation partners (Wu et al., 2020b).

Based on Instance-Based Learning Theory (IBLT) (Gonzalez et al., 2003), our proposed framework captures how people process multiple social relationships through

prototype-based categorization (Tamarit et al., 2018), evaluate actions through social comparison (Chierchia et al., 2017), while at the same time dynamically adjusts their prosocial strategies based on observed reciprocity (D. Balliet & Lindström, 2023; Kleiman-Weiner et al., 2016a). Our research demonstrates how these cognitive mechanisms, dynamic prosociality, category learning, and contrast effects, enable effective multi-agent strategic decision making despite inherent human cognitive limitations. We validated our model using data from online group experiments involving strategic social dilemmas, demonstrating its ability to reproduce human behavior patterns without parameter fitting. The detailed analysis of the dynamic prosociality parameter  $(\alpha)$  reveals significant practical implications beyond theoretical understanding, including how individuals develop differential responses to cooperation partners based on their interaction history. These findings suggest design principles for interventions and systems that better align with natural cognitive processes, potentially improving cooperative outcomes in various domains of interdependent decision making. By bridging cognitive mechanisms with strategic sophistication, our work offers a foundation for understanding and enhancing cooperation in complex social environments where individuals must manage and navigate multiple concurrent relationships.

## Related Work

The study of repeated strategic interactions between interdependent agents has a rich research history. Early work focused on simple strategies with minimal partner modeling. Axelrod's seminal computer tournaments of the Iterated Prisoner's Dilemma (Axelrod, 1984) demonstrated the success of Tit-for-Tat (TFT), which only considers the partner's last action. Although these simple strategies proved to be remarkably effective in structured environments, subsequent research revealed their limitations in noisy or complex settings (Nowak, 2006). This led to increasingly sophisticated approaches incorporating richer agent modeling and learning mechanisms, including modern machine learning methods and cognitively inspired strategies. Two fundamental challenges have emerged in

this progression: the computational demands associated with memory and learning and the complexity of modeling diverse agent strategies. Our work addresses these challenges by incorporating cognitive mechanisms for efficient memory use and agent categorization.

# Learning and Agent Modeling in Interdependent Interactions

A significant theoretical advancement in agent modeling came with Earnest (2013)'s discovery of zero-determinant strategies, which established mathematical boundaries on strategy effectiveness. Zero-determinant strategies allow an agent to enforce specific relationships between its own payoff and that of its partner, effectively controlling the distribution of payoffs without the partner's cooperation. Stewart and Plotkin (2013) extended this work by demonstrating that "generous" variants often outperform purely extortionate strategies in evolutionary settings, highlighting how successful strategies must balance exploitation with mutual benefit.

The development of agent modeling approaches has followed several trajectories. Early work focused on explicit prediction of others' actions through pattern recognition (Carmel & Markovitch, 1995), while later approaches incorporated uncertainty and partial observability (Gmytrasiewicz & Doshi, 2005). Modern machine learning methods, particularly deep reinforcement learning, have demonstrated impressive success in learning implicit representations of agent behavior (Harper et al., 2017; Lowe et al., 2017). These approaches can uncover sophisticated counterstrategies through extensive self-play and experience accumulation, often exceeding hand-crafted strategies in complex environments.

However, the increasing sophistication of learning algorithms has led to an "arms race" in strategy complexity. Neural network-based approaches can provide highly complex patterns in how agents respond to different situations (Leibo et al., 2017), allowing more context-sensitive and adaptive responses, but also making strategies harder to interpret and analyze. This complexity creates challenges for theoretical analysis and raises questions about the robustness of the learned strategies. Some studies suggest that simpler strategies with clear theoretical foundations may be more robust among diverse interaction

partners (Wang et al., 2018).

The tension between strategy complexity and robustness has motivated research into hybrid approaches that combine machine learning with domain knowledge. For example, Crandall et al. (2018) demonstrated how the incorporation of simple mechanisms that promote mutual benefit in learning algorithms can improve generalization between different interaction partners. Similarly, Kleiman-Weiner et al. (2016a) showed that learning algorithms constrained by the principles of game theory often develop more stable and interpretable strategies.

Recent work has increasingly focused on multi-agent scenarios in which agents must simultaneously model and adapt to multiple partners (Lanctot et al., 2017). This setting introduces additional complexities, as agents must balance their responses between different partners while maintaining coherent strategies. The challenge is compounded in settings with incomplete information or when partners may change their strategies over time (Hernandez-Leal et al., 2019).

# Memory Constraints and Cognitive Plausibility

While machine learning approaches have demonstrated impressive performance in agent modeling, they typically assume unlimited memory capacity and computational resources. These approaches often maintain complete interaction histories or complex state representations, enabling sophisticated pattern recognition but diverging significantly from human cognitive constraints. This disconnect raises important questions about the psychological plausibility and practical applicability of such models to emulate human interdependencies in multi-agent scenarios.

Empirical studies reveal clear limitations in human memory use during strategic interactions. Research consistently shows that humans typically access only a handful of previous interactions when making decisions (Moreira et al., 2013), indicating a clear cognitive bottleneck. This limitation reflects broader constraints on working memory capacity, which affects how individuals process and utilize information in dynamic social

situations. Memory traces follow systematic decay patterns (Anderson & Schooler, 1991), with recent interactions more heavily weighted while maintaining the diminishing influence of established patterns, a phenomenon known as the power law of forgetting.

The relationship between memory complexity and strategy performance follows an inverted U-shaped pattern (Hertwig & Erev, 2009), suggesting optimal performance at intermediate levels of memory complexity. This finding has profound implications for the design of strategies. Although too little memory prevents recognition of important behavioral patterns, excessive memory complexity can lead to overfitting and reduced adaptability. This balance reflects the fundamental principles of bounded rationality (Simon, 1990), where cognitive constraints paradoxically contribute to more robust and adaptable decision making.

Recent work has highlighted the critical role of cognitive constraints in shaping human decision-making, especially in multi-agent contexts where individuals must concurrently track and respond to multiple peers. Human cognitive constraints such as limited working memory and attentional resources directly influence how effectively individuals manage and learn from multiple ongoing interactions. For example, empirical studies demonstrate that humans typically rely on simplified cognitive strategies, such as categorization of the type of partners, when faced with the complex task of tracking multiple interaction partners (Macrae & Bodenhausen, 2000; Stevens et al., 2018). This reliance on simplified strategies in multi-agent contexts suggests that cognitive constraints may play a critical role in determining how robust and adaptable human strategies are in group interactions, and that strategies such as categorization of partners may be a coping mechanism to address the cognitive limitations into models of interdependent decision-making.

Overall, these findings suggest that effective strategies should not simply operate within memory constraints but actively leverage them as design principles. Memory limitations can serve as natural regularizers, promoting generalization by preventing

overfitting to specific interaction patterns. This perspective aligns with the ecological rationality frameworks (Todd et al., 2012), which emphasize how cognitive constraints can improve decision making in natural environments.

### Categorical Learning and Contrast Effects

A fundamental challenge in social dilemmas is the wide space of possible peer strategies. As the diversity of peers increases, the complexity of the modeling increases exponentially (Lim et al., 2016), making the modeling of direct strategies computationally intractable. This challenge becomes particularly acute in multi-agent settings where traditional modeling approaches often fail to scale effectively or require unrealistic computational resources.

Humans address this complexity through sophisticated inductive categorical learning mechanisms that enable efficient but flexible social learning. Research shows that people actively form and update categories based on patterns of interdependence in social interactions (Martin et al., 2014). These categories may serve not just as simplifying heuristics, but as predictive models that guide future cooperation decisions. For example, when individuals identify patterns of reciprocity or exploitation, they develop categorical representations that help them anticipate and respond to similar behaviors in new interactions (Kleiman-Weiner et al., 2016b).

What makes categorical learning particularly powerful is its ability to balance efficiency with effectiveness. Although categorization reduces the granularity of social information, it paradoxically allows more sophisticated responses by capturing essential behavioral patterns (Chierchia et al., 2017). People continually refine these categories based on new experiences, maintaining a dynamic equilibrium between stable categorical knowledge and adaptability to novel patterns. This process of category refinement is strongly influenced by the social context: Individuals' classifications of "cooperative" versus "non-cooperative" behavior emerge relative to their broader social experience (Gonzalez et al., 2015; Wu et al., 2020b).

These categorical learning mechanisms have been demonstrated in various social dilemmas. In cybersecurity information-sharing networks, Mermoud et al. Mermoud et al. (2019) found that defenders naturally categorize their peers into "regular sharers" and "free-riders" based on sharing patterns, using these categories to guide their own sharing decisions even with new peers. Similarly, in organizational contexts, studies of group-based resource allocation show that managers develop categorical representations of "reciprocators" versus "opportunists" that influence future resource-sharing decisions (Hámornik & Krasznay, 2017).

The power of categorical learning is particularly evident in repeated interaction settings. For example, in public goods games, participants rapidly develop categories for "consistent contributors" and "strategic free-riders," with these categories shaping not only direct responses but also reputation sharing within groups (Fehr & Schurtenberger, 2019). These categories prove to be remarkably stable - once an individual is categorized as a reliable cooperator, isolated defections are often discounted as anomalies rather than prompting immediate category reassignment (Andreoni & Miller, 1993).

Experimental studies of group cooperation reveal how categorical learning enables efficient decision making under time pressure. When faced with multiple potential cooperation peers, participants may not track detailed histories, but instead maintain broader categorical assessments like "trustworthy," "unpredictable," or "exploitative" (Kelley & Stahelski, 1970). These categorical judgments are particularly influential in early interactions with new peers, where they serve as default expectations until individual-specific evidence accumulates (Zhang et al., 2019).

The categorical perception of peers introduces systematic contrast effects in behavior evaluation. Rather than evaluating each peer's actions in isolation, individuals evaluate behaviors relative to their experiences with other peers (Young et al., 2019). These contrast effects are particularly pronounced between categorically distinct peers. For instance, Kirchkamp et al. (2016) demonstrated that players' responses to moderately

cooperative behavior become more positive when they simultaneously interact with clearly non-cooperative peers, suggesting that categorical boundaries enhance behavioral discrimination.

The sophistication of categorical human learning extends beyond simple classification. Successful players develop hierarchical category structures, with broad behavioral types that contain subtypes that capture more nuanced patterns (Rand et al., 2020). This hierarchical organization allows players to balance computational efficiency with strategic sophistication. Moreover, these learned categories are effectively transferred between different economic games (Peysakhovich & Lerer, 2018), suggesting that categorical learning captures fundamental aspects of strategic behavior.

### Cognitive Approaches to Social Learning

Category learning represents a fundamental cognitive mechanism that influences how individuals perceive, process, and retain information about social interactions. Huang and Sherman (2018) describe how attentional mechanisms significantly impact social perception through category accentuation, where individuals exaggerate differences between groups while minimizing within-group variations. Sherman et al. (2009) further demonstrate that such cognitive biases can enhance memory for features associated with majority groups while reducing recall of minority group characteristics, highlighting how categorization processes can systematically shape learning outcomes.

The contrast effect, a key phenomenon in category-based perception, is significantly influenced by the social distance between learners and their interaction partners. In social networks, the principle of homophily—the tendency of individuals to associate with others who are similar to themselves—plays a crucial role in shaping these interactions (Centola, 2011; McPherson et al., 2001). Research indicates that as social distance increases, individuals are more likely to focus on abstract goals rather than concrete behaviors (Hansen et al., 2016; Kalkstein et al., 2016). In contrast, individuals who perceive themselves as similar to their interaction partners are more likely to closely observe,

imitate, and learn specific behaviors from these partners, consistent with homophily-driven interactions (Centola, 2011). This relationship between social distance and learning style suggests that the position of individuals relative to their interaction partners fundamentally shapes how they process and internalize social information.

Multiple theoretical frameworks have been proposed to explain how individuals navigate and learn from social interactions. Reinforcement learning models (Erev & Roth, 2006) focus on outcome-based behavioral adjustments but often struggle with the dynamic nature of social environments. Bayesian approaches (Baker et al., 2011) represent uncertainty through probabilistic beliefs about others' intentions, but frequently assume unrealistic inferential capabilities. Theoretical frameworks of the mind (Yoshida et al., 2008) emphasize meta-representational abilities but may overestimate typical cognitive capacities in complex scenarios. Heuristic approaches (Gigerenzer & Gaissmaier, 2011) propose that simple decision rules can achieve effective social coordination despite limited information processing.

Among these various frameworks, the Instance-Based Learning Theory (IBLT) (Gonzalez et al., 2003) offers a particularly compelling account of social learning under cognitive constraints. Unlike approaches that oversimplify cognitive processes or assume unrealistic computational capabilities, IBLT provides a psychologically grounded explanation for how individuals learn from specific experiences while respecting memory limitations. Through mechanisms like activation decay and similarity-based retrieval, IBLT naturally explains how categorical thinking emerges from interactive experiences. Memory constraints guide attention toward meaningful patterns rather than exhaustive details, leading to more robust and generalizable learning (Hertwig & Erev, 2009). However, the use of IBLT in social dilemmas and multi-agent interactions is significantly under-developed.

Gonzalez et al. (2015) proposed a model of dyadic interdependence using the Prisoner's Dilemma as an example (named IBL-PD). The IBL-PD model captures how

individuals learn to cooperate or defect in repeated two-player interactions by retrieving experiences stored as instances in memory. Decisions depend on activation and blending mechanisms that estimate the value of actions based on past outcomes. A unique aspect of this model is its dynamic prosociality parameter  $(\alpha)$ , which quantifies how much a player values the outcomes of their interaction partner relative to their own outcomes, adapting based on the history of interactions.

However, the original IBL-PD formulation is insufficient to model larger groups because it focuses exclusively on dyadic (two-player) interactions, which do not fully capture multi-agent social environments. When interactions expand beyond two individuals, managing cognitive load and maintaining effective cooperation becomes significantly more challenging, demanding mechanisms beyond simple memory retrieval and outcome-weighting for a single other agent.

Thus, we build upon the IBL-PD model by extending it with two additional cognitive mechanisms—category learning and contrast effects. Category learning helps individuals efficiently organize multiple partners by grouping similar peers into manageable behavioral prototypes, thus reducing cognitive complexity. The contrast effect sharpens the perceptual distinctions between these categories, enhancing the discriminative capacity of individuals to respond appropriately to various interaction partners. By integrating these new components, our extended model addresses the fundamental challenge of managing numerous social relationships within realistic cognitive constraints, allowing for a more accurate representation and prediction of human behavior in multi-agent strategic interactions.

#### Instance-Based Learning Theory for Dyadic Interactions

IBLT provides a formal framework to model how humans learn from experience while respecting cognitive constraints (Gonzalez, 2024). In social dilemmas, instances represent specific interaction experiences, storing not only outcomes but also information about the social context and others' actions. The IBL-PD model captures how individuals

learn to cooperate or defect in the Prisoner's Dilemma through repeated two-player interactions by retrieving and blending memory traces of past outcomes (Gonzalez et al., 2015). For each dyadic interaction, agents store instances in memory of the form [PeerIndex, MyAction, PeerAction, MyOutcome, PeerOutcome], which allows them to track both their own and their partner's behavior over time.

The core mechanism of IBL-PD lies in its activation and blending processes. Instances i are stored in memory  $\mathcal{M}$  and become more or less accessible depending on their frequency, recency, and similarity to the current context. Activation for an instance i is computed as:

$$A_i(t) = \ln\left(\sum_{t' \in \mathcal{T}_i(t)} (t - t')^{-d}\right) + \mu \sum_{j \in \mathcal{F}} \omega_j (S_{ij} - 1) + \sigma \xi \tag{1}$$

The parameters are: d for decay,  $\mu$  for mismatch penalty,  $\omega_j$  for feature weights, and  $\sigma$  for the scale of Gaussian noise  $\xi \sim \mathcal{N}(0,1)$ . Default values are  $(d,\mu,\omega_j,\sigma) = (0.5,1,1,0.25)$ . The similarity term  $S_{ij}$  reflects how closely the stored instance matches the current context.

IBL-PD uses a blending process to compute the expected value of each action option k based on retrieved instances. Crucially, Gonzalez et al. (2015) incorporated social preferences into this blended value by introducing a dynamic prosociality parameter  $\alpha$ , which represents the degree to which the agent values the peer's outcomes:

$$V_k(t) = \sum_{i=1}^n P_{ik}(x_{self} + \alpha(t) \cdot x_{other})$$
(2)

Here,  $x_{self}$  and  $x_{other}$  are the outcomes for self and peer in instance i, and  $P_{ik}$  is the retrieval probability for instance i under action k (computed using a softmax function over the activation values). The dynamic prosociality parameter  $\alpha(t)$  is updated after each interaction:

$$Gap(t) = Abs(V_k(t-1) - (x_{self} + \alpha(t)x_{other}))$$
(3)

$$\alpha(t+1) \leftarrow (1-\eta) \,\alpha(t) + \,\eta \,(1 - \hat{Gap}(t)) \tag{4}$$

Where  $\hat{Gap}(t) \in [0, 1]$  is the normalized gap and  $\eta$  is the learning rate. This learning rule captures the intuition that agents reduce regard for peers when outcomes deviate from expectations.

# A Correction to Dynamic Prosociality

While intuitive, this formulation fails to distinguish between "good" and "bad" surprises. For example, if a peer consistently defects, they become predictable, resulting in a low gap and a higher  $\alpha$  even though they are not cooperative. This creates a problematic learning dynamic in which predictability is mistaken for trustworthiness.

To address this limitation, we propose a refined update rule that conditions the update on whether the peer exceeded expectations or fell short:

$$Gap(t) = V_k(t-1) - (x_{self} + \alpha(t)x_{other})$$
(5)

$$\alpha(t+1) \leftarrow \begin{cases} (1-\eta)\,\alpha(t) + \eta \, \max\left(\alpha(t), \, Gap(t)\right), & \text{if } Gap(t) \ge 0, \\ (1-\eta)\,\alpha(t) + \eta \, Gap(t), & \text{if } Gap(t) < 0. \end{cases}$$
(6)

This rule ensures that  $\alpha$  increases only when a peer exceeds expectations and otherwise decreases. The exponential moving average structure preserves gradual adaptation while differentiating between positive and negative deviations from expectations.

To demonstrate the limitations of the original  $\alpha$  formulation in Gonzalez et al. (2015) and validate our proposed correction, we conducted a controlled simulation experiment. We created two sets of 100 identical IBL-PD agents, differing only in their  $\alpha$  update mechanism - one set using the original formulation from Gonzalez et al. (2015) and

the other using our revised formulation. Each agent played 50 rounds of the Prisoner's Dilemma against two fixed-strategy peers simultaneously: a fully cooperative peer (100% cooperation rate) and a fully defective peer (0% cooperation rate). We tracked the evolution of each agent's  $\alpha$  value toward each peer over these 50 rounds, averaging results across all 100 agents for each formulation. Figure 1 presents these average  $\alpha$  trajectories. Under the original formulation (Figure 1a), agents develop similarly high  $\alpha$  values for both peers despite their opposing behaviors. With our revised formulation (Figure 1b), agents appropriately discriminate between peers, developing high  $\alpha$  values (approximately 0.9) for cooperative peers while maintaining low values (approximately 0.1) for defective peers.

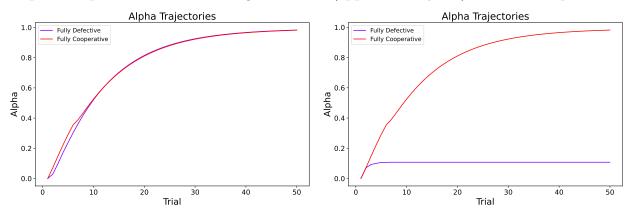


Figure 1

Alpha trajectories for a focal agent interacting with a fully cooperative (red) and a fully defective (purple) peer. (Left) Under the original formulation, the agent converges to similar  $\alpha$  values for both peers. (Right) With the revised formulation, the agent assigns a higher  $\alpha$  to the cooperative peer while maintaining a near-zero  $\alpha$  for the defective peer.

To verify that our revised  $\alpha$  formulation preserves the core behavioral dynamics captured in Gonzalez et al. (2015), we paired 100 IBL-PD agents using our revised  $\alpha$  formulation to play 200 rounds of the Prisoner's Dilemma against each other, as done by Gonzalez and colleagues. Each agent received complete information on the actions and outcomes of both players after each round. All model parameters except the  $\alpha$  update rule remained identical to those in Gonzalez et al. (2015): default values of d=5 for memory

decay and r = 1.5 for noise. Figure 2 shows the mutual cooperation rate averaged across all agent pairs in each trial. The simulation reproduces the characteristic pattern observed in human data: an initial cooperation rate around 0.55 that quickly declines to approximately 0.20 by trial 50, followed by a gradual increase to about 0.40 by trial 200. This U-shaped cooperation curve closely matches the human behavior reported in the original study, confirming that our modification to the  $\alpha$  update mechanism preserves the model's ability to capture dynamic cooperation patterns.

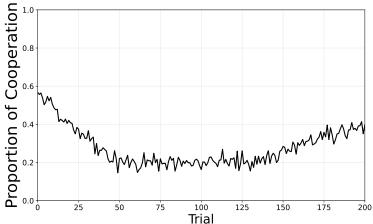


Figure 2

Cooperation behavior over 200 trials in the dynamic prisoner's dilemma task. The curve illustrates our model's predictions showing the characteristic pattern observed in human behavior: initial high cooperation that quickly declines as strategic adaptation occurs, followed by a gradual increase in later trials. This pattern successfully replicates the findings reported by (Gonzalez et al., 2015), validating that our revised  $\alpha$  formulation preserves the core behavioral dynamics.

# IBL-Group: A New Model of Interdependent Decisions in Groups

While the dynamic prosociality parameter  $\alpha$  enables discriminatory decision making between two peers, scaling to larger groups requires more sophisticated cognitive mechanisms to process multiple relationships. Our IBL-group model addresses this challenge with three cognitive mechanisms: (1) a revised dynamic prosociality mechanism that adapts  $\alpha$  based on individual experiences (explained above), (2) category learning to organize social experiences, and (3) contrast effects to improve perceptual distinctions between behavior types.

Category learning allows agents to efficiently process multiple relationships by organizing peers into behavioral prototypes. Each prototype represents a distinct behavioral pattern characterized by a multidimensional feature vector that captures the essential dimensions of strategic behavior. As illustrated in Figure 3, category learning guides memory retrieval by identifying similar peers (for example, P1 and P3 share Prototype 1), allowing the agent to access experiences from multiple related peers while making decisions about one of them.

Contrast effect then sharpens how we perceive differences between categories. In social perception, this means that distinctive behaviors become more prominent when contrasted with different types of behavior. When evaluating the instances accumulated through interaction with a partner who consistently chose action A (80% of the time in the past), the presence of another partner showing predominantly opposite behavior will lead to relatively higher activations for instances with action A through the contrast effect, as shown by the darker shading in the activated instances. These mechanisms work sequentially: category learning determines which experiences are retrieved, while the contrast effect modulates how strongly each retrieved experience influences the current decision through activation adjustment.

## Category Learning

Each behavioral category is defined by a five-dimensional feature vector that includes action tendency, entropy, responsiveness, recovery propensity, and volatility. The action tendency indicates the proportion of times an agent chooses a particular action from their available options. Entropy reflects the unpredictability or randomness in an agent's action sequence. Responsiveness measures how much an agent's current action is influenced by their partner's previous action, capturing the reactive nature of strategic interactions. Recovery propensity represents how quickly an agent returns to a consistent behavior

pattern following a deviation. Volatility characterizes the frequency with which an agent changes their action choices over time. Together, these features provide a comprehensive, yet parsimonious description of strategic behavior observed in various repeated interaction contexts (Axelrod & Hamilton, 1981). Although additional features (e.g., memory length) could be considered, this five-dimensional framework strikes a balance between capturing the essential elements of decision making and keeping the model computationally manageable. Moreover, the partial overlap among these features helps filter out random noise and highlights subtle differences in behavior that might be overlooked by less transparent statistical methods like principal component analysis.

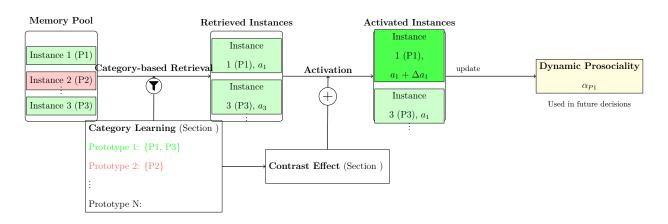


Figure 3

Illustration of the decision-making process for peer P1. Instances associated with P1 and related peers (e.g., P3) that belong to Prototype 1 are retrieved from memory, then selectively enhanced through contrast effects. The resulting activation levels directly update the dynamic prosociality parameter  $\alpha_{P1}$ , which shapes future cooperation decisions. Activation intensity is represented by shading.

The model employs prototype-based categorization with dynamic category learning. Initial classification attempts to match agents with predefined prototypes. Unmatched agents enter a clustering process, where stable clusters can be promoted to new prototypes.

Existing work has shown psychological evidence that humans approach interactions with pre-existing behavioral schemas (Fiske & Taylor, 2020). These initial categories serve as cognitive anchors that facilitate the rapid assessment of interaction partners. This matches the findings that people first assess others using broad dispositional categories before developing more nuanced impressions through experience (Gilbert & Malone, 1995). We choose to start with two contrasting behavioral prototypes that serve as natural reference points. This anchoring approach allows the system to bootstrap learning while maintaining the flexibility to discover intermediate clusters (Kelley & Stahelski, 1970). So we can capture how cognitive agents initially categorize others along fundamental behavioral dimensions before developing more sophisticated representations.

When computing blended values for potential actions, the agent retrieves instances not only from direct interactions with the target agent but also from all agents within the same behavioral category. This categorization-based retrieval affects the surprise calculation through the blended value  $V_k$ , which now incorporates the instances of categorically similar agents. Consequently, the gap between expected and actual outcomes reflects a deviation from category-level expectations rather than purely individual-level predictions. The resulting surprise value updates the values of  $\alpha$  for all agents within the same category, capturing how cognitive systems could adjust their regard to groups of similar actors rather than processing each relationship in isolation.

As detailed in Algorithm 1, interaction experiences with peers are processed through fingerprint-based categorization. Using a sliding window of size w, the algorithm calculates a five-dimensional fingerprint vector (lines 2-3). This fingerprint is compared with existing prototypes using cosine similarity, with a confidence threshold  $\theta_{conf}$  determining classification (lines 4-8).

For these unclassified agents, Algorithm 2 is applied to dynamically discover new behavioral categories. When the number of unclassified agents exceeds the minimum cluster size m, hierarchical clustering is used to group similar behavior patterns (lines 2-3).

### Algorithm 1 Individual Agent Categorization

10: **end if** 

11: **return** c, f

```
Require: Action sequence, prototypes P, window size w, confidence threshold \theta_{conf}
 1: Initialize agent fingerprint f, category c
 2: if sequence length \geq w then
         f \leftarrow \text{ComputeFingerprint(sequence[-w:])} \triangleright \text{Action tendency, entropy, responsiveness}
 3:
         (strategy, conf) \leftarrow \text{MatchPrototypes}(f, P)
 4:
        if conf \geq \theta_{conf} then
 5:
 6:
             c \leftarrow strategy
         else
 7:
             c \leftarrow "unclassified"
 8:
 9:
        end if
```

The algorithm maintains a record of historical fingerprints - computed from previous classification attempts - to monitor cluster stability over time (lines 4-5). Once a cluster exhibits consistent behavior across  $\theta_{stab}$  observations, it is promoted to prototype status by computing the centroid of its member fingerprints (lines 6-10). This two-phase process, where Algorithm 1 categorizes individual agents and Algorithm 2 refines and expands the set of prototypes, ensures efficient categorization while retaining adaptability to novel behavioral patterns.

The divisive clustering algorithm (Algorithm 3) balances two competing cognitive requirements: the need to form meaningful behavioral categories that guide future interactions, and the constraint of human working memory capacity (Miller's 1956 magic number  $7\pm2$ ) (Miller, 1956). This constraint motivates our limit of 9 total groups, as research shows that humans struggle to maintain and effectively utilize more complex categorization schemes in real-time strategic interactions (Stevens et al., 2018). At each decision point, the algorithm selects the split metrics based on their discriminative power,

## Algorithm 2 Cluster Management and Prototype Promotion

15: return  $P_{new}$ 

```
Require: Unclassified fingerprints F, min cluster size m, stability threshold \theta_{stab}
 1: Initialize cluster history H, prototype updates P_{new}
 2: if |F| \ge m then
         labels \leftarrow Cluster(F)
                                                                                     ▶ Hierarchical clustering
 3:
         for each valid cluster k in labels do
 4:
             H[k] \leftarrow H[k] \cup \{\text{fingerprints with label } k\}
 5:
             if |H[k]| \ge \theta_{stab} then
 6:
                  stability \leftarrow ComputeStability(H[k])
 7:
                  if stability \ge \theta_{stab} then
 8:
                      p_{new} \leftarrow \text{Mean}(H[k])
 9:
                      P_{new} \leftarrow P_{new} \cup \{p_{new}\}
10:
                  end if
11:
             end if
12:
         end for
13:
14: end if
```

measured by the range of observed values  $(r = \max(F[:, m]) - \min(F[:, m]))$ . This approach reflects how natural categories often form around observable behavioral variations that meaningfully distinguish between different strategies (Rosch, 1978). The algorithm identifies potential category boundaries by looking for natural gaps in behavioral metrics (differences > r/10), similar to how humans tend to form categories around clusters of similar experiences rather than through arbitrary divisions (Martin et al., 2014).

The algorithm terminates when: reaching the maximum of 9 groups (cognitive capacity constraint), having too few agents to meaningfully divide (< 4 per group) or exhausting behaviorally meaningful splits ( $r \le threshold$ ). This helps ensure that the resulting categorization scheme remains both cognitively manageable and strategically

useful (Hertwig & Erev, 2009). If the natural clustering process ends with a tree depth less than 2, we force a binary split using the action\_tendency metric. This mechanism ensures that categorization maintains at least a basic level of strategic discrimination (Macrae & Bodenhausen, 2011).

# Algorithm 3 Divisive Clustering for Agent Categorization

```
1: function DIVIDE(agents, metrics, depth)
       if CountLeaves() \ge 9 or |agents| < 4 then
 2:
           return LeafNode(agents)
 3:
       end if
 4:
       m^* \leftarrow \text{SelectMetricWithMaxRange(metrics)}
 5:
       r \leftarrow \text{ComputeRange}(m^*, \text{ agents})
 6:
       if r < threshold then
 7:
           if depth < 2 and action_tendency \in metrics then
 8:
               return ForceBinarySplit(agents, action tendency)
 9:
10:
           end if
           return LeafNode(agents)
11:
       end if
12:
       gaps \leftarrow FindNaturalGaps(m^*, agents)
13:
       subgroups \leftarrow SplitOnGaps(agents, gaps)
14:
       metrics' \leftarrow metrics \setminus \{m^*\}
15:
       for group in subgroups do
16:
           DIVIDE(group, metrics', depth + 1)
17:
18:
       end for
       return BranchNode(m^*, subgroups)
19:
20: end function
```

#### Contrast Effect

Contrast Effect is rooted in the well-established psychological phenomenon, where the perception of a stimulus is influenced by the context in which it is presented (Kenrick & Gutierres, 1980). We implement this effect using inspiration from a fundamental mechanism in cognitive architectures where an activated memory trace spreads activation to related traces and influences their accessibility (Anderson, 1983).

Building on this theoretical framework, we propose that memory activation spreads between agents in strategic interactions - when evaluating one agent, experiences with other agents influence memory activation levels. This effect amplifies memories that highlight behavioral contrasts. For example, when an agent with a particular behavioral tendency is evaluated with the existence of an agent with an opposing tendency, memories that emphasize the distinctive characteristics of the first agent become more strongly activated.

To formalize this process, we introduce an algorithm that computes the contrast effect using retrieved instances (Algorithm 4). activation (B) represents the initial activation level of each instance. For each retrieved instance i associated with a peer belonging to category  $c_i$ , we calculate a stereotype score  $s_i$  as:

$$s_i = \sum_{m \in M} w_m (1 - |v_{i,m} - p_{c_{i,m}}|)$$

where M is the set of behavioral metrics,  $v_{i,m}$  is instance i's value for metric m, and  $p_{c_i,m}$  is category  $c_i$ 's prototype value for metric m. This formulation ensures that instances are scored based on how closely they match their category's prototype. We then calculate the contrast between category  $c_i$  and each other category  $c_j$  using weighted cosine similarity between their prototypes:

$$C(p_{c_i}, p_{c_j}) = \frac{\sum_{m \in M} w_m \cdot p_{c_i, m} \cdot p_{c_j, m}}{\sqrt{\sum_{m \in M} w_m \cdot p_{c_i, m}^2} \cdot \sqrt{\sum_{m \in M} w_m \cdot p_{c_j, m}^2}}$$

where  $p_{c_i}$  and  $p_{c_j}$  are the prototypes of categories  $c_i$  and  $c_j$  respectively, and  $w_m$  is the discriminative power coefficient for metric m. The final activation value for instance i

combines its activation with contrast effects from all other categories:

Augmented-Activation
$$(i, c_i) = A_i + \sum_{c_j \neq c_i} cf \cdot C(p_{c_i}, p_{c_j}) \cdot s_i$$

where cf = 1/(n-1) is the contrast factor (with n being the number of categories) that determines the relative contribution of contrast against every other category. The activations of instances that strongly exemplify their category's characteristic patterns are amplified through this process.

Algorithm 4 Calculate Contrast-Augmented Activation for Retrieved Instances

**Require:** Retrieved instances R from category  $c_i$ , Set of other categories  $\{c_j\}_{j\neq i}$ , activation A (Eq. 1)

- 1: **for** each instance i in R **do**
- 2: Calculate stereotype score  $s_i = \sum_{m \in M} w_m (1 |v_{i,m} p_{c_{i,m}}|)$
- 3: Set contrast augmentation  $CA_i = 0$
- 4: **for** each other category  $c_j$  where  $j \neq i$  **do**
- 5: Calculate contrast strength  $C(p_{c_i}, p_{c_j})$  using weighted cosine similarity
- 6:  $CA_i = CA_i + cf \cdot C(p_{c_i}, p_{c_i}) \cdot s_i$
- 7: end for
- 8: Total activation:  $AA_i = A_i + CA_i$
- 9: end for

# **Emergence of Categorical Social Learning**

To demonstrate how our proposed mechanisms—category learning and contrast effects—influence an agent's development of prosociality regarding each peer  $(\alpha)$ , we instantiated our cognitive model in a repeated Prisoner's Dilemma setting. In this setup, a focal agent interacts one-on-one with multiple peers, independently across rounds. Importantly, peers don't interact with each other (only with the focal agent), as these peers serve as fixed-strategy agents while we focus on the focal agent's learning behavior.

In this specific context, our general behavioral metrics take domain-specific forms: the agent's general action tendency manifests as cooperation rate, responsiveness corresponds to reciprocal behavior, and recovery propensity reflects willingness to cooperate again after experiencing defection. We initialized our model with two fundamental categories for PD: a purely cooperative category and a purely defective one, which serve as cognitive anchors for the category learning process.

We compared learning trajectories under three conditions while simultaneously interacting with peers who showed systematically varied cooperation rates (0.1-1.0). This choice was motivated by two key considerations. First, since our agent can learn from experiences, it naturally develops tit-for-tat-like behavior, which would elicit predominantly mutual cooperation from classic iterated prisoner's dilemma strategies (as shown in Figure A1's TitForTat panel in the Appendix), masking the differential effects of our cognitive mechanisms. Second, cooperation rate is the most discriminative characteristic for strategy categorization. By systematically varying this key metric through fixed-rate peers, we obtain more discriminable behavioral patterns, allowing clearer observation of category formation and contrast effects.

In the baseline condition (Figure 4a), each line represents the  $\alpha$  value trajectory for a single peer, showing how the agent's prosociality toward each individual develops over time. We observe that  $\alpha$  values gradually differentiate based on the cooperation rates of the individual peers, with higher cooperators generally receiving higher weights. This demonstrates a basic form of reciprocity, where the agent learns to adjust its prosociality based on each peer's individual behavior.

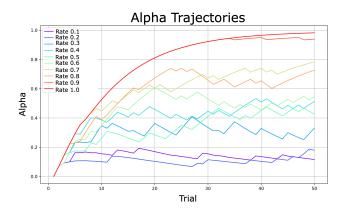
When category learning is enabled (Figure 4b), the agent no longer maintains separate  $\alpha$  values for each individual peer. Instead, it clusters peers into behavioral categories and develops shared  $\alpha$  values for all peers assigned to the same category. Each line in this panel represents the  $\alpha$  trajectory for a distinct category, not an individual peer. Two categories were initialized with cooperative and defective prototypes, serving as

cognitive anchors. Peers whose behavior closely matches these prototypes (cooperation rates near 1.0 or 0.1) are assigned to these predefined categories. Additionally, "learned categories" emerge when the agent encounters peers whose behavior patterns fall between these extremes but demonstrate internal consistency (e.g., peers with cooperation rates of 0.4-0.6). This categorical learning leads to more stable trajectories and a clearer separation between behavioral types, demonstrating how categorization can help manage multiple relationships more efficiently.

With both category learning and contrast effects active (Figure 4c), these category-based  $\alpha$  trajectories become even more distinct. Each line again represents a category, not an individual peer. The agent develops consistently higher  $\alpha$  values for the cooperative category and lower values for the defective category, with the learned categories maintaining stable intermediate values. For example, learned category 1 receives intermediate but relatively high  $\alpha$  values because it includes peers whose behavior is generally cooperative but not consistent enough to qualify for the cooperative prototype. This enhanced separation demonstrates how contrast effects strengthen categorical perception: the presence of clearly cooperative or defective peers makes behavioral differences more salient, leading to more pronounced differentiation between categories while maintaining stable within-category treatment.

Figure 5 shows the focal agent's cooperation rates when interacting with peers of varying cooperation probabilities under the three conditions. The cooperation rates positively correspond to the prosociality ( $\alpha$ ) values shown in Figure 4. In the baseline condition, the agent's cooperation gradually adapts to match the cooperation level of each peer, demonstrating basic reciprocity. When category learning is enabled, we observe more pronounced differentiation in cooperation rates across peers, with higher rates for cooperative peers and lower rates for defective ones, reflecting the agent's ability to recognize and respond to distinct behavioral patterns. The addition of contrast effects further amplifies this differentiation, showing how the agent's cooperation becomes more

selective—maintaining high cooperation with cooperative peers while reducing cooperation more dramatically with defective ones. This pattern demonstrates how our cognitive mechanisms enable more sophisticated and adaptive social behavior beyond simple reciprocity.



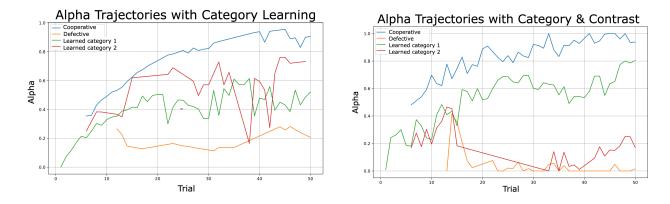
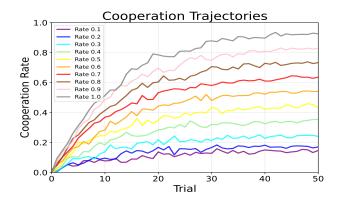
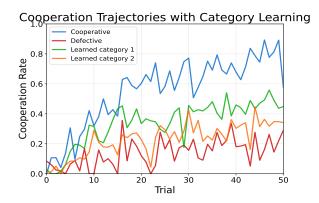


Figure 4

Evolution of alpha values under different cognitive mechanisms across 50 trials, all tracking interactions with peers of varying cooperation rates (0.1–1.0). (Top) Baseline condition showing individual alpha trajectories for each peer. (Bottom-left) With category learning enabled, alpha trajectories reflect learning of behavioral categories rather than individuals. (Bottom-right) The combined effect of category learning and contrast mechanisms demonstrates enhanced separation between learned behavioral categories. Higher alpha values indicate greater weight given to a peer's outcomes in decision-making.





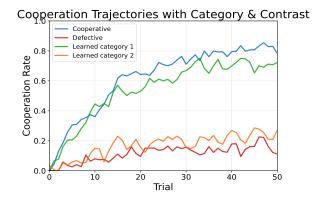


Figure 5

Cooperation rates of the focal agent when interacting with peers of varying cooperation probabilities (0.1–1.0) across 50 trials under different cognitive mechanisms. (Top) Baseline condition showing adaptation toward matching peer cooperation levels. (Bottom-left) With category learning enabled, the focal agent demonstrates more differentiated cooperation patterns based on peer categories. (Bottom-right) With both category learning and contrast effects, the focal agent shows amplified selective cooperation—maintaining high rates with cooperative peers while reducing cooperation with defectors.

### **Empirical Validation with Human Data**

In this section, we will validate our IBL-Group model on a dataset collected from a controlled online experiment involving 150 participants (50 three-person groups) (Du et al., 2025). In this experiment, players play a repeated three-player game in the context of cybersecurity information sharing. In each round of the game, the players can choose to share information with other players or not. The players need to play strategically, and their outcomes are interdependent based on their joint actions. More specifically, the game consists of 50 rounds. Each player began with a one-time 1,000-point endowment and a 53% chance of being breached in the first round. In each round, participants first learned whether they had been attacked (incurring a 30-point penalty) and then chose whether to share that attack status with one or both groupmates; sharing cost the sender 15 points and granted 35 points to each recipient. Following each round, players saw updated point totals, their groupmates' sharing decisions, and attack outcomes, enabling strategic adaptation. Before starting, participants provided informed consent, reviewed detailed instructions, and passed a comprehension quiz to ensure full understanding of the task mechanics

To evaluate IBL-group's fit to human behavior, we ran 50 simulation replicates—one per experimental triad—each with three agents (150 agents total) interacting over T=50 rounds under the same payoffs (breach penalty = 30, share cost = 15, share benefit = 35). Since IBL-group model extends the individual IBL-PD model to a group setting, we first set individual agent's IBL-PD model parameters values following existing work (Gonzalez et al., 2015): memory decay d=0.5, mismatch penalty  $\mu=1$ , feature weights  $\omega_m=1$ , retrieval noise  $\sigma=0.25$ , initial prosociality  $\alpha(0)=0.5$ , and learning rate  $\eta=0.1$ . For the new parameters introduced by IBL-group, we choose the parameters based on insights from closely related cognitive model literature. Specifically, we set the sliding window w=5 in category learning, as it reflects the working memory capacity, the prototype matching threshold  $\theta_{\text{conf}}=0.8$ , minimum cluster size m=2, and stability threshold  $\theta_{\text{stab}}=5$ .

Figure 6A plots the overall cooperation rate—defined as the proportion of rounds in which a participant shared with at least one groupmate—across all 50 rounds. The IBL-group model closely matches the human data, yielding a mean squared deviation of MSD = 0.02 and a Pearson correlation coefficient of r = 0.86 (p < 0.001). Both curves experience their steepest decline over approximately the first 15 rounds, after which the rate of change diminishes, indicating an initial phase of strategic learning followed by more stable behavior.

The model also closely matched the distribution of the sharing decisions, as shown in Figure 6. It reproduced the shifts from initially sharing with both peers to increasingly sharing with just one partner over time. The MSD values were 0.01 and 0.02 for the proportions of sharing with both and one, respectively. The correlations were also strong for sharing with both (r = 0.90, p < 0.001), although weaker for sharing with one (r = 0.35, p < 0.05), potentially due to more variability in human strategies.

One divergence appears in the "shared with one" category: participants did so in 27% of rounds, while the IBL-group model predicted 25%, an absolute gap of 2%. The "shared with both" rates were nearly identical (human: 29%, model: 30%). This minor underprediction likely reflects stochastic variability in human sharing decisions around the model's central tendency. Despite the small divergence, the model successfully reproduced two principal empirical patterns: (1) the declining overall cooperation rate across 50 rounds (Figure 6A) and (2) the shift from predominantly "shared with both" to increasing "shared with one" behavior over time (Figure 6B). This quantitative alignment confirms that the model's combination of dynamic instance-based retrieval and category-learning mechanisms captures the core temporal and distributional features of human sharing behavior.

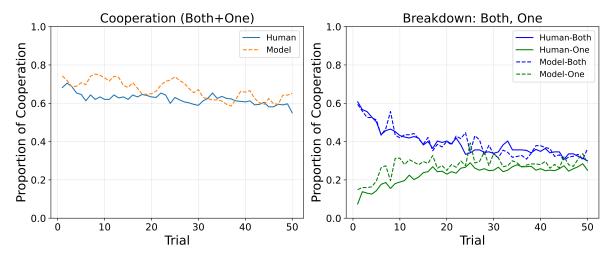


Figure 6

Model predictions versus human behavior in the information-sharing game. (A) Overall cooperation rates over 50 rounds. (B) Proportions of participants sharing information with both, one, or neither of their peers over time. Shaded regions represent 95% confidence intervals.

# Psychological interpretation of $\alpha$ with model-tracing

To understand how the prosociality parameter  $\alpha$  develops in human participants, we use model tracing (Anderson & Schooler, 1991). This approach applies our IBL-Group model to each participant's actual sequence of decisions and outcomes. For each participant, we create an agent that follows our cognitive model with the default parameters. As the participant progresses through the experiment, we populate the agent's memory with instances directly from the human's behavioral data. Each instance contains the participant's sharing decision, their attack outcome, and the corresponding peer responses observed in that trial. The agent processes these instances and updates its internal  $\alpha$  values according to our learning mechanism, without generating new decisions. This trial-by-trial memory population produces a trajectory of  $\alpha$  values that reflects how each participant likely valued their peers' outcomes throughout the interaction sequence. This procedure follows the decision-from-experience model-tracing method (Anderson & Schooler, 1991) and its adaptation to cybersecurity research (Cranford et al., 2019).

### **Dyadic Analysis**

Our model proposes that players develop their prosociality regarding each other's outcomes ( $\alpha$ ) through repeated interactions. Although previous analysis has shown how a focal agent develops different  $\alpha$  values for each of the peers in a triad with distinct behavioral patterns, here we examine the bidirectional nature of the dynamics of  $\alpha$  between pairs of players. This analysis reveals how mutually sharing decisions shape the evolution of prosociality in dyadic relationships. We categorize dyadic relationships based on players' sharing rates over the 50-round interaction period. High sharing is defined as sharing information in more than 70% of opportunities, while low sharing is defined as a lower share rate than 30%. We then analyze the correlation between the paired  $\alpha$  values and their temporal development patterns.

Figure 7 illustrates these relationship patterns by presenting the average trajectories of the dynamic prosociality parameter  $(\alpha)$  within dyadic pairs, categorized based on their sharing rates across 50 rounds: mutual high-sharing (both players shared more than 70% of the time), mutual low-sharing (both shared less than 30%), and asymmetric relationships (one shared frequently, while the other did not). Within each dyad, each agent maintains a separate  $\alpha$  value representing their prosocial regard toward their opponent. To summarize these patterns clearly, we first identified within each dyad the agent who ended with the higher final  $\alpha$  (greater prosocial regard toward their opponent) and the agent who ended with the lower final  $\alpha$ . We then averaged these "higher" and "lower"  $\alpha$  trajectories separately across all dyads within each relationship type. In mutual high-sharing pairs, both players exhibit closely synchronized trajectories, rapidly converging to similarly high  $\alpha$  values, indicating mutual responsiveness and reciprocal prosociality. In mutual low-sharing pairs, both players'  $\alpha$  values simultaneously decline and stabilize at low levels, reflecting mutual defection or lack of reciprocity. In contrast, asymmetric pairs show distinctly diverging trajectories, with the more cooperative partner maintaining high  $\alpha$ values while the less cooperative partner's values decline, capturing imbalanced

responsiveness and limited reciprocity. This clear differentiation highlights the model's capability to capture the interdependent evolution of prosocial behavior based explicitly on observed patterns of cooperation.

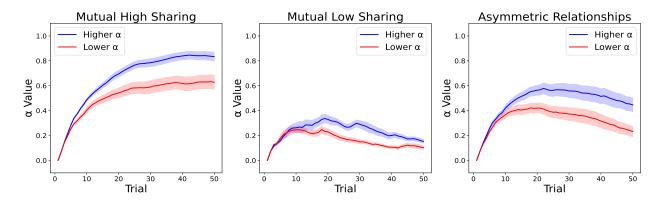


Figure 7

Trajectories of  $\alpha$  for different dyadic relationship types across 50 trials (shaded areas represent standard errors). Left: Mutual high-sharing pairs (>70% sharing rate). Middle: Mutual low-sharing pairs (<30% sharing rate). Right: Asymmetric relationships. For each relationship type, trajectories are shown for each pair's relatively higher and lower alpha values.

Analysis of  $\alpha$  trajectories reveals systematic patterns in different types of dyadic relationships. In pairs of mutual high sharing (N=17), the values of  $\alpha$  showed a strong positive correlation  $(r=.91,\,p<.001)$  and converged to similar high levels (mean = .87, SD = .04) at trial 30. The mean absolute difference between the values of paired  $\alpha$  in these relationships remained small (mean = .06, SD = .04), indicating synchronized prosociality. Mutual pairs of low-sharing (N=9) demonstrated a coordinated decline in  $\alpha$  values, stabilizing at lower levels (mean = .19, SD = .08) with moderate correlation between paired values  $(r=.61,\,p<.001)$ . The temporal pattern shows a rapid initial decline followed by stabilization around trial 25. The asymmetric relationships (N=22) produced the most divergent  $\alpha$  trajectories. The high-sharing players maintained significantly higher values of  $\alpha$  (mean = .82, SD = .08) compared to their low-sharing peers (mean = .28, SD

= .11; t(21) = 15.33, p < .001). The correlation between the values of the pairs'  $\alpha$  in these relationships was weak (r = .29, p = .196), reflecting the disconnect in prosociality.

### Triadic Analysis

Beyond dyadic relationships, we examine how prosociality develops at the group level, specifically testing whether strong cooperative relationships between some members influence cooperation with others. This analysis addresses a key question: does the emergence of trust between two players create spillover effects that facilitate cooperation throughout the entire group? Understanding these indirect effects is crucial for predicting when and how cooperative clusters emerge in larger social networks. We categorize triads based on how quickly the dyads of players achieve mutual high  $\alpha$  values. Specifically, we define a pair as having mutual high"  $\alpha$  if both players'  $\alpha$  values exceed 0.7. Triads are thus classified as early-forming" if at least one pair achieves mutual high  $\alpha$  within the first 10 trials, and as "gradual" if no pairs reach this level within that initial period.

The analysis reveals two key patterns. First, when a pair of players establish mutually high  $\alpha$  values, this successful relationship shapes how the third player develops their  $\alpha$  values for both peers. The third player tends to develop relatively high and stable  $\alpha$  values (around 0.6) for both members of the strong pair, suggesting that  $\alpha$  development is influenced not just by direct interactions but also by observed relationship strength between peers. Second, triads with early-forming strong relationships show distinctly different evolution patterns from gradual triads. In early-forming triads, the initial strong relationship appears to create positive conditions for the remaining relationships, leading to higher and more stable  $\alpha$  values throughout the triad. In contrast, gradual triads show lower  $\alpha$  values (0.2-0.4) and more differentiation between pairs, indicating that the absence of an early strong relationship leads to more tentative cooperation throughout the group. Figure 8 illustrates these patterns by tracking the development of the third player  $\alpha$  towards each partner in both types of triads.

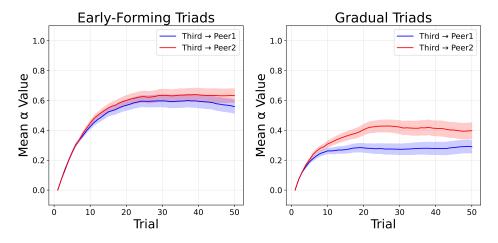


Figure 8

Development of  $\alpha$  values from the third player's perspective in different triad formations over 50 trials (shaded areas represent standard errors). Left: In early-forming triads (where one pair establishes mutual cooperation  $\alpha > 0.6$  within 10 trials), the third player develops relatively high and similar  $\alpha$  values for both peers. Right: In gradual triads (no early high mutual  $\alpha$ ), the third player's  $\alpha$  values remain lower and more differentiated between peers.

#### Discussion

Our findings offer a cognitively grounded account of how individuals make interdependent decisions in group contexts, revealing how dynamic prosociality, category learning, and contrast effects represent adaptive cooperation under cognitive constraints. The proposed IBL-Group model advances the literature by expanding the IBL-PD beyond dyadic interactions, addressing a gap identified in previous work (Gonzalez et al., 2015) and responding to long-standing concerns regarding the scalability and cognitive plausibility of agent models in multiparty strategic environments (Macrae & Bodenhausen, 2000; Stiller & Dunbar, 2007).

The proposed IBL-Group model successfully reproduces core patterns of human social behavior observed in multiagent experimental settings. Specifically, we show that dynamic prosociality enables agents to calibrate their cooperative tendencies based on observed reciprocation, producing distinct  $\alpha$  trajectories toward cooperative peers versus defective peers.

Category learning allows efficient representation of multiple peers by clustering them into prototypes, preserving essential behavioral distinctions without exceeding cognitive limits. The integration of category learning with instance-based memory provides a cognitively plausible explanation for how humans efficiently process multiple relationships despite memory limitations. The model's ability to match human cooperation rates (r=0.86) without parameter tuning underscores the model's explanatory power. The results suggests that categorical processing is not merely a coping mechanism for cognitive constraints but a fundamental aspect of social learning. This aligns with recent research showing that cognitive limitations can paradoxically foster robust social strategies (Stevens et al., 2018).

Contrast effects sharpen those distinctions, enhancing the discrimitation between behavioral categories and improving responsiveness to diverse partner types. This mechanism offers new insights into how the social context shapes cooperative decisions. Model-tracing further reveals that prosociality develops not only as a function of direct interaction, but also in relation to peer dynamics and social structure within triads. The early emergence of strong dyadic cooperation fosters greater respect from third-party observers, highlighting the social signaling role of reciprocal ties. Our tracing analysis reveals that agents develop significantly different  $\alpha$  values for categorically distinct peers (mean difference = 0.75, p < 0.001), consistent with empirical findings that humans evaluate cooperation relative to broader social experiences (Wu et al., 2020a). This contextual processing may explain why cooperative clusters emerge in network simulations, as local reference points reinforce categorical boundaries between cooperative and non-cooperative regions.

These results underscore the utility of cognitive mechanisms not only as approximations of human limitations but as strategic assets. Our findings challenge the dominant assumption in the computational modeling literature that rich agent modeling requires exhaustive memory or complete partner tracking (e.g., (Leibo et al., 2017; Lowe

et al., 2017). Instead, we demonstrate that cognitively efficient processes—specifically categorical abstraction and contrastive amplification—can yield robust, generalizable social strategies.

Despite its strengths, our model has some limitations. First, the model assumes that agents have perfect and equal access to others' actions and outcomes. This idealization overlooks the informational asymmetries and noise common in real-world environments (e.g., organizations, cybersecurity, or social network settings). Second, we currently assume identical cognitive parameters across agents. This limits the exploration of individual differences in prosociality, working memory, or learning rate, all of which likely influence the dynamics of group-level cooperation. Finally, the model operates in static triads. It remains an open question how these mechanisms scale to fluid social environments, with changing memberships, reputations, and social ties.

Several avenues for future research emerge from this work. Introducing uncertainty in feedback and differing observational abilities could test the robustness of categorical and contrast-based learning mechanisms under more realistic conditions. Incorporating heterogeneous agent traits would also allow for a detailed exploration of group composition effects. Future work should extend this model to larger networks and dynamic group memberships. Specifically, understanding how categorical processing adapts to fluid group boundaries and how contrast effects operate across multiple reference groups will be critical for applications in cybersecurity information sharing and organizational collaboration. Finally, based on the model results, future work can design and test behavioral interventions to promote cooperation. by leveraging human categorization and social evaluation tendencies.

### Acknowledgments

This research was sponsored by the Army Research Office and accomplished under Australia-US MURI Grant Number W911NF-20-S-000.

#### References

- Anderson, J. R. (1983). A spreading activation theory of memory. *Journal of Verbal Learning and Verbal Behavior*, 22(3), 261–295.
- Anderson, J. R., & Schooler, L. J. (1991). Reflections of the environment in memory.

  Psychological Science, 2(6), 396–408.
- Andreoni, J., & Miller, J. H. (1993). Rational cooperation in the finitely repeated prisoner's dilemma: Experimental evidence. *The economic journal*, 103(418), 570–585.
- Axelrod, R. (1984). The evolution of cooperation. *Basic Books*.
- Axelrod, R., & Hamilton, W. D. (1981). The evolution of cooperation. *science*, 211(4489), 1390–1396.
- Baker, C., Saxe, R., & Tenenbaum, J. B. (2011). Bayesian theory of mind: Modeling joint belief-desire attribution. *Proceedings of the Annual Meeting of the Cognitive Science Society*, 33, 2469–2474.
- Balliet, D., & Van Lange, P. A. M. (2013). Trust, conflict, and cooperation: A meta-analysis. *Psychological Bulletin*, 139(5), 1090–1112. https://doi.org/10.1037/a0030939
- Balliet, D. P., Tybur, J. M., & Van Lange, P. A. M. (2017). Functional interdependence theory: An evolutionary account of social situations. *Personality and Social Psychology Review*, 21(4), 361–388. https://doi.org/10.1177/1088868316657965
- Balliet, D., & Lindström, B. (2023). Inferences about interdependence shape cooperation.

  Trends in Cognitive Sciences, 27(4), 289–301.
- Carmel, D., & Markovitch, S. (1995). Opponent modeling in multi-agent systems. *IJCAI*, 95, 40–45.
- Centola, D. (2011). An experimental study of homophily in the adoption of health behavior. *Science*, 334 (6060), 1269–1272.
- Chierchia, G., et al. (2017). Integrating social and nonsocial cognitive control. *Cognition*, 163, 42–60.

- CISA. (2023). Information sharing | cybersecurity and infrastructure security agency [Accessed: 2025-04-12].
  - https://www.cisa.gov/topics/cyber-threats-and-advisories/information-sharing
- Cowan, N. (2001). The magical number 4 in short-term memory: A reconsideration of mental storage capacity. *Behavioral and Brain Sciences*, 24(1), 87–114. https://doi.org/10.1017/S0140525X01003922
- Crandall, J. W., et al. (2018). Cooperating with machines. Nature Communications, 9(1), 1-12.
- Cranford, E. A., Lebiere, C., Rajivan, P., Aggarwal, P., & Gonzalez, C. (2019). Modeling cognitive dynamics in end-user response to phishing emails. *Proceedings of the 17th ICCM*.
- Du, Y., Singh, K., Aggarwal, P., Fang, F., & Gonzalez, C. (2025). Emergent cooperative decision-making in triadic prisoner's dilemma: Effects of incentives and information [Accessed: 2025-02-13].
- Dunbar, R. I. M. (1998). The social brain hypothesis. *Evolutionary Anthropology: Issues, News, and Reviews*, 6(5), 178–190.
- Dziura, S. L., Hosangadi, A., Shariq, D., Merchant, J. S., & Redcay, E. (2023). Partner similarity and social cognitive traits predict social interaction success among strangers.

  Social Cognitive and Affective Neuroscience, 18(1), nsad045.
- Earnest, M. J. (2013). Extortion and evolution in the iterated prisoner's dilemma.
- Erev, I., & Roth, A. E. (2006). Learning in extensive-form games: Experimental data and simple dynamic models in the intermediate term. *Games and Economic Behavior*, 55(1), 1–29. https://doi.org/10.1016/j.geb.2005.03.001
- Fehr, E., & Schurtenberger, I. (2019). Cooperative phenotypes predict reciprocal behavior. Science, 365(6457).
- Fiske, S. T. T., & Taylor, S. E. (2020). Social cognition: From brains to culture.

- Gallotti, R., & Grujić, J. (2018). A drift-diffusion model of the cooperative decision-making process. *Scientific Reports*, 8(1), 1–10. https://doi.org/10.1038/s41598-018-21924-2
- Gigerenzer, G., & Gaissmaier, W. (2011). Heuristic decision making. *Annual Review of Psychology*, 62, 451–482. https://doi.org/10.1146/annurev-psych-120709-145346
- Gilbert, D. T., & Malone, P. S. (1995). The correspondence bias. *Psychological bulletin*, 117(1), 21.
- Gmytrasiewicz, P. J., & Doshi, P. (2005). A framework for reasoning about others: The theory of mind approach. *Journal of Artificial Intelligence Research*, 24, 1–35.
- Gonzalez, C. (2024). Building human-like artificial agents: A general cognitive algorithm for emulating human decision-making in dynamic environments. *Perspectives on Psychological Science*, 19(5), 860–873.
- Gonzalez, C., Ben-Asher, N., Martin, J. M., & Dutt, V. (2015). A cognitive model of dynamic cooperation with varied interdependency information. *Cognitive science*, 39(3), 457–495.
- Gonzalez, C., Lerch, J. F., & Lebiere, C. (2003). Instance-based learning in dynamic decision making. *Cognitive Science*, 27(4), 591–635.
- Hámornik, B. P., & Krasznay, C. (2017). A team-level perspective of human factors in cyber security: Security operations centers. *International Conference on Applied Human Factors and Ergonomics*, 224–236.
- Hansen, J., Alves, H., & Trope, Y. (2016). Psychological distance reduces literal imitation: Evidence from an imitation-learning paradigm. *Journal of Experimental Psychology Human Perception & Performance*. https://doi.org/10.1037/xhp0000150
- Harper, M., et al. (2017). Reinforcement learning produces dominant strategies for the iterated prisoner's dilemma. *PloS One*, 12(12), e0188046.
- Hernandez-Leal, P., et al. (2019). A survey of learning in multiagent environments: Dealing with non-stationarity. arXiv preprint arXiv:1709.02779.

- Hertwig, R., & Erev, I. (2009). Decisions from experience and statistical probabilities: Why they trigger different choices than a priori probabilities. *Journal of Behavioral Decision Making*, 22(1), 1–35.
- Huang, L. M., & Sherman, J. W. (2018). Attentional processes in social perception. https://doi.org/10.1016/bs.aesp.2018.03.002
- Kalkstein, D. A., Kleiman, T., Wakslak, C., Liberman, N., & Trope, Y. (2016). Social learning across psychological distance. *Journal of Personality and Social Psychology*. https://doi.org/10.1037/pspa0000042
- Kelley, H. H., & Stahelski, A. J. (1970). Social interaction basis of cooperators' and competitors' beliefs about others. *Journal of personality and social psychology*, 16(1), 66.
- Kenrick, D. T., & Gutierres, S. E. (1980). Contrast effects and judgments of physical attractiveness: When beauty becomes a social problem. *Journal of Personality and Social Psychology*, 38(1), 131.
- Kirchkamp, O., et al. (2016). Behavioral spillovers and cooperation. *Journal of Economic Behavior & Organization*, 130, 160–175.
- Kleiman-Weiner, M., et al. (2016a). Coordinate to cooperate or compete: Abstract goals and joint intentions in social interaction. *Topics in Cognitive Science*, 8(2), 424–435.
- Kleiman-Weiner, M., et al. (2016b). Learning to cooperate: The evolution of social rewards in repeated interactions. *Topics in Cognitive Science*, 8(4), 860–877.
- Knight, V., et al. (2015). Axelrod library tutorials [Accessed: 2025-02-13]. https://axelrod.readthedocs.io/en/stable/tutorials/index.html
- Lanctot, M., et al. (2017). A unified game-theoretic approach to multiagent reinforcement learning. Advances in Neural Information Processing Systems, 4193–4206.
- Leibo, J. Z., et al. (2017). Multi-agent reinforcement learning in sequential social dilemmas.

  Proceedings of AAMAS.

- Li, X., Richter, A., & Lehtonen, J. (2023). Modeling social evolution: A review of evolutionary game theory. *Journal of Evolutionary Biology*, 36(5), 799–824. https://doi.org/10.1111/jeb.14159
- Lim, S., et al. (2016). Opponent modeling in deep reinforcement learning. *International Conference on Machine Learning*.
- Lowe, R., et al. (2017). Multi-agent actor-critic for mixed cooperative-competitive environments. Advances in Neural Information Processing Systems, 6379–6390.
- Macrae, C. N., & Bodenhausen, G. V. (2000). Social cognition: Thinking categorically about others. *Annual review of psychology*, 51(1), 93–120.
- Macrae, C. N., & Bodenhausen, G. V. (2011). Categories in context: How features of perceived similarity and category membership guide social evaluation. *Social cognition*, 29(5), 547–562.
- Martin, J. M., Gonzalez, C., Juvina, I., & Lebiere, C. (2014). A description–experience gap in social interactions: Information about interdependence and its effects on cooperation.

  Journal of Behavioral Decision Making, 27(4), 349–362.
- McPherson, M., Smith-Lovin, L., & Cook, J. M. (2001). Birds of a feather: Homophily in social networks. *Annual review of sociology*, 27(1), 415–444.
- Mermoud, A., Keupp, M. M., Huguenin, K., Palmié, M., & Percia David, D. (2019). To share or not to share: A behavioral perspective on human participation in security information sharing. *Journal of Cybersecurity*, 5(1), tyz006.
- Miller, G. A. (1956). The magical number seven, plus or minus two: Some limits on our capacity for processing information. *Psychological Review*, 63(2), 81–97.
- Moreira, J. A., et al. (2013). Social dilemmas and cooperation. *Physics of Life Reviews*, 10(4), 208–243.
- Nowak, M. A. (2006). Five rules for the evolution of cooperation. *Science*, 314 (5805), 1560–1563.

- Peysakhovich, A., & Lerer, A. (2018). Consequentialist learning in repeated interactions.

  AAMAS.
- Rand, D. G., et al. (2020). Cluster analysis reveals a binary effect of storage on drosophila courtship conditioning. *Nature Communications*, 11(1), 1–12.
- Rosch, E. (1978). Principles of categorization. Cognition and Categorization, 27–48.
- Sherman, J. W., Kruschke, J. K., Sherman, S. J., Percy, E. J., Petrocelli, J. V., & Conrey, F. R. (2009). Attentional processes in stereotype formation: A common model for category accentuation and illusory correlation. *Journal of Personality and Social Psychology*. https://doi.org/10.1037/a0013778
- Shum, M., Kleiman-Weiner, M., Littman, M. L., & Tenenbaum, J. B. (2019). Theory of minds: Understanding behavior in groups through inverse planning. arXiv preprint arXiv:1901.06085.
- Simon, H. A. (1974). How big is a chunk? by combining data from several experiments, a basic human memory unit can be identified and measured. *Science*, 183(4124), 482–488.
- Simon, H. A. (1990). Bounded rationality. MIT Press.
- Stevens, J. R., et al. (2018). Cognitive constraints in strategic decision making. Nature Human Behaviour, 2(3), 213–221.
- Stewart, A. J., & Plotkin, J. B. (2013). From extortion to generosity, evolution in the iterated prisoner's dilemma. *PNAS*, 110(38).
- Stiller, J., & Dunbar, R. I. (2007). Perspective-taking and memory capacity predict social network size. *Social Networks*, 29(1), 93–104.
- Tamarit, I., Cuesta, J. A., Dunbar, R. I., & Sánchez, A. (2018). Cognitive resource allocation determines the organization of personal networks. *Proceedings of the National Academy of Sciences*, 115(33), 8316–8321.
- Todd, P. M., et al. (2012). Ecological rationality: Intelligence in the world. Oxford University Press.

- Van Lange, P. A. M., Balliet, D. P., Parks, C. D., & Van Vugt, M. (2014). Social dilemmas:

  The psychology of human cooperation. Oxford University Press.
- Van Lange, P. A., Joireman, J., Parks, C. D., & Van Dijk, E. (2013). The psychology of social dilemmas: A review. *Organizational Behavior and Human Decision Processes*, 120(2), 125–141.
- Wang, Z., et al. (2018). Evolving strategies for the iterated prisoner's dilemma. *Scientific Reports*, 8(1), 1–8.
- Wu, J., et al. (2020a). Social networks and cooperation. *Nature Human Behaviour*, 4(1), 96–104.
- Wu, J., et al. (2020b). Too much or too little? a meta-analysis of the social comparison effects on cooperation. *Psychological Bulletin*, 146(7), 651–679.
- Yoshida, W., Dolan, R. J., & Friston, K. J. (2008). Game theory of mind. *PLoS Computational Biology*, 4(12), e1000254. https://doi.org/10.1371/journal.pcbi.1000254
- Young, J., et al. (2019). Attention and strategic decision making. *Cognitive Science*, 43(4), e12721.
- Zhang, D., Lin, Y., Jing, Y., Feng, C., & Gu, R. (2019). The dynamics of belief updating in human cooperation: Findings from inter-brain erp hyperscanning. *NeuroImage*, 198, 1–12.

### Appendix

## Behavioral Clustering of Classical IPD Strategies

To demonstrate that peer strategies can be effectively categorized through their behavioral patterns, we use three different focal agents, TitForTat, Cooperator, and Defector, to play against a set of peers using classic strategies in iterated prisoner's dilemma (Knight et al., 2015). As shown in Figure A1, the clusters depend on the interaction history and thus are different for each focal agent. When interacting with TitForTat (left panel), the strategies showed a clear separation between highly cooperative behaviors (cooperation rate > 0.8, high reciprocity) and more defensive or reactive strategies (moderate cooperation rates with varying reciprocity). The Cooperator focal agent (middle panel) elicited the most pronounced cooperative behaviors, with many strategies showing both high cooperation rates and high reciprocity. In contrast, the Defector focal agent (right panel) demonstrated how strategies adapt to persistent defection, with most strategies showing lower cooperation rates and a negative correlation between cooperation rate and reciprocity. Some strategies, particularly cyclic (Cyc) and random (Ran) strategies, showed distinctive behavioral patterns across the three reference agents, as their behavior is not dependent on the peer.

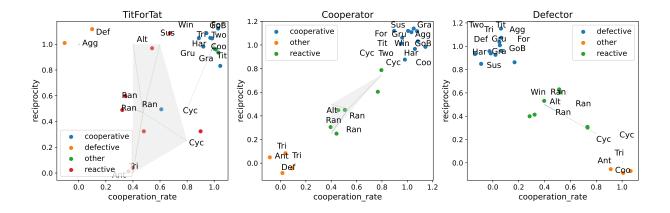
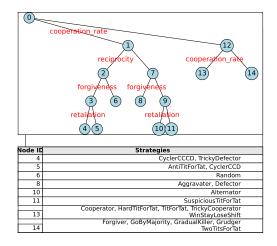
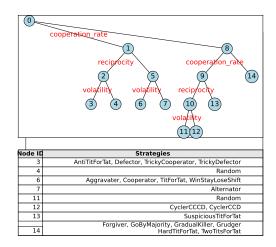


Figure A1

Behavioral clustering of IPD strategies based on their cooperation rate and reciprocity after 50 rounds of interaction (using a sliding window of 5 rounds for behavioral metrics). Each panel shows the same set of peer strategies interacting simultaneously with a different focal agent (TitForTat, Cooperator, or Defector). Strategies are colored by their classified behavioral type. The varying distribution of strategies across panels demonstrates how peer behavior depends on the focal agent's strategy.





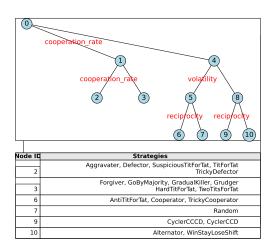


Figure A2

Hierarchical clustering of peer strategies when interacting with different focal agents:

TitForTat (Top-left), Cooperator (Top-right), and Defector (Bottom). The decision trees show how strategies are categorized based on behavioral metrics, with red labels indicating the splitting criteria at each node. Each terminal node (blue) corresponds to a group of strategies with similar behavioral patterns, as detailed in the accompanying tables. The varying tree structures across focal agents demonstrate how peer classification depends on the focal agent's strategy, with more complex distinctions possible against TitForTat and simpler categorizations emerging against persistent defection.

Figure A2 shows the hierarchical clustering process. In all three cases, the cooperation rate serves as the primary splitting criterion at the root, reflecting its fundamental role in strategy differentiation. Subsequent levels incorporate reciprocity, with forgiveness and retaliation metrics enabling finer distinctions at lower levels. Against TitForTat, this creates nuanced groupings, for instance, distinguishing Anti-TitForTat and Random strategies from Cycler-CCD, while grouping Cooperator with TrickyCooperator based on their similar behavioral patterns. When facing a Cooperator, the tree maintains multiple reciprocity-based splits, effectively separating strategies like WinStayLoseShift from pure cooperators. With a Defector as the focal agent, the tree becomes more compact, requiring fewer splits to categorize strategies, as persistent defection tends to elicit more uniform responses. These trees demonstrate how behavioral metrics can be systematically applied to classify strategies in a context-dependent manner.